
PRACTICAL SPEEDUP OF BAYESIAN INFERENCE OF SPECIES PHYLOGENIES BY RESTRICTING THE SPACE OF GENE TREES

A PREPRINT

Yaxuan Wang*
Computer Science Department
Rice University
Houston, TX 77005
yaxuan.wang@rice.edu

Huw A. Ogilvie
Computer Science Department
Rice University
Houston, TX 77005
huw.a.ogilvie@rice.edu

Luay Nakhleh
Computer Science Department
Rice University
Houston, TX 77005
nakhleh@rice.edu

September 15, 2019

ABSTRACT

Species tree inference from multi-locus data has emerged as a powerful paradigm in the post-genomic era, both in terms of the accuracy of the species tree it produces as well as in terms of elucidating the processes that shaped the evolutionary history. Bayesian methods for species tree inference are desirable in this area as they have been shown to yield accurate estimates, but also to naturally provide measures of confidence in those estimates. However, the heavy computational requirements of Bayesian inference have limited the applicability of such methods to very small data sets.

In this paper, we show that the computational efficiency of Bayesian inference under the multispecies coalescent can be improved in practice by restricting the space of the gene trees explored during the random walk, without sacrificing accuracy as measured by various metrics. The idea is to first infer constraints on the trees of the individual loci in the form of unresolved gene trees, and then to restrict the sampler to consider only resolutions of the constrained trees. We demonstrate the improvements gained by such an approach on both simulated and biological data.

Keywords species tree · multispecies coalescent · Bayesian MCMC · efficiency

1 Introduction

Species tree inference under the multispecies coalescent (MSC) model has gained much attention due to the fact that it allows for modeling gene tree heterogeneity that arises across the genome from incomplete lineage sorting (ILS; Degnan and Rosenberg, 2009). A wide array of methods that assume or are inspired by the MSC have been devised (Liu et al., 2010; Liu and Yu, 2011; Mirarab et al., 2014; Chifman and Kubatko, 2014; Wang and Nakhleh, 2018), including the Bayesian methods of Ogilvie et al. (2017); Flouri et al. (2018); Zhu et al. (2018). The MSC was recently extended to the multispecies network coalescent (MSNC) to account for reticulation (in addition to ILS; see Yu et al., 2014) and Bayesian methods for inference under this model have been devised (Wen et al., 2016; Wen and Nakhleh, 2017; Wen et al., 2018; Zhu et al., 2018; Zhang et al., 2017; Bouckaert et al., 2019).

*Corresponding author

The power of Bayesian methods lies in their ability to incorporate prior knowledge, infer values of parameters beyond the tree topology, and provide measures of confidence in the inference based on the posterior that they sample (Huelsenbeck et al., 2001). However, for these Bayesian methods to converge onto the true posterior distribution, they demand significant computational resources, an issue that has thus far limited their applicability to only small data sets in terms of both the number of taxa and number of loci in the data set (Ogilvie et al., 2016). This is why ILS-aware methods that have been proven to be statistically consistent under the MSC and, at the same time very efficient computationally, are used to infer large-scale species trees (Mirarab et al., 2014; Liu et al., 2010; Liu and Yu, 2011; Chifman and Kubatko, 