

Phase Resolution of Heterozygous Sites in Diploid Genomes is Important to Phylogenomic Analysis under the Multispecies Coalescent Model

JUN HUANG^{1,2†}, JEREMY BENNETT,^{1,3,†} TOMÁŠ FLOURI,¹, ADAM D. LEACHÉ⁴ AND ZIHENG YANG^{1,*}

¹*Department of Genetics, Evolution and Environment, University College London, Gower Street, London WC1E 6BT, UK*

²*Department of Mathematics, Beijing Jiaotong University, Beijing, 100044, P.R. China*

³*Department of Ecology and Evolutionary Biology, University of Connecticut, 75 N. Eagleville Road, Unit 3043, Storrs, CT 06269-3043, USA*

⁴*Department of Biology & Burke Museum of Natural History and Culture, University of Washington, Seattle, WA 98195-1800, USA*

**Ziheng Yang, Department of Genetics, Evolution and Environment, University College London, Gower Street, London WC1E 6BT, UK (orcid: 0000-0003-3351-7981).
Phone: +44 20 76794379 (z.yang@ucl.ac.uk)*

† *Those authors contributed equally to the project.*

ABSTRACT

1 Genome sequencing projects routinely generate haploid consensus sequences from diploid
2 genomes, which are effectively chimeric sequences with the phase at heterozygous sites resolved
3 at random. The impact of phasing errors on phylogenomic analyses under the multispecies
4 coalescent (MSC) model is largely unknown. Here we conduct a computer simulation to evaluate
5 the performance of four phase-resolution strategies (the true phase resolution, the diploid
6 analytical integration algorithm which averages over all phase resolutions, computational phase
7 resolution using the program PHASE, and random resolution) on estimation of the species tree
8 and evolutionary parameters in analysis of multi-locus genomic data under the MSC model. We
9 found that species tree estimation is robust to phasing errors when species divergences were
10 much older than average coalescent times but may be affected by phasing errors when the species
11 tree is shallow. Estimation of parameters under the MSC model with and without introgression is
12 affected by phasing errors. In particular, random phase resolution causes serious overestimation
13 of population sizes for modern species and biased estimation of cross-species introgression
14 probability. In general the impact of phasing errors is greater when the mutation rate is higher, the
15 data include more samples per species, and the species tree is shallower with recent divergences.
16 Use of phased sequences inferred by the PHASE program produced small biases in parameter
17 estimates. We analyze two real datasets, one of East Asian brown frogs and another of Rocky
18 Mountains chipmunks, to demonstrate that heterozygote phase-resolution strategies have similar
19 impacts on practical data analyses. We suggest that genome sequencing projects should produce
20 unphased diploid genotype sequences if fully phased data are too challenging to generate, and
21 avoid haploid consensus sequences, which have heterozygous sites phased at random. In case the
22 analytical integration algorithm is computationally unfeasible, computational phasing prior to
23 population genomic analyses is an acceptable alternative.

24 *Key words:* BPP, introgression, multispecies coalescent, phase, species tree

25 1. INTRODUCTION

26 Next-generation sequencing technologies have revolutionized population genetics and
27 phylogenetics by making it affordable to sequence whole genomes or large portions of the
28 genome, even for non-model organisms. Many phylogenomic studies use the approach of
29 reduced representation library to maximize their DNA sequencing efforts on a small subset of the
30 genome. These strategies can generate thousands of genomic segments (called loci in this paper
31 irrespective of whether they are protein-coding) with high coverage, and target sequences can be
32 assembled with confidence. Examples include restriction site-associated DNA sequencing
33 (RADseq), which is used frequently to identify single nucleotide polymorphisms (SNPs) for
34 population genetic and phylogeographic studies (Andrews *et al.*, 2016; Leaché and Oaks, 2017),
35 although it has also been applied to address phylogenetic questions at deeper timescales (Eaton
36 *et al.*, 2017). A more common approach for phylogenomic studies is targeted sequence capture,
37 generating so-called reduced-representation datasets, with typically longer sequences for
38 distantly related species than with RADseq data. Examples include exome sequencing,
39 ultraconserved elements (UCEs, Faircloth *et al.*, 2012), anchored hybrid enrichment (AHE,
40 Lemmon *et al.*, 2012), conserved nonexonic elements (CNEEs, Edwards *et al.*, 2017), or rapidly
41 evolving long exon capture (RELEC, Karin *et al.*, 2020).

42 Typical sequencing technologies produce short fragments of sequenced DNA called
43 ‘reads’ that are either *de novo* assembled or mapped to a pre-existing reference genome. This
44 leads to chromosomal positions being sequenced a variable number of times across the genome
45 (usually referred to as the sequencing depth). A common practice in genome sequencing projects
46 has been to produce the so-called “haploid consensus sequence” for a diploid individual, which
47 uses the most common nucleotide at any heterozygous site to produce one genomic sequence.
48 Assemblers like Velvet (Zerbino and Birney, 2008), ABySS (Simpson *et al.*, 2009), and Trinity
49 (Grabherr *et al.*, 2011), pick up only one of the two nucleotide bases at any heterozygous site and
50 essentially resolve the phase of heterozygous sites at random, producing chimeric sequence that
51 may not exist in nature. Suppose a diploid individual is heterozygous at two sites in a genomic
52 region, so that the diploid genotype may be represented Y...R, with two heterozygous sites Y (for
53 T/C) and R (for A/G) (Fig. 1). Suppose the reads are $14 \times T$ and $6 \times C$ at the first site, and $7 \times A$,
54 $10 \times G$, and $1 \times T$ at the second (with the single T to be most likely a sequencing error). The
55 haploid consensus sequence is constructed as T...G. In effect a heterozygote site with high quality
56 scores for the two nucleotides is represented as one consensus nucleotide with a low quality
57 score. Because it is largely pure chance which of the two nucleotides at a heterozygous site has
58 the greater number of reads, this strategy is equivalent to resolving the phase at random and using
59 only one of the constructed sequences. The resulting haploid consensus sequence may not be a
60 real biological sequence and may not represent the biology of the diploid individual. Besides loss
61 of information, a more serious problem is that the artefactual phased haploid sequence may be
62 unusually divergent from other sequences in the sample, potentially introducing systematic biases
63 in downstream inference. Currently constructing true diploid *de novo* assemblies is expensive. A
64 sequencing platform has been developed in combination with bioinformatic algorithms to
65 determine the true diploid genome sequence but the strategy still involves high cost (Weisenfeld

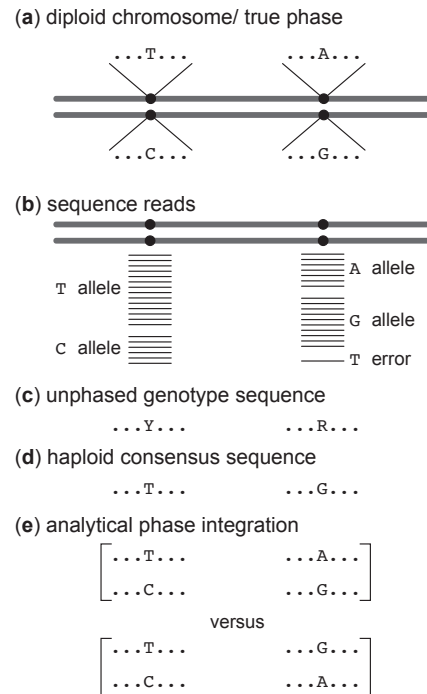


Fig. 1. Example of heterozygote phase resolution. (a) A hypothetical diploid chromosome with two heterozygous sites (T/C and A/G). The true haploid genotypes are T...A and C...G. (b) Sequence reads around the two heterozygous sites, assuming that they are far apart on the chromosome so that they are not present on any single read (in which case phase would be determined) while they are close enough to be on one locus. In this case genome assemblers should produce the unphased genotype sequence (c), using the IUPAC ambiguity codes to represent heterozygote sites, but instead they produce the so-called 'haploid consensus sequence' (d), picking up the most common nucleotide at each heterozygote site (T...G since T and G are by chance the most common sequence reads at the two sites), which may not match either of the true haploid sequences. (e) Analytical integration of phase resolution takes the unphased genotype sequences as data and averages over all possible phase resolutions, weighting each one appropriately according to their relative likelihood based on the whole sequence alignment at the locus.

66 *et al.*, 2017). If a read is long and fully covers a locus, multiple heterozygous sites in the same
 67 locus will be naturally phased. However, if the reads are short, and the two heterozygous sites do
 68 not occur in the same read, their genotypic phase resolution will become an issue.

69 How the heterozygote phase is resolved may have a significant impact on population
 70 genomic and phylogenomic inference using genomic sequence data. Phase information is
 71 well-known to be important for relating genotype to phenotype in human disease mapping
 72 (Tewhey *et al.*, 2011). Similarly, Gronau *et al.* (2011) found that use of an analytical integration
 73 method (which averages over all possible phase resolutions) leads to nearly identical performance
 74 as the use of true phase resolutions for estimating population parameters, and that random phase
 75 resolution produced unreliable estimates. Andermann *et al.* (2019) developed a bioinformatics
 76 pipeline to recover allelic sequences from sequence capture data, and found it to produce more
 77 accurate estimation of species divergence times under the MSC model (Rannala and Yang, 2003)
 78 than other strategies such as use of consensus haploid sequences, random phasing, or ambiguity
 79 encoding. Overall little is known about the effects of heterozygote phase resolution on many
 80 inference problems using multilocus genomic sequence data under the MSC model, including
 81 species tree estimation, estimation of population sizes and species divergence times, and
 82 inference of cross-species introgression/hybridization.

83 We have implemented in BPP (Flouri *et al.*, 2018) an analytical integration algorithm to

84 handle unphased diploid sequences, developed by Gronau *et al.* (2011) in their G-PhoCS
85 program, which is an orthogonal extension of an earlier version of BPP (Rannala and Yang, 2003;
86 Burgess and Yang, 2008). Previously Kuhner and Felsenstein (2000) implemented an Markov
87 chain Monte Carlo (MCMC) algorithm to average over different phase resolutions in the
88 likelihood calculation for estimating θ under the single-population coalescent. The algorithm was
89 found to mix slowly even for small datasets. The analytical integration algorithm uses a
90 data-augmentation strategy, in which the unknown fully resolved haploid sequences constitute the
91 complete data or latent variables, and enumerates and averages over all possible phase
92 resolutions, weighting them according to their likelihoods based on the whole sequence
93 alignment. For example, if a diploid sequence has two heterozygous sites, Y...R, the approach
94 will average over both phased genotypic resolutions: (i) T...A and C...G versus (ii) T...G and
95 C...A (Fig. 1). Note that there may be rich information about the phase resolution of any
96 unphased sequence in an alignment of many sequences, either from the same species or from
97 different but closely related species. Consider for example the phase resolutions for a human
98 diploid sequence Y...R (Fig. 1). If we observe in the chimpanzee fully resolved sequences T...A
99 and C...G (e.g., in an individual homozygous at both sites, with genotypes T/T...A/A) and never
100 observe sequences T...G and C...A, then very likely the human diploid sequence has the haploid
101 genotypes T...A and C...G. Our implementation of the algorithm works with all four analyses
102 under the MSC model in BPP (Yang, 2015; Flouri *et al.*, 2018, 2020b), including species tree
103 estimation (Yang and Rannala, 2014; Rannala and Yang, 2017) and species delimitation through
104 Bayesian model selection (Yang and Rannala, 2010, 2014; Leaché *et al.*, 2019). We also
105 implemented the algorithm under the multispecies-coalescent-with-introgression (MSci) model
106 (Flouri *et al.*, 2020a).

107 Here we use computer simulation to evaluate different phase-resolution strategies in terms
108 of their precision and accuracy in Bayesian species tree estimation under the MSC and in
109 parameter estimation under both the MSC and MSci models. In addition to using the true phase
110 resolution, which is generated during the simulation and is known with certainty, we also include
111 analytical phase integration (Gronau *et al.*, 2011; Flouri *et al.*, 2018), phase resolution using the
112 program PHASE (Stephens *et al.*, 2001; Stephens and Donnelly, 2003), and random resolution.
113 The strategy of random resolution is largely equivalent to the common method of using haploid
114 consensus sequences. The PHASE program was developed for population data from the same
115 species, but is here applied to unphased sequences from both within and between species. We
116 note that a number of computational phasing algorithms have been developed, such as
117 Haplotyper (Niu *et al.*, 2002) and fastPHASE (Scheet and Stephens, 2006). These are mostly
118 developed to improve the computational efficiency and to handle long sequences (Choi *et al.*,
119 2018), and are expected to produce similar results to PHASE in analysis of short sequences.

120 2. METHODS AND MATERIALS

121 *Simulation to Estimate Species Trees*

122 We use the program MCCOAL in BPP3.4 (Yang, 2015) or the simulate switch of BPP4.3 (Flouri
123 *et al.*, 2020b) to simulate gene trees and multi-locus sequence data using four fixed species trees
124 for eight species (Figs. 2a, a', b, & b'). The trees have very short branches, mimicking
125 challenging species trees generated during radiative speciation events. In the two deep trees

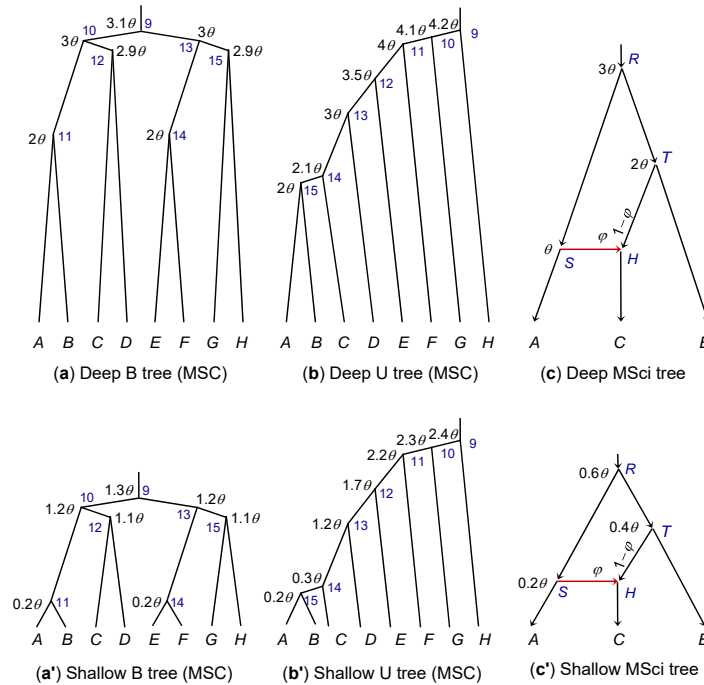


Fig. 2. (a & a') Deep and shallow balanced species trees, and (b & b') deep and shallow unbalanced species trees for eight species used for simulating data under the MSC model. (c & c') Deep and shallow species trees with introgression used to simulate data under the MSci model. The ages of internal nodes (τ s) are shown next to the nodes, with $\theta = 0.01$ (high rate) or 0.001 (low rate). The blue indexes at internal nodes of the tree are used to identify the parameters.

126 species divergences are much older than average coalescent times ($\theta/2$). In the two shallow trees,
 127 species divergences are very recent relative to coalescent times, mimicking different populations
 128 of the same species. Note that in this study, we make no distinction between species and
 129 populations. The MSC model has two sets of parameters: the species divergence times (τ s) and
 130 the population size parameters (θ s). Both are measured by the expected number of
 131 mutations/substitutions per site. For each species/population, $\theta = 4N\mu$, where N is the effective
 132 population size and μ is the mutation rate per site per generation. We consider two mutation
 133 rates, with $\theta = 0.001$ (low rate) or 0.01 (high rate), respectively, for all populations on the tree.
 134 The species divergence times (τ s) are given as multiples of θ . We consider 10, 20, 50, or 100
 135 loci, with each locus having 500 sites. On average there should be 0.5 and 5 heterozygous sites
 136 between the two sequences of any individual at the low and high rates, respectively. We sample
 137 $S = 2$ or 4 haploid sequences (or 1 or 2 diploid individuals) per species at each locus. Gene trees
 138 with branch lengths (coalescent times) are generated independently among loci using the MSC
 139 density given the species tree and parameters (Rannala and Yang, 2003). The JC model (Jukes
 140 and Cantor, 1969) is then used to 'evolve' the sequences along the gene tree to generate the
 141 sequence alignments at the tips of the tree. Analysis of this full dataset by BPP is strategy 'F'.

142 To simulate unphased diploid sequences, two sequences from the same species are
 143 combined into one diploid sequence, using the International Union of Pure and Applied
 144 Chemistry (IUPAC) ambiguity characters to represent heterozygous sites (for example, Y means
 145 a T/C heterozygote) (Fig. 1c). The data of unphased diploid sequences are analyzed using the
 146 diploid or phase option of the BPP program (strategy 'D'), which analytically averages over all

147 possible phase resolutions (Gronau *et al.*, 2011). With strategy ‘P’, we use the program PHASE
148 (Stephens *et al.*, 2001) to resolve the phase, and then analyze the phased sequences using BPP
149 (with 16 or 32 sequences in the alignment per locus for $S = 2$ and 4, respectively). Lastly, we use
150 random phase resolution, referred to as strategy ‘R’. The simulation program automatically
151 generates the sequence alignments for strategies F, D, and R. For strategy P, we ran PHASE 2.1
152 (Stephens *et al.*, 2001) to reconstruct the phased sequences for each locus, and used the PERL
153 program SeqPhase (Flot, 2010) to convert files.

154 The number of replicate datasets is 100. With four trees, two mutation rates ($\theta = 0.001$ or
155 0.01), two sampling configurations ($S = 2$ or 4), four numbers of loci ($L = 10, 20, 50, 100$), we
156 generated in total $4 \times 2 \times 2 \times 4 \times 100 = 6400$ datasets, each of which is analyzed using the four
157 strategies. The BPP program (Flouri *et al.*, 2018) was used in the analysis. Inverse-gamma priors
158 are assigned on parameters under the MSC model, with the shape parameter 3 so that the priors
159 are diffuse and with the mean to be close to the true value. We use $\theta \sim \text{IG}(3, 0.02)$ with mean
160 0.01 and $\tau_0 \sim \text{IG}(3, 0.08)$ with mean 0.04 for the age of the root of the species tree for data
161 simulated with the high rate ($\theta = 0.01$). For data of the low rate ($\theta = 0.001$), the priors are $\theta \sim$
162 $\text{IG}(3, 0.002)$ with mean 0.001 and $\tau_0 \sim \text{IG}(3, 0.008)$ with mean 0.004. The prior means for τ_0 are
163 close to the true values for the deep trees but are larger than the true values for the shallow trees,
164 although the priors are diffuse. For species tree estimation, we integrate out θ s analytically
165 through the use of the conjugate inverse-gamma priors. We conducted pilot runs to determine the
166 chain lengths needed for convergence. The final settings for the MCMC are 20,000 iterations for
167 burn-in, then taking 2×10^5 samples, sampling every 2 iterations.

168 Strategy P requires running the Bayesian MCMC program PHASE L times if there are L
169 loci in the dataset, to generate the fully resolved sequence alignments at the loci. This is
170 somewhat expensive if there is a large number of loci and the mutation rate is high resulting in
171 many heterozygous sites at each locus. After the datasets are generated, the BPP analysis of each
172 dataset by strategies F, P, and R involves about the same amount of computation. Strategy D is
173 more expensive as the method averages over all possible phase resolutions, which may involve
174 likelihood calculation for many site patterns, especially if there are many sequences per locus
175 with many heterozygous sites.

176 For species tree estimation (A01 analysis in Yang, 2015), we calculated the proportion
177 (among the 100 replicates) with which each node on the true species tree is found in the *maximum*
178 *a posteriori* (MAP) species tree in the BPP analysis. This is a measure of accuracy since the MAP
179 tree is the best ‘point estimate’ of the species tree (Rannala and Yang, 1996). We examined the
180 size and coverage probability of the 95% credibility set of species trees. The coverage probability
181 is the proportion among the 100 replicate datasets in which the credibility set includes the true
182 species tree. The size of the set indicates the precision or power of the method, but the method is
183 considered reliable only if the coverage probability exceeds the nominal 95%.

184 *Simulation to Estimate Parameters under the MSC Model*

185 The same data simulated under the MSC model for species tree estimation are analyzed using the
186 four phase-resolution strategies to estimate parameters in the MSC model (θ s and τ s), with the
187 species tree fixed. This is the A00 analysis in Yang (2015). We calculated the posterior means and
188 the 95% HPD CI intervals for each parameter and examine the relative root mean square error

189 (rRMSE), using the posterior means as point estimates. This is defined as

$$\text{rRMSE} = \frac{1}{\omega} \left[\frac{1}{R} \sum_{i=1}^R (\hat{\omega}_i - \omega)^2 \right]^{\frac{1}{2}}, \quad (1)$$

190 where ω is the true value of any parameter, and $\hat{\omega}_i$ its estimate (posterior mean) in the i th
191 replicate dataset, with $i = 1, \dots, R$ over the $R = 100$ replicates. For example, $\text{rRMSE} = 0.1$ means
192 that the mean square error is 10% of the true value. The rRMSE is a combined measure of bias
193 and variance.

194 *Simulation to Estimate Parameters under the MSci Model*

195 The MSci models for three species of Figures 2c&c' are assumed to generate gene trees and
196 sequence alignments using the `simulate` option of BPP4.3 (Flouri *et al.*, 2020a). The three
197 species have the phylogeny $(A, (C, B))$, but there was introgression from A to C at the time
198 $\tau_H = \tau_S$, with the introgression probability $\varphi = 0.1$ and 0.3. Other settings are the same as above
199 for the simulation under the MSC model. We consider two mutation rates (with $\theta = 0.001$ and
200 0.01) and four datasizes (with $L = 10, 20, 50$, and 100 loci), with each locus having 500 sites. We
201 sample either $S = 2$ or 4 sequences per species per locus. The JC model is used both to simulate
202 and to analyze the data.

203 For data simulated at the high rate ($\theta = 0.01$), the priors are $\theta \sim \text{IG}(3, 0.02)$ and $\tau_0 \sim$
204 $\text{IG}(3, 0.06)$ for the root age. At the low rate ($\theta = 0.001$), the priors are $\theta \sim \text{IG}(3, 0.002)$ and $\tau_0 \sim$
205 $\text{IG}(3, 0.006)$. A $\mathbb{U}(0, 1)$ prior is used for the introgression probability φ .

206 *Analyses of two real datasets*

207 We applied different phase-resolution strategies (D, P, and R) to analyze two previously
208 published datasets, one of East Asian brown frogs (Zhou *et al.*, 2012) and another of Rocky
209 Mountains chipmunks (Sarver *et al.*, 2021), to demonstrate that the effects discovered in the
210 simulations apply to real data analysis. With real data, the option of true phase resolution (F) is
211 unavailable, and the analytical phase resolution (D) is expected to perform the best. In addition,
212 we include an approach of treating heterozygote sites in the alignment as ambiguity characters in
213 the likelihood calculation, and refer to it as strategy 'A' (for ambiguity). This is considered a
214 mistaken approach of handling the data and is not included in our simulation, but we use it in the
215 real data analysis to illustrate its effects.

216 We re-analyzed a dataset of five nuclear loci from the East Asia brown frogs in the *Rana*
217 *chensinensis* species complex (Zhou *et al.*, 2012) to infer the species tree (the A01 analysis) and
218 to estimate the parameters under the MSC on the MAP tree (the A00 analysis). There are three
219 morphologically recognized species or four populations: *R. chensinensis* (clades C and L), *R.*
220 *kukunoris* (K) and *R. huanrensis* (H) (Fig. 3a). The dataset was previously analyzed by Yang
221 (2015), treating heterozygotes as ambiguities (strategy A). Each locus has 20-30 sequences, with
222 sequence lengths to be 285–498 sites. We assign inverse-gamma priors on parameters: $\theta \sim \text{IG}(3,$
223 $0.002)$ with mean 0.001 and $\tau_0 \sim \text{IG}(3, 0.004)$ with mean 0.002 for the root age. We used a
224 burnin of 8000 iterations, then taking 10^5 samples, sampling every two iterations. The same
225 analysis was run at least twice to confirm consistency between runs. This is a small dataset and

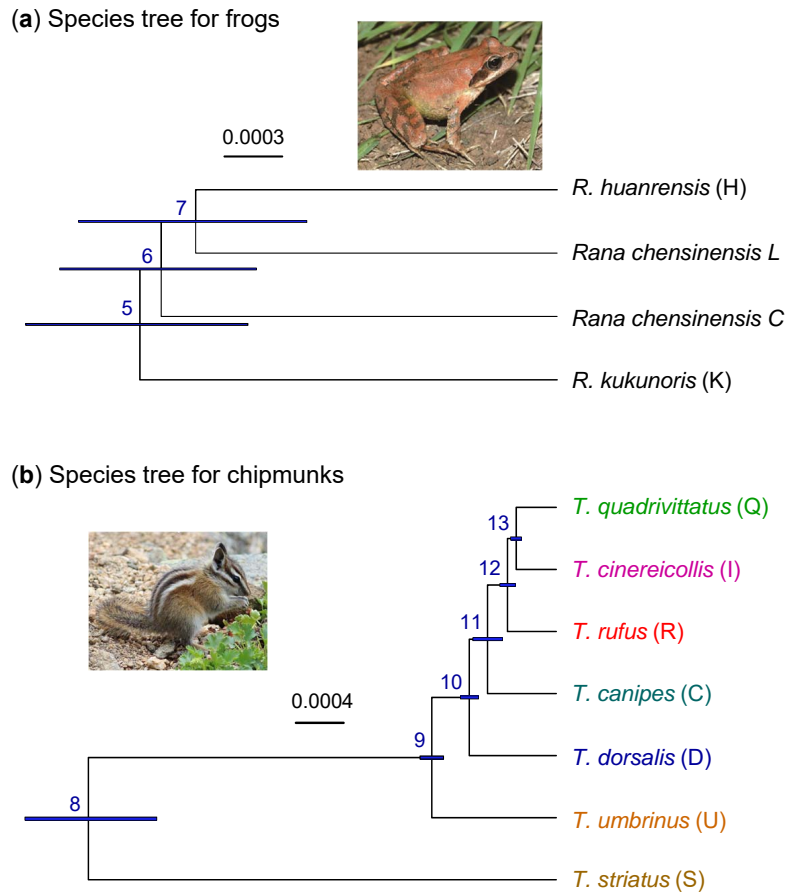


Fig. 3. Inferred species trees (a) for East Asian brown frogs and (b) for Rocky Mountains chipmunks. Branch lengths reflect the posterior means of divergence times, with branch bars representing the 95% HPD intervals, obtained under the MSC using the analytical phase integration algorithm (strategy D). Estimates of other parameters are in table 6.

226 the MCMC algorithm mixes well.

227 The second dataset consist of nuclear loci from six species of Rocky Mountains
228 chipmunks in the *Tamias quadrivittatus* group: *T. canipes* (C), *T. cinereicollis* (I), *T. dorsalis* (D),
229 *T. quadrivittatus* (Q), *T. rufus* (R), and *T. umbrinus* (U) (Fig. 3b). Sarver *et al.* (2021) used a
230 targeted sequence-capture approach to sequence 51 Rocky Mountains chipmunks from those six
231 species. As a reference genome assembly was lacking, reads were assembled iteratively into
232 contigs using an approach called “assembly by reduced complexity”. A dataset of 1060 nuclear
233 loci was compiled for molecular phylogenomic and introgression analyses, including 3
234 individuals from an outgroup species, *T. striatus*. High-quality heterozygotes, judged by mapping
235 quality and read depth, are represented in the alignments using the IUPAC ambiguity codes. The
236 filters applied by the authors suggest that the loci may be mostly coding exons or conserved parts
237 of the genome. The majority of loci have ≤ 5 variable sites (including the outgroup). We used the
238 first 500 loci in our analyses to infer the species tree and to estimate parameters under the MSC
239 model. We assigned inverse-gamma priors on parameters: $\theta \sim \text{IG}(3, 0.002)$ with mean 0.001 and
240 $\tau_0 \sim \text{IG}(3, 0.01)$ with mean 0.005 for the root age. In the A01 analysis (species tree estimation),
241 we used a burnin of 16000 iterations, then taking 2×10^5 samples, sampling every two iterations.

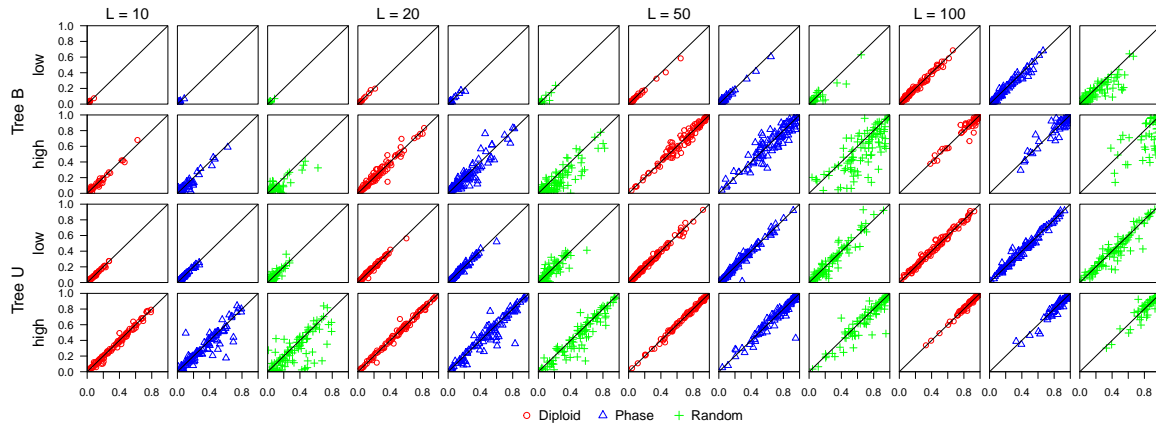


Fig. 4. (A01 under MSC, shallow tree, $S = 4$) Posterior probability for the true species tree for phase-resolution strategies D (diploid), P (PHASE) and R (random) plotted against the probability for strategy F (full data). The data are simulated under the MSC models with species trees Shallow B and Shallow U (Figs. 2a' & b'), with $S = 4$ sequences sampled per species. Each plot has 100 scatter points, for the 100 replicate datasets, with the x -axis to be the posterior probability for strategy F while the y -axis is for strategies D, P, or R. 'Low' ($\theta = 0.001$) and 'high' ($\theta = 0.01$) refer to the mutation rate, and $L (= 10, 20, 50, 100)$ is the number of loci. Results for other simulation settings are in Figures S1-S3.

242 The A00 analysis (parameter estimation on the MAP tree) used the same settings except that only
 243 10^5 samples were collected. The same analysis was run at least twice to confirm consistency
 244 between runs.

245 3. RESULTS

246 *Species Tree Estimation under the MSC Model*

247 Bayesian analysis of each replicate dataset using each of the four strategies produced a sample
 248 from the posterior distribution of the species trees, which we summarized to identify the
 249 maximum *a posteriori* probability (MAP) tree, and construct the 95% credibility set of species
 250 trees. The proportion, among the 100 replicates, with which the clades represented by those short
 251 branches were recovered in the MAP tree are shown in tables 1, S1-S3. Other clades on the trees,
 252 represented by longer branches, were recovered with probability near 100%, even for the low
 253 mutation rate and 10 loci. We also plotted the posterior probabilities for the true tree for the
 254 different phasing strategies in Figures 4, S1-S3. Strategy F, the analysis of the fully resolved
 255 haploid data, is expected to have the best performance and is thus the gold standard, against
 256 which the other strategies are compared.

257 In data simulated using the two deep trees (Deep B and Deep U) (Figs. 2a&b), the four
 258 phase-resolution strategies produced similar probabilities for recovering the true clades, with the
 259 differences among methods not being larger than the random sampling errors due to the limited
 260 number of replicates (tables S1 & S3). The different strategies most often produced the same
 261 MAP tree, although the posterior probability attached to the MAP tree varies somewhat among
 262 methods, but the differences are comparable to MCMC sampling errors. This can be seen in
 263 Figures S1 & S3, where the posterior for the true tree is plotted. Even random resolution (R)
 264 produced very similar results to the use of the fully resolved data (F). Note that in data simulated
 265 at the high rate, there are very likely to be two or more heterozygote sites in the diploid genotype

266 of each individual at any locus, and the switching error rate for random phase resolution, which is
267 the average proportion of heterozygous sites mis-assigned relative to the previous heterozygous
268 site (Stephens and Donnelly, 2003), is 50%. Even the PHASE program generates substantial
269 errors of phase resolution at the high mutation rate (table 2). Species tree estimation is thus robust
270 to considerable phasing errors when species divergences are much older than average coalescent
271 times.

272 For the two shallow trees (Figs. 2a' & b'), large differences were found among the four
273 strategies (tables 1 & S2, Figures 4 & S2). While strategy D produced results very similar to use
274 of the full data (F), both strategies P and R had poorer performance, especially at the high rate,
275 when strategy R produced larger CI set, with lower coverage than strategies F and D.

276 Thus phasing errors have different effects on species tree estimation depending on
277 whether the species tree is deep or shallow. We suggest that this may be explained by the
278 probability that the sequences from the same species coalesce before they reach the time of
279 species divergence, when one traces the genealogical history at each locus backwards in time. For
280 example, the probability that $S = 2$ sequences from species A coalesce before reaching the
281 common ancestor of A and B is $\mathbb{P}\{t_{\text{mrca}} < \tau_{AB}\} = 1 - e^{-4} \approx 0.982$ in the two deep trees and
282 $1 - e^{-0.4} \approx 0.330$ in the two shallow trees, while the corresponding probabilities for $S = 4$
283 sequences are 0.967 and 0.077 for the deep and shallow trees, respectively (Fig. S4). In the deep
284 trees, there is a high chance for all sequences from the same species to coalesce before reaching
285 species divergence, and then the problem will be similar to using the ancestral sequence for each
286 species (which is mostly determined by the most common nucleotides at the individual sites;
287 Yang *et al.*, 1995) for species tree estimation, a process that is not expected to be sensitive to
288 phasing errors. In the shallow species trees, there are high chances that sequences from the same
289 species may not have coalesced before reaching the time of species divergence, and sequences
290 with phasing errors will enter ancestral populations, interfering with species tree estimation.

291 While our main objective in this study is to evaluate the impacts of different phasing
292 strategies, it is worth noting the effects of other major factors on species tree estimation that are
293 obvious from our results (Figs. 4, S1–S3 and tables 1, S1–S3). By design species tree B is harder
294 to recover than tree U because tree B has four short branches (for clades C_{10} , C_{12} , C_{13} , and C_{15})
295 while tree U has only three (for clades C_{10} , C_{11} , and C_{15}) (Fig. 2). Thus tree B is recovered with
296 much lower probability than tree U by all methods in all parameter settings. We note that the
297 individual clades in tree B are recovered with lower probabilities than those in tree U (tables 1,
298 S1–S3). We speculate that this may be due to the fact that the four short branches in tree B are
299 close together (so that 945 trees around them are nearly equally good) while the three short
300 branches in tree U are far apart (so that only $3 \times 15 = 45$ trees around them are nearly equally
301 good). Because of the symmetry in tree B, the probabilities of recovering clades C_{10} and C_{13}
302 should be equal, as are those for C_{12} and C_{15} . Differences within each pair reflect the random
303 sampling errors due to our use of only 100 replicates. (Note that clades C_{11} and C_{14} were always
304 recovered in the simulation.)

305 The mutation rate had a dramatic impact on the precision and accuracy of species tree
306 estimation. At the higher rate (with $\theta = 0.01$ vs. 0.001), the credibility set was smaller, its
307 coverage was higher, and the MAP tree matched the true species tree with higher probability. In
308 our species trees, species divergence times (τ) are proportional to θ . This allows us to compare
309 the two values of θ , mimicking the use of conserved or variable regions of the genome for species

Table 1. (MSC A01, shallow, $S = 4$) Probabilities of recovering true clades and the size and coverage of the 95% credibility set of species trees when the true species tree is Shallow B and Shallow U (Figs. 2a' & b') and $S = 4$ sequences are sampled per species

Key	Species tree B							Species tree U					
	C_{10}	C_{12}	C_{13}	C_{15}	tree	CI cover	CI size	C_{10}	C_{11}	C_{15}	tree	CI cover	CI size
Low mutation rate													
F, 10L	0.34	0.25	0.27	0.22	0.00	0.66	233.8	0.49	0.47	0.63	0.12	0.96	84.4
D, 10L	0.35	0.25	0.25	0.22	0.00	0.66	233.4	0.47	0.44	0.61	0.09	0.96	82.7
P, 10L	0.34	0.25	0.25	0.20	0.00	0.68	235.5	0.47	0.47	0.61	0.12	0.96	82.9
R, 10L	0.36	0.24	0.26	0.24	0.01	0.66	225.7	0.47	0.46	0.61	0.13	0.94	81.8
F, 20L	0.46	0.29	0.34	0.26	0.00	0.73	178.3	0.52	0.53	0.74	0.22	0.97	33.5
D, 20L	0.45	0.29	0.34	0.22	0.01	0.73	175.4	0.54	0.52	0.72	0.23	0.97	33.8
P, 20L	0.46	0.27	0.38	0.26	0.01	0.73	178.2	0.55	0.53	0.76	0.26	0.97	33.1
R, 20L	0.47	0.26	0.31	0.26	0.02	0.70	168.8	0.53	0.50	0.74	0.22	0.97	32.6
F, 50L	0.56	0.44	0.56	0.46	0.07	0.90	86.3	0.60	0.65	0.95	0.40	0.97	11.4
D, 50L	0.56	0.45	0.53	0.47	0.06	0.90	83.3	0.59	0.61	0.95	0.40	0.95	11.5
P, 50L	0.61	0.45	0.52	0.48	0.08	0.87	89.0	0.59	0.61	0.93	0.40	0.98	11.7
R, 50L	0.52	0.37	0.57	0.46	0.06	0.86	80.8	0.65	0.63	0.91	0.41	0.96	11.6
F, 100L	0.72	0.75	0.81	0.74	0.33	0.99	25.5	0.75	0.77	0.99	0.60	0.99	7.0
D, 100L	0.75	0.74	0.81	0.76	0.34	0.98	25.7	0.75	0.76	1.00	0.59	0.99	6.8
P, 100L	0.74	0.71	0.80	0.75	0.34	0.97	26.4	0.74	0.78	1.00	0.59	0.98	7.1
R, 100L	0.73	0.66	0.75	0.72	0.26	0.96	27.4	0.75	0.74	0.98	0.57	0.99	7.9
High mutation rate													
F, 10L	0.68	0.58	0.68	0.49	0.19	0.92	91.5	0.70	0.76	0.96	0.53	1.00	11.2
D, 10L	0.70	0.58	0.66	0.50	0.18	0.92	95.3	0.72	0.75	0.94	0.52	1.00	11.8
P, 10L	0.66	0.58	0.63	0.48	0.14	0.91	94.6	0.68	0.76	0.92	0.53	0.98	12.5
R, 10L	0.53	0.54	0.64	0.45	0.12	0.89	108.4	0.71	0.77	0.78	0.44	0.97	13.3
F, 20L	0.91	0.74	0.90	0.72	0.43	0.99	22.2	0.80	0.85	1.00	0.72	1.00	5.7
D, 20L	0.92	0.75	0.88	0.72	0.43	1.00	23.3	0.81	0.85	1.00	0.72	1.00	6.0
P, 20L	0.91	0.70	0.89	0.72	0.41	1.00	27.2	0.77	0.86	0.97	0.68	1.00	6.6
R, 20L	0.86	0.71	0.79	0.66	0.31	0.98	30.1	0.79	0.84	0.93	0.64	0.99	7.3
F, 50L	1.00	0.97	1.00	0.94	0.91	1.00	4.1	0.90	0.97	1.00	0.87	1.00	2.6
D, 50L	1.00	0.97	1.00	0.94	0.91	1.00	4.1	0.90	0.97	1.00	0.87	1.00	2.6
P, 50L	1.00	0.94	1.00	0.92	0.86	1.00	4.3	0.91	0.97	1.00	0.88	1.00	2.7
R, 50L	0.98	0.91	0.94	0.90	0.76	1.00	5.6	0.92	0.97	1.00	0.89	0.99	2.9
F, 100L	1.00	0.99	1.00	1.00	0.99	1.00	1.6	1.00	0.98	1.00	0.98	1.00	1.7
D, 100L	1.00	0.99	1.00	1.00	0.99	1.00	1.6	1.00	0.98	1.00	0.98	1.00	1.6
P, 100L	1.00	0.99	1.00	0.99	0.98	1.00	1.6	0.99	0.98	1.00	0.97	1.00	1.7
R, 100L	1.00	0.98	1.00	0.99	0.97	1.00	2.0	1.00	0.98	1.00	0.98	1.00	1.7

Note.— The two mutation rates are low ($\theta = 0.001$) and high ($\theta = 0.01$), while 10L, 20L, 50L, 100L are the number of loci. C_{10} , C_{12} , etc. are probabilities of recovering the true clades on the species trees, while ‘tree’ is the probability of recovering the whole tree. ‘CI size’ is the number of species trees in the 95% credibility set and ‘CI cover’ is the probability that the set contains the true species tree. Results for other simulation settings are in tables S1-S3.

310 tree estimation. Our study focuses on closely related species with highly similar sequences, and
 311 data simulated at the high rate contain more variable sites and more phylogenetic information.

312 The number of loci similarly had a huge impact on species tree estimation. With more
 313 loci, inference became more precise (with smaller credibility set) and more accurate (with the
 314 MAP tree matching the true tree with greater probability). Increasing the number of loci by 10
 315 fold improves performance for all strategies more than increasing the mutation rate by the same
 316 factor.

317 The number of sequences sampled per species had consistent but relatively small effects
 318 on species tree estimation. Changing $S = 2$ to 4 improved the probabilities of recovering the true

Table 2. Average switching error rate for datasets simulated under the MSC and MSci models in this study

Model	PHASE (P)				Random (R)			
	low		high		low		high	
	$S = 2$	$S = 4$	$S = 2$	$S = 4$	$S = 2$	$S = 4$	$S = 2$	$S = 4$
MSC, Deep B	0.485	0.327	0.499	0.371	0.505	0.504	0.499	0.501
MSC, Deep U	0.489	0.332	0.501	0.370	0.488	0.488	0.498	0.499
MSC, Shallow B	0.448	0.370	0.459	0.349	0.488	0.495	0.498	0.498
MSC, Shallow U	0.390	0.304	0.430	0.331	0.489	0.505	0.501	0.502
MSci, Deep ($\varphi = 0.1$)	0.480	0.317	0.492	0.363	0.500	0.492	0.500	0.502
MSci, Deep ($\varphi = 0.3$)	0.482	0.311	0.494	0.360	0.520	0.490	0.501	0.499
MSci, Shallow ($\varphi = 0.1$)	0.402	0.342	0.461	0.346	0.496	0.489	0.492	0.498
MSci, Shallow ($\varphi = 0.3$)	0.402	0.331	0.454	0.337	0.498	0.502	0.502	0.501

Note.— Data of $L = 100$ loci are used in the calculation although the error rate does not depend on the number of loci. The same data generated under the MSC model are used in the A01 (species tree estimation) and A00 (parameter estimation) analyses. Note that the error rate for random phase resolution (R) is expected to be 0.5.

319 clades in the true species tree, reduced the CI set size, and improved the coverage of the CI set,
 320 but the improvements are in general small.

321 It is noteworthy that the coverage of the 95% CI set was below the nominal 95% in small
 322 or uninformative datasets while above 95% in large and informative datasets. In the case of 10
 323 loci at the low rate for tree Deep B, coverage was even below 50% (table S1). Even though the set
 324 included nearly 500 trees, more than a half of the CI sets failed to include the true tree. In
 325 contrast, at the high mutation rate and with 50 or 100 loci, CI coverage was often 100%. The
 326 method is over-confident in small and uninformative datasets and conservative in large and
 327 informative ones. The same pattern was noted in a previous simulation examining the information
 328 content in phylogenomic datasets (Huang *et al.*, 2020, table 3). Note that in our simulation, the
 329 replicate datasets are generated under a fixed model (species tree) and fixed parameter values, so
 330 that we are evaluating the Frequentist properties of Bayesian model selection, and a match is not
 331 expected (Huelsenbeck and Rannala, 2004; Yang and Rannala, 2005). Yet the large discrepancies
 332 are striking.

333 *Estimation of Divergence Times and Population Sizes under the MSC Model*

334 **The impact of the phasing strategies.** The same datasets simulated for species tree estimation
 335 were analyzed to estimate the parameters in the MSC model (θ s and τ s) with the species tree
 336 fixed (Figs. 2a, a', b & b'). The posterior means and 95% HPD CI for the 100 replicates are
 337 plotted in Figures 5, S5–S11, while the relative root mean square errors (rRMSE) are presented in
 338 tables S4–S11. Whereas the rRMSE reflects both biases and variances in parameter estimation,
 339 the datasets generated by the four phase-resolution strategies have about the same size in terms of
 340 the number of loci, the number of sequences per locus, and the number of sites per sequence, so
 341 that the sampling errors or variances are similar among methods and the differences in rRMSE
 342 mainly reflect differences in bias. Furthermore, we may use the symmetry of species tree B to
 343 gauge the magnitude of random sampling errors due to our use of 100 replicates: for instance,
 344 rRMSE should be equal for θ_A , θ_B , θ_E and θ_F , and for τ_{10} and τ_{13} , on the balanced trees.

345 The four phase-resolution strategies (F, D, P, and R) performed similarly for the Deep
 346 trees at the lower rate and when only $S = 2$ sequences (or one individual) are sampled per species.

IMPORTANCE OF HETEROZYGOTE PHASE RESOLUTION

13

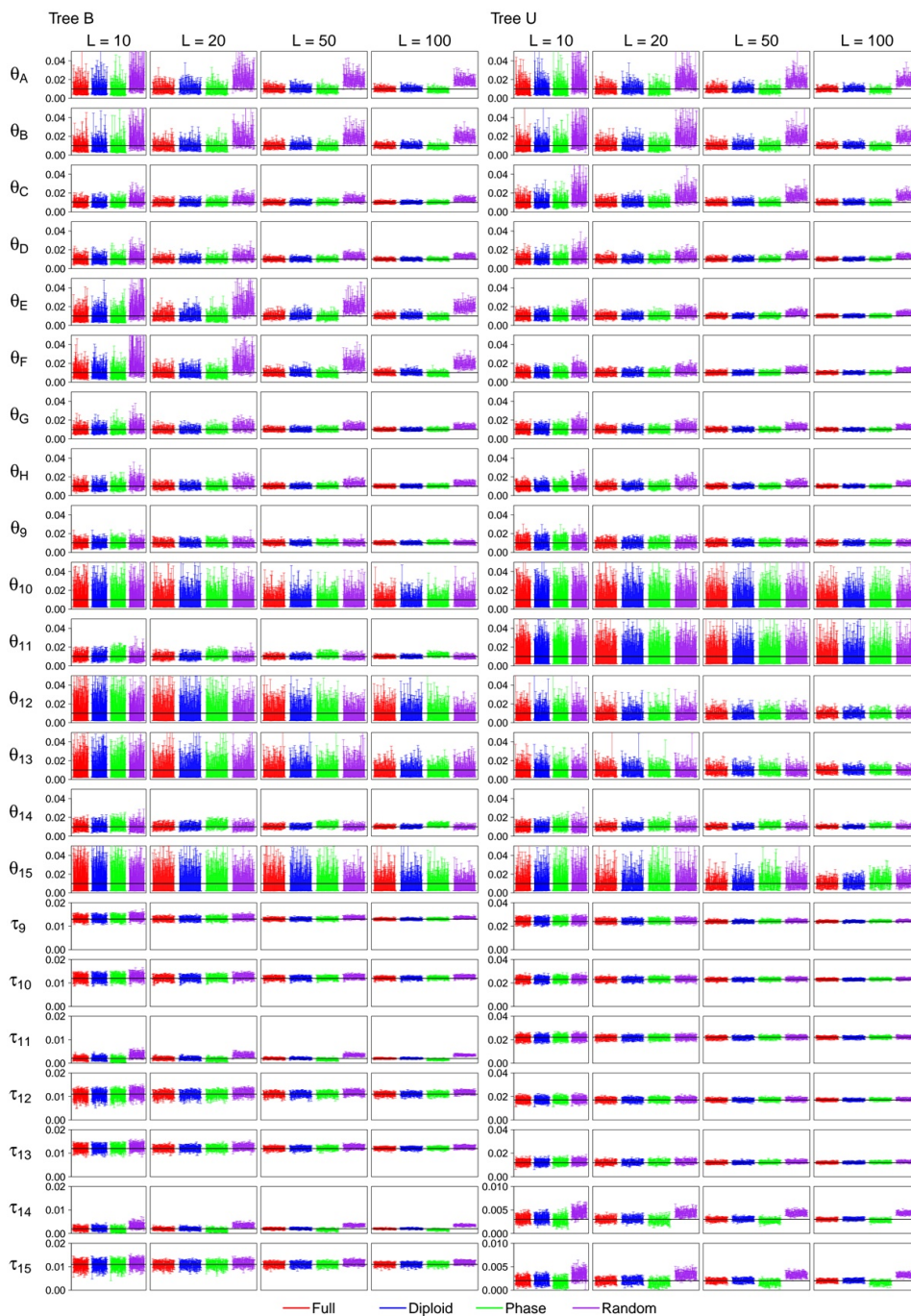


Fig. 5. (MSC, high rate, shallow, $S = 4$) The 95% HPD CIs for parameters for four phase-resolution strategies: F (the full data), D (diploid), P (PHASE), and R (random) in 100 replicate datasets simulated under MSC model trees Shallow B and Shallow U (Figs. 2a' & b'), at the high mutation rate ($\theta = 0.01$) and $S = 4$ sequences per species. The horizontal black lines indicate the true values. Results for other simulation settings are in Figures S5-S11.

347 We note that with $S = 2$ and at the low mutation rate (with heterozygosity at $\theta = 0.001$), there
 348 will be on average 0.5 heterozygous sites at the same locus, and the probability of having two or
 349 more heterozygous sites is $1 - 0.999^{500} - 500 \cdot 0.999^{499} \cdot 0.001 = 0.0901$. Then phase resolution

will not be a serious issue, and all four strategies examined in the study will be nearly equivalent.

At the high mutation rate ($\theta = 0.01$) for the Shallow trees, differences were noted among the strategies even for $S = 2$ sequences (Fig. S6 and tables S5 & S7). The PHASE program produced underestimates for the youngest species divergence times (τ_{11} and τ_{14} on Shallow B and τ_{15} on Shallow U) (Fig. 2a' & b'). The biases became more pronounced when $S = 4$ sequences per species are in the sample (Fig. 5 and tables S9– S11). At the high rate, there are on average 5 heterozygotes per locus in the individual and the probability of having two or more heterozygotes at the locus is 96%. Two factors may be responsible for the bias. First the PHASE program may have inferred heterozygote phase incorrect (indeed the error rate is comparable to that of random phasing with $S = 2$). Second PHASE is an MCMC program generating a distribution of different phase resolutions but we used only the optimal resolution, which may lead to underestimation of sequence divergences.

At the high rate and for shallow trees, random phasing (R) also created serious biases, but the biases are in the opposite direction. Random phasing overestimated the youngest species divergence times (τ_{11} and τ_{14} on Shallow B and τ_{15} on Shallow U), and overestimated θ for all modern species. The underestimation of modern θ is most striking, and occurred for both deep and shallow species trees at the high rate and is more dramatic with more sequences ($S = 4$ rather than 2) or more loci.

We examined the number of distinct site patterns in the alignment at each locus for the high-rate data (Fig. S12). Site patterns are compressed for the JC model, so that one site pattern is constant while the others are variable (Yang, 2006, p.144), and the number is thus an indication of the level of sequence divergence. At almost every locus, the PHASE program (P) produced alignments with fewer distinct site patterns than the true phase resolution (for example, with the mean to be 36.07 compared with the true value 38.51 on tree B), apparently because we used the optimal phase resolution inferred by the program and ignored the less likely ones. Random resolution produced about the same number of site patterns as the true number (average 38.36 vs. 38.51 for tree Deep B). The number of site patterns is thus not the reason for the poor performance of random phasing.

Note that calculation of the heterozygosity for each diploid individual, which is simply the proportion of heterozygous sites in the two sequences at the locus, does not rely on phase resolution. If we calculate the heterozygosity for each diploid individual and then average over individuals of the same species, we will get a reasonably good estimate of θ for that species. However, in the gene-tree based analysis conducted in BPP, each randomly phased haploid sequence is compared not only with the other sequence from the same individual, but also with sequences from other individuals through the use of a gene tree relating all phased haploid sequences at the locus. While the true haploid sequences may all be closely related, random phase resolution may generate chimeric sequences that are very different from naturally occurring fully resolved sequences, inflating apparent coalescent times and genetic diversity in the population. This effect is expected to be more serious when more individuals are included in the sample.

Estimation of θ for a single species. To explore this interpretation, we conducted a small simulation sampling independent loci from a single species to estimate the only parameter θ (Fig. 6, table 3). With $S = 2$ sequences per locus (one diploid individual), the four phase-resolution strategies are equivalent. However, with the increase of S , the strategy of random phase resolution becomes increasingly biased. Previously Felsenstein (1992) examined

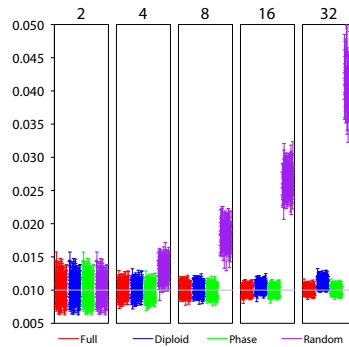


Fig. 6. The 95% HPD CIs for parameter θ in the single-population coalescent model in 100 replicate datasets using four phase-resolution strategies: F (the full data), D (diploid), P (PHASE), and R (random). There are 100 independent loci in each dataset, and at each locus there are S sequences of 500 sites (or $S/2$ diploid individuals), with $S = 2, 4, 8, 16,$ and 32 . The true parameter value is 0.01.

394 the efficiency of two summary methods based on the number of segregating (variable) sites ($\hat{\theta}_S$;
395 Watterson, 1975) and the average pairwise distance ($\hat{\theta}_\pi$; Tajima, 1983), relative to the maximum
396 likelihood (ML) method based on gene genealogies. He found that the summary methods ($\hat{\theta}_S$ and
397 $\hat{\theta}_\pi$) were much less efficient than the ML estimate, with orders-of-magnitude differences in the
398 variance in large samples (Felsenstein, 1992, tables 1 and 2), indicating that there is much
399 information about θ in the genealogical histories. The ML method should be very similar to BPP
400 here as both are full likelihood methods. Here we note that the number of segregating sites does
401 not depend on phase resolutions, and similarly the average proportion of different sites, averaged
402 over all the $S(S-1)/2$ pairwise comparisons, depends on the site configurations at each variable
403 site (such as 10 Ts and 4 Cs) but not on the genotypic phase between different heterozygous sites.
404 Both Watterson's estimator and the average pairwise distance are thus unaffected by phasing
405 errors. It is also noteworthy that those two simple methods are not affected by recombination
406 within the locus, while coalescent-based methods are (Felsenstein, 2019). While it is not
407 unexpected that a full likelihood method may be more sensitive to certain errors in the model or
408 in the data than heuristic methods, in this case it is striking that the systematic bias is so large
409 (with estimates to be several times larger than the true value) when the coalescent-based method
410 is applied to randomly phased sequences.

411 Felsenstein's (1992) analysis, as mentioned above, assumed knowledge of the true gene
412 trees and coalescent times (or equivalently infinitely long sequences at each locus). Here BPP is
413 applied to sequence alignments and accommodates uncertainties in the genealogical trees. The
414 different methods then have much more similar performance (table 3, $\hat{\theta}_S$, $\hat{\theta}_\pi$ and BPP strategy F),
415 suggesting that the uncertainties in the genealogical trees due to mutational variations in the
416 sequences have eroded much of the information in the gene trees. The summary methods (in
417 particular, $\hat{\theta}_\pi$) have larger variances than the BPP estimates, especially in large samples of $S = 32$
418 sequences, but the differences are relatively small. We also note that analytical phase integration
419 (D) produced variances that are nearly identical to those for the use of the full data (F).

420 **Impacts of other factors on parameter estimation under the MSC model.** We note
421 that different parameters are estimated with very different precision and accuracy, reflecting the
422 different amount of information in the data. Population size parameters (θ s) for modern species
423 are well estimated, as well as θ_0 for the root population, but θ s for other ancestral species,

Table 3. Mean and standard deviation ($\times 10^{-3}$) of estimates of θ for a single population (true value is 0.01) from a sample of S sequences using BPP with different strategies of phase resolution and two summary methods

Method	$S = 2$	$S = 4$	$S = 8$	$S = 16$	$S = 32$
BPP (F)	10.06 \pm 1.02	10.06 \pm 0.61	10.03 \pm 0.52	9.96 \pm 0.36	10.03 \pm 0.34
BPP (D)	10.06 \pm 1.02	10.05 \pm 0.62	10.17 \pm 0.53	10.50 \pm 0.43	11.19 \pm 0.47
BPP (P)	10.06 \pm 1.02	9.80 \pm 0.61	9.84 \pm 0.51	9.94 \pm 0.37	10.05 \pm 0.34
BPP (R)	10.06 \pm 1.02	12.86 \pm 0.91	18.13 \pm 1.32	26.43 \pm 1.65	41.27 \pm 3.22
Watterson ($\hat{\theta}_S$)	9.94 \pm 1.01	9.92 \pm 0.61	9.85 \pm 0.55	9.76 \pm 0.40	9.82 \pm 0.36
Pairwise distance ($\hat{\theta}_\pi$)	9.94 \pm 1.01	9.94 \pm 0.63	9.87 \pm 0.63	9.78 \pm 0.50	9.93 \pm 0.55
Pairwise distance ($\hat{\theta}'_\pi$)	10.01 \pm 1.03	10.01 \pm 0.64	9.94 \pm 0.64	9.84 \pm 0.50	9.99 \pm 0.56

Note.— Watterson’s estimate ($\hat{\theta}_S$) and the average pairwise distance ($\hat{\theta}_\pi$) do not depend on phase resolutions. JC correction is applied in calculation of $\hat{\theta}'_\pi$.

424 especially those represented by very short branches (e.g., $\theta_{10}, \theta_{13}, \theta_{12}, \theta_{15}$ in tree B) have large
 425 errors (Figs. 5, S5–S11). Species divergence times are all well estimated, with rRMSE to be even
 426 much smaller than those for population size parameters for modern species (tables S4–S11).

427 Both the mutation rate and the number of loci had a major impact on the estimation of the
 428 parameters. For all phasing strategies increasing the number of loci by 10 fold improves
 429 performance more than increasing the mutation rate by the same factor (Figs. 5, S5–S11, tables
 430 S4–S11).

431 *Estimation of Introgression Probability under the MSci Model*

432 We used the MSci models of Figure 2c&c' to simulate sequence data and used BPP to analyze
 433 them to estimate parameters in the MSci model. We are in particular interested in whether the
 434 different strategies of heterozygote phase resolution may lead to biases in the estimation of the
 435 timing (τ_H) and strength of the introgression (ϕ). The results are summarized in Figures 7 &
 436 S13–S19 and tables 4 & S12–S18.

437 As before, the diploid strategy (D) produced results almost indistinguishable from the use
 438 of the full data (F) in all parameter settings. The performance of the PHASE program (P) and
 439 random phasing (R) depends on the mutation rate and, to an lesser extent, on the number of
 440 sequences per species S . At the low rate, and in particular with only $S = 2$ sequences per species,
 441 all four strategies have similar performance, but large differences were found at the high mutation
 442 rate. Strategy R overestimates the modern θ and the species divergence times (τ) at the high rate,
 443 with the bias being more serious for $S = 4$ sequences than for $S = 2$. This is the same behavior as
 444 discussed earlier in the simulation under the MSC model. Strategy R also tends to overestimate
 445 ϕ , but the bias is small. Strategy P had the opposite bias and produced underestimates of modern
 446 θ and species divergence times when the mutation rate is high, with smaller biases than for
 447 strategy R. Strategy P also underestimates the introgression probability (ϕ).

448 An interesting question is whether each method detects introgression. We calculated the
 449 proportion of replicates in which the lower limit of the 95% HPD CI for ϕ exceeds a small value,
 450 set somewhat arbitrarily at 0.001. If the CI excludes the small value, we may take it as evidence
 451 that $\phi = 0$ is ruled out so that there is significant evidence for introgression. By this measure of
 452 power of the Bayesian ‘test’, strategies D and P had nearly identical power as the use of the full
 453 data (F), while random resolution (R) had reduced power at the high mutation rate (tables 5 &

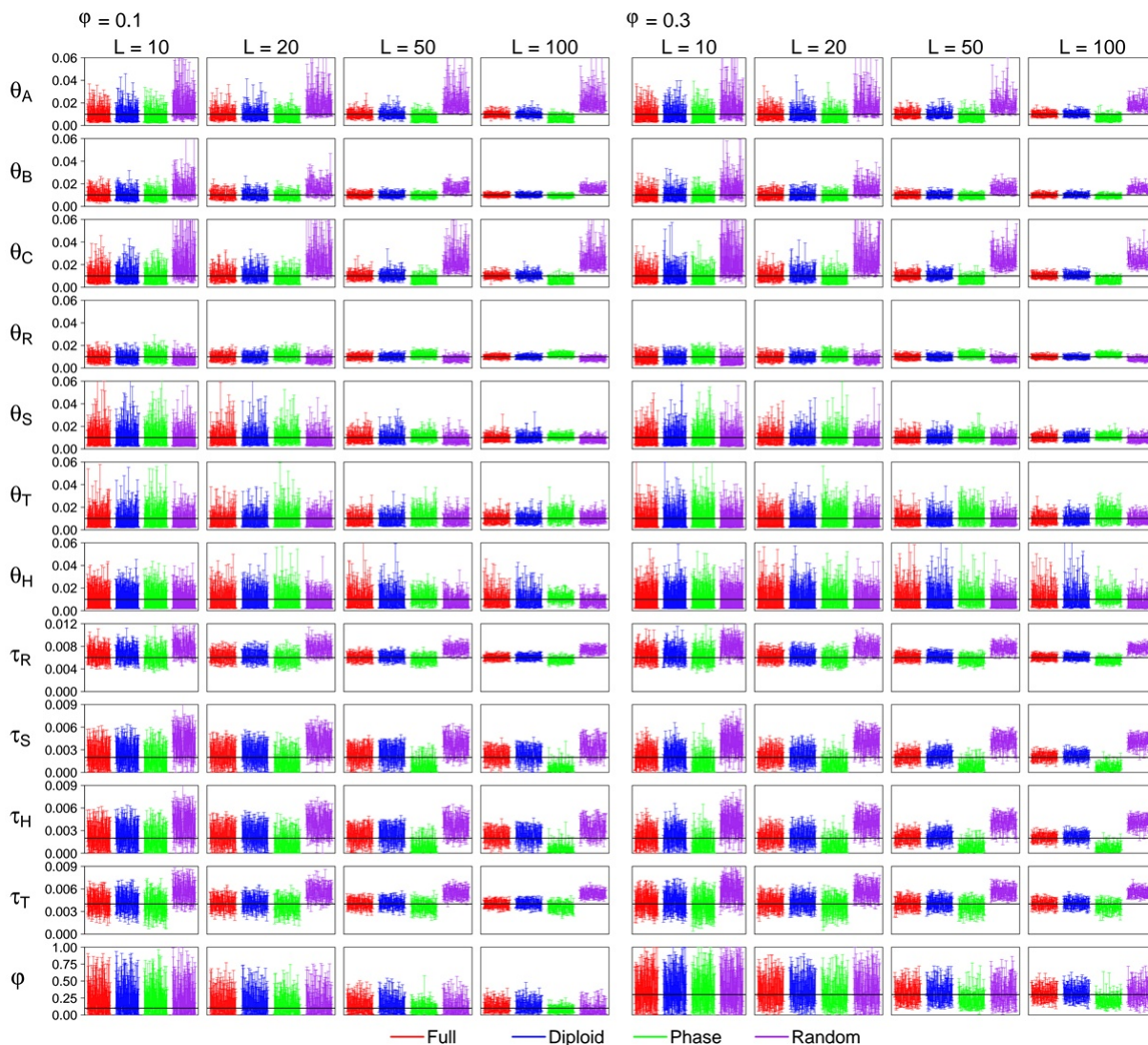


Fig. 7. (MSci model, high rate, Shallow, $S = 4$) The 95% HPD CIs for parameters under the MSci model of Figure 2c' when $S = 4$ sequences are sampled per species. Results for $S = 2$ are in Figure S14. See legend to Figure 5.

454 S19). Overall, power was very high even with only 10-20 loci and at the low mutation rate.
 455 Having more sequences is noted to boost the power of the test for all phase-resolution strategies.

Running Time for Different Analyses

457 The running time for the A01 analysis under the MSC model (species tree estimation) for the four
 458 phasing strategies (F, D, P, and R), averaged over the 100 replicates, is plotted against the number
 459 of loci in Figure S20. Running time increases nearly linearly with the number of loci, with the
 460 slope being steeper when $S = 4$ sequences are sampled per species than for $S = 2$. The diploid
 461 integration algorithm (D) has the longest running time. Note that the number of parameters in the
 462 MSC model, the number of loci, the number of sequences etc. are identical for the four strategies,

Table 4. (MSci A00 $S = 4$, high rate, shallow) Relative root mean square error (rRMSE) for parameter estimates under the Deep MSci model (fig. 2c') with $\phi = 0.1$ or 0.3 at the high mutation rate and $S = 4$

Truth	θ_A 1	θ_B 1	θ_C 1	θ_R 1	θ_S 1	θ_T 1	θ_H 1	τ_R 3	τ_S 1	τ_T 2	0.1	ϕ 0.3
$\phi = 0.1$												
F, 10L	0.28	0.24	0.34	0.24	0.45	0.36	0.25	0.15	0.45	0.17	0.42	-
D, 10L	0.31	0.25	0.36	0.24	0.48	0.36	0.22	0.16	0.46	0.18	0.41	-
P, 10L	0.27	0.22	0.28	0.27	0.50	0.51	0.27	0.15	0.45	0.22	0.43	-
R, 10L	0.99	0.86	1.49	0.31	0.37	0.25	0.23	0.36	1.12	0.46	0.47	-
F, 20L	0.24	0.21	0.28	0.18	0.40	0.35	0.28	0.10	0.41	0.13	0.46	-
D, 20L	0.26	0.22	0.31	0.18	0.43	0.33	0.26	0.11	0.43	0.14	0.45	-
P, 20L	0.24	0.20	0.26	0.26	0.39	0.49	0.38	0.11	0.46	0.17	0.54	-
R, 20L	1.09	0.73	1.75	0.27	0.32	0.25	0.22	0.30	1.04	0.41	0.53	-
F, 50L	0.17	0.12	0.18	0.13	0.24	0.27	0.39	0.07	0.33	0.08	0.53	-
D, 50L	0.18	0.13	0.20	0.13	0.27	0.27	0.32	0.07	0.38	0.09	0.51	-
P, 50L	0.31	0.13	0.34	0.22	0.21	0.41	0.29	0.09	0.60	0.14	0.68	-
R, 50L	1.18	0.63	1.50	0.22	0.23	0.26	0.24	0.26	0.90	0.38	0.64	-
F, 100L	0.12	0.08	0.13	0.09	0.18	0.23	0.32	0.04	0.27	0.06	0.59	-
D, 100L	0.14	0.08	0.16	0.09	0.19	0.24	0.30	0.04	0.33	0.06	0.57	-
P, 100L	0.39	0.11	0.42	0.21	0.15	0.42	0.23	0.09	0.72	0.14	0.74	-
R, 100L	1.27	0.59	1.60	0.19	0.20	0.21	0.23	0.23	0.73	0.35	0.73	-
$\phi = 0.1$												
F, 10L	0.30	0.27	0.31	0.21	0.35	0.35	0.30	0.15	0.34	0.16	-	0.48
D, 10L	0.31	0.29	0.41	0.22	0.40	0.35	0.28	0.16	0.35	0.17	-	0.49
P, 10L	0.27	0.24	0.27	0.22	0.48	0.55	0.35	0.15	0.48	0.21	-	0.44
R, 10L	1.03	0.93	1.83	0.32	0.37	0.28	0.24	0.36	1.11	0.49	-	0.57
F, 20L	0.24	0.18	0.26	0.18	0.35	0.28	0.37	0.09	0.30	0.12	-	0.39
D, 20L	0.28	0.20	0.33	0.19	0.33	0.27	0.33	0.10	0.33	0.13	-	0.40
P, 20L	0.26	0.17	0.29	0.21	0.45	0.47	0.42	0.10	0.52	0.18	-	0.37
R, 20L	1.04	0.69	1.69	0.30	0.30	0.26	0.24	0.31	1.08	0.44	-	0.51
F, 50L	0.18	0.12	0.16	0.11	0.22	0.28	0.48	0.05	0.21	0.10	-	0.27
D, 50L	0.20	0.12	0.19	0.11	0.21	0.29	0.49	0.06	0.25	0.10	-	0.30
P, 50L	0.28	0.13	0.34	0.20	0.26	0.53	0.56	0.09	0.60	0.17	-	0.33
R, 50L	0.97	0.61	1.56	0.24	0.21	0.24	0.25	0.29	1.07	0.41	-	0.38
F, 100L	0.11	0.09	0.11	0.08	0.15	0.20	0.42	0.04	0.15	0.08	-	0.20
D, 100L	0.12	0.10	0.14	0.08	0.14	0.20	0.49	0.04	0.19	0.08	-	0.21
P, 100L	0.34	0.12	0.41	0.20	0.19	0.47	0.37	0.09	0.66	0.15	-	0.34
R, 100L	0.87	0.59	1.50	0.22	0.18	0.17	0.30	0.27	1.06	0.39	-	0.32

Note.— Truth represents the true parameter values used in the simulation; values for θ and τ are $\times 10^{-2}$. Results for other simulation settings are in tables S12-S18.

463 so that their computational load is proportional to the number of site patterns. As strategy D
 464 enumerates all possible phase resolutions (including the true resolution), which may result in
 465 many distinct site patterns, it is more expensive than the other methods. The running time for
 466 each BPP analysis on a single core ranged from ~ 20 minutes for 10 loci to ~ 5 hours for strategy
 467 D with data of 100 loci. Strategy P involves running the Bayesian MCMC program PHASE for
 468 each of the L loci. At the low mutation rate with very few heterozygous sites per locus, this
 469 requires minimal computation (Fig. S21), but at the high rate and with $S = 4$ sequences per
 470 species, the running time can be comparable with running the subsequent BPP analyses.

471 The running time for the A00 analysis (parameter estimation) under the MSC and MSci
 472 models is shown in Figures S22–S25. The A00 analysis under the MSC involves less computation
 473 than the A01 analysis as there is no MCMC moves to change the species tree. Overall, the same
 474 patterns are observed as discussed above for the A01 analysis. Note that the computer cluster used

Table 5. (MSci test, shallow) Power of the Bayesian test for introgression (measured by the proportion of replicates in which the lower limit of the 95% HPD CI for ϕ is > 0.001) when the true model is the Shallow MSci tree

	low				high			
	10L	20L	50L	100L	10L	20L	50L	100L
$\phi = 0.1$								
$S = 2$ seqs per species								
F	0.42	0.41	0.42	0.49	0.48	0.52	0.64	0.89
D	0.48	0.40	0.45	0.46	0.38	0.35	0.61	0.80
P	0.49	0.45	0.49	0.49	0.55	0.60	0.83	0.99
R	0.51	0.44	0.47	0.48	0.27	0.25	0.40	0.57
$S = 4$ seqs per species								
F	0.58	0.47	0.55	0.58	0.59	0.71	0.98	0.99
D	0.56	0.44	0.54	0.60	0.56	0.66	0.94	0.99
P	0.57	0.43	0.44	0.49	0.60	0.65	0.94	1.00
R	0.56	0.45	0.50	0.57	0.44	0.46	0.71	0.81
$\phi = 0.3$								
$S = 2$ seqs per species								
F	0.68	0.74	0.87	0.94	0.86	0.97	1.00	1.00
D	0.64	0.81	0.89	0.93	0.84	0.91	1.00	1.00
P	0.68	0.78	0.85	0.89	0.86	0.94	0.99	1.00
R	0.66	0.72	0.87	0.92	0.74	0.82	0.96	0.99
$S = 4$ seqs per species								
F	0.79	0.88	0.97	1.00	0.95	1.00	1.00	1.00
D	0.83	0.86	0.95	1.00	0.92	0.99	1.00	1.00
P	0.84	0.84	0.95	1.00	0.94	0.97	1.00	1.00
R	0.84	0.82	0.94	0.99	0.75	0.89	1.00	1.00

Note.— Results for the Deep MSci model are in table S19.

475 in this work consist of computers with different processors, so there may be random fluctuations
 476 in running time due to the different jobs being assigned to different processors. For example the
 477 differences in Figures S21 & S23 reflect this random fluctuation as the data were the same.

478 Analysis of two real datasets

479 We analyzed two real datasets using four different phase-resolution strategies: D (diploid), P
 480 (PHASE), R (random), and A (ambiguity). With real data, the option of true phase resolution (F)
 481 is unavailable, and the analytical phase resolution (D) is expected to have the best performance,
 482 against which we compare the other strategies.

483 **East Asia brown frogs.** We re-analyzed a dataset of five nuclear loci from the East Asia
 484 brown frogs in the *Rana chensinensis* species complex (Zhou *et al.*, 2012) (Fig. 3a). This dataset
 485 was previously analyzed by Yang (2015) using strategy A. The number of site patterns at each
 486 locus is 18–26 for strategy A, and 22–102 for strategy D. Running time using one thread on our
 487 server was 3 mins for A, 7-8 mins for P and R, and 12 mins for D.

488 In the A01 analysis (species tree estimation), the four strategies (D, P, R, and A) produced
 489 the same MAP tree (Fig. 3a): ((H, L), C), K), with the posterior to be 0.29 for D, 0.36 for P, 0.35
 490 for R, and 0.21 for A. The analysis of Yang (2015) produced a different MAP tree, ((H, L), (C,
 491 K)). The difference is due to the use of different priors: Yang (2015) used BPP3.1, with gamma
 492 priors on the parameters (θ s for all populations and τ for the root), whereas here inverse gamma
 493 priors are used in BPP4.3. Note that the species trees have low support in both analyses.

494 In the A00 analysis (parameter estimation under MSC with the MAP species tree fixed),
 495 the posterior means and 95% HPD intervals are shown in table 6a. Strategy P (PHASE) produced

Table 6. Posterior means and 95% HPD CIs for parameters under the MSC model for the east Asian brown frogs and for the chipmunks

	Diploid (D)	PHASE (P)	Random (R)	Ambiguity coding (A)
(a) East Asian brown frogs (Fig. 3a)				
θ_K	4.94 (2.62, 7.65)	6.35 (3.32, 9.80)	6.81 (3.64, 10.48)	2.78 (1.06, 4.99)
θ_C	20.57 (11.78, 30.66)	22.84 (12.98, 34.29)	32.00 (18.11, 48.37)	5.65 (2.41, 9.71)
θ_L	10.82 (6.32, 15.94)	10.12 (5.98, 14.76)	12.33 (7.45, 17.79)	6.73 (2.37, 12.29)
θ_H	3.73 (1.42, 6.73)	3.35 (1.29, 5.96)	5.39 (1.83, 10.16)	1.18 (0.29, 2.51)
θ_5	5.13 (2.20, 8.43)	5.49 (2.54, 8.82)	4.41 (1.50, 7.49)	4.56 (1.74, 7.84)
θ_6	2.21 (0.21, 6.53)	2.00 (0.20, 5.76)	2.65 (0.22, 7.87)	1.85 (0.21, 5.32)
θ_7	1.72 (0.20, 4.61)	1.43 (0.21, 3.62)	1.54 (0.23, 3.94)	1.28 (0.19, 3.15)
τ_5	2.14 (1.57, 2.73)	2.00 (1.50, 2.53)	2.49 (1.89, 3.17)	1.37 (0.86, 1.93)
τ_6	2.03 (1.53, 2.54)	1.91 (1.45, 2.38)	2.30 (1.79, 2.83)	1.23 (0.78, 1.70)
τ_7	1.85 (1.28, 2.44)	1.77 (1.27, 2.27)	2.15 (1.60, 2.72)	1.11 (0.65, 1.61)
(b) Rocky Mountains chipmunks (Fig. 3b)				
θ_Q	0.81 (0.70, 0.93)	0.83 (0.72, 0.94)	0.93 (0.81, 1.05)	0.39 (0.31, 0.47)
θ_I	0.78 (0.67, 0.89)	0.81 (0.69, 0.91)	0.94 (0.81, 1.07)	0.26 (0.21, 0.32)
θ_R	0.36 (0.30, 0.41)	0.36 (0.30, 0.41)	0.37 (0.31, 0.42)	0.32 (0.25, 0.39)
θ_C	0.47 (0.38, 0.55)	0.47 (0.39, 0.55)	0.50 (0.42, 0.58)	0.48 (0.39, 0.56)
θ_D	1.79 (1.61, 1.98)	1.79 (1.60, 1.97)	2.05 (1.84, 2.26)	0.67 (0.57, 0.77)
θ_U	1.04 (0.93, 1.15)	1.04 (0.93, 1.15)	1.06 (0.95, 1.17)	0.83 (0.73, 0.94)
θ_S	0.79 (0.67, 0.90)	0.79 (0.67, 0.90)	0.84 (0.71, 0.96)	0.34 (0.25, 0.43)
θ_8	9.94 (8.31, 11.54)	10.03 (8.61, 11.45)	9.95 (8.32, 11.55)	10.05 (8.62, 11.43)
θ_9	1.24 (1.02, 1.47)	1.24 (1.02, 1.45)	1.24 (1.01, 1.46)	1.24 (1.01, 1.46)
θ_{10}	1.01 (0.65, 1.39)	1.06 (0.68, 1.44)	0.99 (0.64, 1.34)	1.06 (0.63, 1.50)
θ_{11}	4.33 (0.33, 9.43)	5.13 (0.77, 10.38)	2.87 (0.35, 5.87)	2.08 (0.20, 5.90)
θ_{12}	2.43 (0.50, 4.62)	1.84 (0.34, 3.68)	2.16 (0.57, 3.74)	2.45 (0.69, 4.38)
θ_{13}	0.51 (0.21, 0.86)	0.54 (0.19, 0.91)	0.49 (0.21, 0.80)	0.90 (0.34, 1.54)
τ_8	3.83 (3.30, 4.50)	3.80 (3.30, 4.24)	3.85 (3.30, 4.48)	3.70 (3.19, 4.23)
τ_9	1.04 (0.95, 1.14)	1.04 (0.95, 1.13)	1.04 (0.95, 1.13)	0.92 (0.82, 1.01)
τ_{10}	0.74 (0.67, 0.81)	0.72 (0.65, 0.79)	0.75 (0.68, 0.81)	0.59 (0.52, 0.67)
τ_{11}	0.58 (0.45, 0.71)	0.52 (0.42, 0.63)	0.60 (0.49, 0.71)	0.55 (0.44, 0.64)
τ_{12}	0.41 (0.35, 0.46)	0.41 (0.35, 0.46)	0.43 (0.38, 0.49)	0.34 (0.26, 0.43)
τ_{13}	0.33 (0.29, 0.38)	0.34 (0.29, 0.37)	0.37 (0.33, 0.41)	0.21 (0.16, 0.26)

Note.— All values are multiplied by 1000.

496 similar results to strategy D. Strategy R (random) produced overestimates of θ s for modern
 497 species, while strategy A (ambiguity) produced serious underestimates of θ s for modern species
 498 and divergence times. The results are consistent with our findings from the simulation.

499 **Rocky Mountains chipmunks.** In the A01 analysis (species tree inference) of the 500
 500 nuclear loci for Rocky Mountains chipmunks, strategies D, P, and R produced the same MAP
 501 tree, shown in Figure 3b, with the posterior for every node ~ 1.0 . This is also the species tree
 502 inferred by Sarver *et al.* (2021) using summary methods, although the authors obtained lower
 503 support values even with all 1060 loci used. The difference may be due to the higher power of the
 504 BPP analysis, which uses the full data rather than data summaries (e.g., Shi and Yang 2018; Kim
 505 and Degnan, 2020; Zhu and Yang, 2021). Strategy A (ambiguity) produced a different MAP
 506 species tree from the other strategies (Fig. 3b), with the relationship (C, (D, (IQR))) instead of
 507 (D, (C, (IQR))), with the posterior at 0.94. The running time for the A01 analysis, using eight
 508 cores on a server with Intel Xeon Gold 6154 3.0GHz processors, was 9 hours for strategy A, and
 509 16-17 hours for strategies D, P, and R, with strategy D having slightly longer running time. The
 510 number of site patterns at the 1060 loci for strategy D is shown in Fig. S26. Strategy P also
 511 needed the additional time for running the PHASE program, which was 33 mins to phase all 1060

512 loci using one thread on the server.

513 In the A00 analysis (parameter estimation), strategy P (PHASE) produced nearly identical
514 results to strategy D (diploid) (table 6). Compared with strategy D, strategy R (random) produced
515 overestimates of θ s for modern species, while divergence times for recent nodes were also
516 over-estimated very slightly. Strategy A (ambiguity) produced serious underestimates of θ s for
517 modern species, with divergence times, especially of recent nodes, to be underestimated as well.
518 Those results mimic our findings about the relative performance of the different strategies in the
519 simulated data. Running time for the A00 analysis was 2.5 hours for strategy A, and 5-6 hours for
520 strategies D, P, and R. Note that in the A00 analysis the chain is only half as long as in the A01
521 analysis.

522 4. DISCUSSION

523 *The Impact of Phasing Errors Depends on the Inference Problem*

524 We have used simulation to examine the performance of four different strategies for handling
525 heterozygote phase in genomic sequence data: F (full phased data), D (diploid analytical phase
526 integration), P (PHASE), and R (random). Inference problems examined have included species
527 tree estimation under the MSC model and parameter estimation under the MSC and MSci
528 models. We found that the different strategies, including random phase resolution (or equivalently
529 the use of haploid consensus sequences), did not affect species tree estimation when the species
530 divergences are much older than the coalescent times. The different phasing strategies may be
531 expected to have even less impact on inference of deep phylogenies, where within-species
532 polymorphism is much lower than between-species divergence. However, species tree estimation
533 is affected by phasing errors if the species tree is shallow and between-species divergence is
534 similar to within-species polymorphism, if the mutation rate is high so that there are many
535 heterozygote sites in the sequence, and if many sequences are sampled from each species.
536 Phasing errors are clearly important when genomic data are used to infer the divergence history
537 of populations of the same species.

538 We found that estimation of parameters in the MSC and MSci models is more sensitive to
539 phasing errors than is species tree estimation. In particular, population sizes for modern species
540 are seriously overestimated under the MSC and MSci models when random phasing or haploid
541 consensus sequences are used. Our analysis of the simple case of estimating θ under the
542 single-population coalescent suggests that the bias is caused mainly by the unusual sequences
543 generated by random phase resolution (Fig. 6 and table 3). Estimates of the introgression
544 probability and introgression time under the MSci model may also be biased by errors in random
545 phasing. The biases are more serious when the mutation rate is high so that there are multiple
546 heterozygote sites at each locus and when multiple sequences are sampled per species. Those
547 results are consistent with Gronau *et al.* (2011), who also found that random phase resolution
548 affected parameter estimation in their analysis of genomic sequence data from different human
549 populations.

550 *Limitations of our Simulation and Implications to Practical Data Analysis*

551 Here we note a few limitations of our study. First we have examined only one inference method,
552 the Bayesian method implemented in the BPP program. Our results may be expected to apply to
553 other full likelihood implementations such as STARBEAST (Ogilvie *et al.*, 2017; Zhang *et al.*,
554 2018) or PHYLONET-SEQ (Wen and Nakhleh, 2018), but may not apply to summary methods.
555 Similarly we considered only a few inference problems under the MSC and MSci models using
556 genomic sequence data. We have not examined the impact of phasing errors on inference of
557 population demographic changes or on inference of migration/introgression histories (our
558 simulation under the MSci model assumed a fixed introgression event).

559 Given those caveats, we discuss the implications of our simulation results to practical data
560 analysis. First, our simulation as well as those of Gronau *et al.* (2011) and Andermann *et al.*
561 (2019) suggest that random phase resolution or the use of haploid consensus sequences should be
562 avoided. Strategy R never performed better than computational phasing (strategy P) in our
563 simulations. Similarly strategy A (ambiguity) should not be recommended. Virtually all
564 phylogenetic likelihood programs accommodate ambiguities in a sequence alignment
565 representing undetermined nucleotides using a data augmentation algorithm in the likelihood
566 calculation (Felsenstein, 2004, pp.255–6; Yang, 2014, pp.110-112). As heterozygotes (with, e.g.,
567 Y meaning both T and C) are not ambiguities (with Y meaning either T or C), this approach
568 misinterprets the data, and has the obvious effect of underestimating the heterozygosity or θ for
569 the modern species. Bias may also be introduced into estimates of other parameters, such as
570 underestimation of divergence times (Andermann *et al.*, 2019). The approach also underestimates
571 the information content in the data, as it in effect treats two sequences (although unphased) as
572 only one. This mistake in the treatment of the data was made by Rannala and Yang (2003) in the
573 analysis of three human noncoding loci of Zhao *et al.* (2000), Yu *et al.* (2001), and Makova *et al.*
574 (2001), and by Yang (2015) in the analysis of the five nuclear loci from East Asian brown frogs
575 (Zhou *et al.*, 2012). The mistake is easy to see from the occurrence of the same ambiguity
576 character (such as Y) in multiple sequences at the same site in the alignment.

577 Strategy D (diploid analytical integration) produced results that are extremely similar to
578 the use of the full data (F) in all simulation settings of this study (see also Gronau *et al.*, 2011).
579 As the algorithm averages over all possible phase resolutions and constitutes a full likelihood
580 approach to handling missing data, it is the optimal statistical approach when the data consist of
581 unphased diploid sequences, and may thus be recommended in general, even for inference
582 problems that are not examined in our simulation study. As a statistical inference method,
583 strategy D is equivalent to the approach of sampling phase resolutions in a Markov chain Monte
584 Carlo (MCMC) algorithm (Kuhner and Felsenstein, 2000). In small or intermediate datasets,
585 analytical phase integration appears more efficient computationally than MCMC, whereas for
586 large datasets, both may be unfeasible.

587 Note that analyses under the four strategies F, D, P, and R involve the same number of
588 species, the same number of parameters, the same number of loci, the same number of sequences,
589 etc., with the only difference being in the number of site patterns. The relative computational load
590 for the strategies is thus proportional to the number of site patterns. Strategy D performs
591 phylogenetic likelihood calculation (Felsenstein, 1981) for all distinct site patterns that may result
592 from enumerating all possible phase resolutions, which include the true phase resolution. Thus
593 strategy D involves at least as many site patterns as in the full data (strategy F). For the

594 simulations of this study, the number of site patterns for strategy D is less than twice the number
595 for strategy F (Fig. S12). However, if there are many long sequences of high heterozygosity at a
596 locus, enumeration of all phase resolutions may lead to a huge number of site patterns. For
597 example, the three noncoding regions of human DNA analyzed by Rannala and Yang (2003) have
598 about 60 sequences per locus, with $\sim 10^4$ sites. The number of site patterns in the unphased
599 alignments (strategy A) is 50–73, but reaches 1.2–4.4 million for strategy D, rendering the
600 analysis unfeasible. Note that those loci are long genomic segments, which may be affected by
601 recombination, whereas datasets suitable for analysis under the MSC typically involve much
602 shorter genomic segments (e.g., Burgess and Yang, 2008).

603 We suggest that computational phasing (strategy P) should be an acceptable alternative
604 when strategy D is computationally unfeasible. In our analyses of the simulated and real datasets,
605 strategy P produced similar results to the use of full data (F) or the analytical phase integration
606 approach (D), with very small biases. Note that the Bayesian program PHASE assumes a
607 population genetics model and is designed for sequence or allelic data from the same species.
608 However, our use of it to analyze sequence data from multiple species produced relatively small
609 biases in parameter estimation in both simulated data and in the two real datasets, much better
610 than random phase resolution or haploid consensus sequences. We also note that phasing based
611 on reads combined with bioinformatic analysis shows great promise (Andermann *et al.*, 2019). In
612 particular, exciting developments in sequencing technology to provide longer reads, combined
613 with computational algorithms (Porubsky *et al.*, 2020; Zhou *et al.*, 2020; Cheng *et al.*, 2021),
614 may soon make it practical to produce routinely fully phased diploid genomes.

615 5. SUPPLEMENTARY MATERIAL

616 Data available from the Dryad Digital Repository: <https://doi.org/10.5061/dryad.vmcvdcnrd>.

617 6. ACKNOWLEDGMENTS

618 We thank two anonymous reviewers, the AE (Laura Kubatko) and the EiC (Bryan Carstens) for
619 constructive comments. We thank Jiayi Ji for converting the chipmunk alignments of Sarver *et al.*
620 (2021) from the Nexus format to the BPP format.

621 7. FUNDING

622 This study has been supported by a Biotechnology and Biological Sciences Research Council
623 grant (BB/P006493/1) to Z.Y. and a BBSRC equipment grant (BB/R01356X/1). A.D.L. is
624 supported by a National Science Foundation grant (NSF-SBS-2023723). J.H.'s visit to London is
625 supported by China Scholarship Council (CSC).

626 8. SUPPLEMENTARY MATERIAL

627 Data available from the Dryad Digital Repository: <http://dx.doi.org/10.5061/dryad.xxxxxx>.

628

REFERENCES

- 629 Andermann, T., Fernandes, A. M., Olsson, U., Topel, M., Pfeil, B., Oxelman, B., Aleixo, A.,
630 Faircloth, B. C., and Antonelli, A. 2019. Allele phasing greatly improves the phylogenetic
631 utility of ultraconserved elements. *Syst. Biol.*, 68(1): 32–46.
- 632 Andrews, K. R., Good, J. M., Miller, M. R., Luikart, G., and Hohenlohe, P. A. 2016. Harnessing
633 the power of RADseq for ecological and evolutionary genomics. *Nat. Rev. Genet.*, 17(2):
634 81–92.
- 635 Burgess, R. and Yang, Z. 2008. Estimation of hominoid ancestral population sizes under
636 Bayesian coalescent models incorporating mutation rate variation and sequencing errors. *Mol.*
637 *Biol. Evol.*, 25: 1979–1994.
- 638 Cheng, H., Concepcion, G. T., Feng, X., Zhang, H., and Li, H. 2021. Haplotype-resolved de novo
639 assembly using phased assembly graphs with hifiasm. *Nat. Methods*.
- 640 Choi, Y., Chan, A. P., Kirkness, E., Telenti, A., and Schork, N. J. 2018. Comparison of phasing
641 strategies for whole human genomes. *PLoS Genet.*, 14(4): e1007308.
- 642 Eaton, D. A. R., Spriggs, E. L., Park, B., and Donoghue, M. J. 2017. Misconceptions on missing
643 data in RAD-seq phylogenetics with a deep-scale example from flowering plants. *Syst. Biol.*,
644 66(3): 399–412.
- 645 Edwards, S., Cloutier, A., and Baker, A. 2017. Conserved nonexonic elements: a novel class of
646 marker for phylogenomics. *Syst. Biol.*, 66(6): 1028–1044.
- 647 Faircloth, B. C., McCormack, J. E., Crawford, N. G., Harvey, M. G., Brumfield, R. T., and Glenn,
648 T. C. 2012. Ultraconserved elements anchor thousands of genetic markers spanning multiple
649 evolutionary timescales. *Syst. Biol.*, 61(5): 717–726.
- 650 Felsenstein, J. 1981. Evolutionary trees from dna sequences: a maximum likelihood approach. *J.*
651 *Mol. Evol.*, 17: 368–376.
- 652 Felsenstein, J. 1992. Estimating effective population size from samples of sequences: inefficiency
653 of pairwise and segregating sites as compared to phylogenetic estimates. *Genet. Res.*, 59:
654 139–147.
- 655 Felsenstein, J. 2004. *Inferring Phylogenies*. Sinauer Associates, Sunderland, Massachusetts.
- 656 Felsenstein, J. 2019. *Theoretical Evolutionary Genetics*.
657 <https://evolution.genetics.washington.edu/pgbook/pgbook.html>.
- 658 Flot, J. F. 2010. SeqPhase: a web tool for interconverting phase input/output files and FASTA
659 sequence alignments. *Mol. Ecol. Resour.*, 10(1): 162–166.
- 660 Flouri, T., Jiao, X., Rannala, B., and Yang, Z. 2018. Species tree inference with BPP using
661 genomic sequences and the multispecies coalescent. *Mol. Biol. Evol.*, 35(10): 2585–2593.
- 662 Flouri, T., Jiao, X., Rannala, B., and Yang, Z. 2020a. A Bayesian implementation of the
663 multispecies coalescent model with introgression for phylogenomic analysis. *Mol. Biol. Evol.*,
664 37(4): 1211–1223.

- 665 Flouri, T., Rannala, B., and Yang, Z. 2020b. A tutorial on the use of bpp for species tree
666 estimation and species delimitation. In C. Scornavacca, F. Delsuc, and N. Galtier, editors,
667 *Phylogenetics in the Genomic Era*, book section 5.6, pages 5.6.1–16. No Commercial
668 Publisher.
- 669 Grabherr, M. G., Haas, B. J., Yassour, M., Levin, J. Z., Thompson, D. A., Amit, I., Adiconis, X.,
670 Fan, L., Raychowdhury, R., Zeng, Q., Chen, Z., Mauceli, E., Hacohen, N., Gnirke, A., Rhind,
671 N., di Palma, F., Birren, B. W., Nusbaum, C., Lindblad-Toh, K., Friedman, N., and Regev, A.
672 2011. Full-length transcriptome assembly from rna-seq data without a reference genome. *Nat.*
673 *Biotechnol.*, 29(7): 644–652.
- 674 Gronau, I., Hubisz, M. J., Gulko, B., Danko, C. G., and Siepel, A. 2011. Bayesian inference of
675 ancient human demography from individual genome sequences. *Nature Genet.*, 43:
676 1031–1034.
- 677 Huang, J., Flouri, T., and Yang, Z. 2020. A simulation study to examine the information content
678 in phylogenomic datasets under the multispecies coalescent model. *Mol. Biol. Evol.*
- 679 Huelsenbeck, J. and Rannala, B. 2004. Frequentist properties of bayesian posterior probabilities
680 of phylogenetic trees under simple and complex substitution models. *Syst. Biol.*, 53: 904–913.
- 681 Jukes, T. and Cantor, C. 1969. Evolution of protein molecules. In H. Munro, editor, *Mammalian*
682 *Protein Metabolism*, pages 21–123. Academic Press, New York.
- 683 Karin, B. R., Gamble, T., and Jackman, T. R. 2020. Optimizing phylogenomics with rapidly
684 evolving long exons: Comparison with anchored hybrid enrichment and ultraconserved
685 elements. *Mol. Biol. Evol.*, 37(3): 904–922.
- 686 Kim, A. and Degnan, J. 2020. Pranc: MI species tree estimation from the ranked gene trees under
687 coalescence. *Bioinformatics*.
- 688 Kuhner, M. K. and Felsenstein, J. 2000. Sampling among haplotype resolutions in a
689 coalescent-based genealogy sampler. *Genet. Epidemiol.*, 19(Suppl): S15–21.
- 690 Leaché, A. D. and Oaks, J. R. 2017. The utility of single nucleotide polymorphism (SNP) data in
691 phylogenetics. *Ann. Rev. Ecol. Evol. Syst.*, 48: 69–84.
- 692 Leaché, A. D., Zhu, T., Rannala, B., and Yang, Z. 2019. The spectre of too many species. *Syst.*
693 *Biol.*, 68(1): 168–181.
- 694 Lemmon, A. R., Emme, S. A., and Lemmon, E. M. 2012. Anchored hybrid enrichment for
695 massively high-throughput phylogenomics. *Syst. Biol.*, 61(5): 727–744.
- 696 Makova, K. D., Ramsay, M., Jenkins, T., and Li, W. H. 2001. Human dna sequence variation in a
697 6.6-kb region containing the melanocortin 1 receptor promoter. *Genetics*, 158: 1253–1268.
- 698 Niu, T., Qin, Z. S., Xu, X., and Liu, J. S. 2002. Bayesian haplotype inference for multiple linked
699 single-nucleotide polymorphisms. *Am. J. Hum. Genet.*, 70(1): 157–169.
- 700 Ogilvie, H. A., Bouckaert, R. R., and Drummond, A. J. 2017. Starbeast2 brings faster species
701 tree inference and accurate estimates of substitution rates. *Mol. Biol. Evol.*, 34(8): 2101–2114.

- 702 Porubsky, D., Ebert, P., Audano, P. A., Vollger, M. R., Harvey, W. T., Marijon, P., Ebler, J.,
703 Munson, K. M., Sorensen, M., Sulovari, A., Haukness, M., Ghareghani, M., Human Genome
704 Structural Variation, C., Lansdorp, P. M., Paten, B., Devine, S. E., Sanders, A. D., Lee, C.,
705 Chaisson, M. J. P., Korb, J. O., Eichler, E. E., and Marschall, T. 2020. Fully phased human
706 genome assembly without parental data using single-cell strand sequencing and long reads.
707 *Nat. Biotechnol.*
- 708 Rannala, B. and Yang, Z. 1996. Probability distribution of molecular evolutionary trees: a new
709 method of phylogenetic inference. *J. Mol. Evol.*, 43: 304–311.
- 710 Rannala, B. and Yang, Z. 2003. Bayes estimation of species divergence times and ancestral
711 population sizes using DNA sequences from multiple loci. *Genetics*, 164: 1645–1656.
- 712 Rannala, B. and Yang, Z. 2017. Efficient Bayesian species tree inference under the multispecies
713 coalescent. *Syst. Biol.*, 66: 823–842.
- 714 Sarver, B. A. J., Herrera, N. D., Sneddon, D., Hunter, S. S., Settles, M. L., Kronenberg, Z.,
715 Demboski, J. R., Good, J. M., and Sullivan, J. 2021. Diversification, introgression, and rampant
716 cytonuclear discordance in rocky mountains chipmunks (sciuridae: Tamias). *Syst. Biol.*
- 717 Scheet, P. and Stephens, M. 2006. A fast and flexible statistical model for large-scale population
718 genotype data: applications to inferring missing genotypes and haplotypic phase. *Am. J. Hum.*
719 *Genet.*, 78(4): 629–644.
- 720 Shi, C. and Yang, Z. 2018. Coalescent-based analyses of genomic sequence data provide a robust
721 resolution of phylogenetic relationships among major groups of gibbons. *Mol. Biol. Evol.*, 35:
722 159–179.
- 723 Simpson, J. T., Wong, K., Jackman, S. D., Schein, J. E., Jones, S. J., and Birol, I. 2009. ABySS: a
724 parallel assembler for short read sequence data. *Genome Res.*, 19(6): 1117–1123.
- 725 Stephens, M. and Donnelly, P. 2003. A comparison of Bayesian methods for haplotype
726 reconstruction from population genotype data. *Am. J. Hum. Genet.*, 73: 1162–1169.
- 727 Stephens, M., Smith, N. J., and Donnelly, P. 2001. A new statistical method for haplotype
728 reconstruction from population data. *Am. J. Hum. Genet.*, 68: 978–989.
- 729 Tajima, F. 1983. Evolutionary relationship of dna sequences in finite populations. *Genetics*, 105:
730 437–460.
- 731 Tewhey, R., Bansal, V., Torkamani, A., Topol, E. J., and Schork, N. J. 2011. The importance of
732 phase information for human genomics. *Nat. Rev. Genet.*, 12: 215–223.
- 733 Watterson, G. 1975. On the number of segregating sites in genetical models without
734 recombination. *Theor. Popul. Biol.*, 7: 256–276.
- 735 Weisenfeld, N. I., Kumar, V., Shah, P., Church, D. M., and Jaffe, D. B. 2017. Direct
736 determination of diploid genome sequences. *Genome Res.*, 27(5): 757–767.
- 737 Wen, D. and Nakhleh, L. 2018. Coestimating reticulate phylogenies and gene trees from
738 multilocus sequence data. *Syst. Biol.*, 67(3): 439–457.

- 739 Yang, Z. 2006. *Computational Molecular Evolution*. Oxford University Press, Oxford, UK.
- 740 Yang, Z. 2014. *Molecular Evolution: A Statistical Approach*. Oxford University Press, Oxford,
741 England.
- 742 Yang, Z. 2015. The BPP program for species tree estimation and species delimitation. *Curr.*
743 *Zool.*, 61(5): 854–865.
- 744 Yang, Z. and Rannala, B. 2005. Branch-length prior influences bayesian posterior probability of
745 phylogeny. *Syst. Biol.*, 54: 455–470.
- 746 Yang, Z. and Rannala, B. 2010. Bayesian species delimitation using multilocus sequence data.
747 *Proc. Natl. Acad. Sci. U.S.A.*, 107: 9264–9269.
- 748 Yang, Z. and Rannala, B. 2014. Unguided species delimitation using DNA sequence data from
749 multiple loci. *Mol. Biol. Evol.*, 31(12): 3125–3135.
- 750 Yang, Z., Kumar, S., and Nei, M. 1995. A new method of inference of ancestral nucleotide and
751 amino acid sequences. *Genetics*, 141: 1641–1650.
- 752 Yu, N., Zhao, Z., Fu, Y. X., Sambuughin, N., Ramsay, M., Jenkins, T., Leskinen, E., Patthy, L.,
753 Jorde, L. B., Kuromori, T., and Li, W. H. 2001. Global patterns of human dna sequence
754 variation in a 10-kb region on chromosome 1. *Mol. Biol. Evol.*, 18: 214–222.
- 755 Zerbino, D. R. and Birney, E. 2008. Velvet: algorithms for de novo short read assembly using de
756 Bruijn graphs. *Genome Res.*, 18(5): 821–829.
- 757 Zhang, C., Ogilvie, H. A., Drummond, A. J., and Stadler, T. 2018. Bayesian inference of species
758 networks from multilocus sequence data. *Mol. Biol. Evol.*, 35: 504–517.
- 759 Zhao, Z., Jin, L., Fu, Y. X., Ramsay, M., Jenkins, T., Leskinen, E., Pamilo, P., Trexler, M., Patthy,
760 L., Jorde, L. B., Ramos-Onsins, S., Yu, N., and Li, W. H. 2000. Worldwide dna sequence
761 variation in a 10-kilobase noncoding region on human chromosome 22. *Proc. Natl. Acad. Sci.*
762 *U.S.A.*, 97: 11354–11358.
- 763 Zhou, Q., Tang, D., Huang, W., Yang, Z., Zhang, Y., Hamilton, J. P., Visser, R. G. F., Bachem, C.
764 W. B., Robin Buell, C., Zhang, Z., Zhang, C., and Huang, S. 2020. Haplotype-resolved
765 genome analyses of a heterozygous diploid potato. *Nat. Genet.*, 52(10): 1018–1023.
- 766 Zhou, W. W., Wen, Y., Fu, J., Xu, Y. B., Jin, J. Q., Ding, L., Min, M. S., Che, J., and Zhang, Y. P.
767 2012. Speciation in the rana chensinensis species complex and its relationship to the uplift of
768 the qinghai-tibetan plateau. *Mol. Ecol.*, 21(4): 960–973.
- 769 Zhu, T. and Yang, Z. 2021. Complexity of the simplest species tree problem. *Mol. Biol. Evol.*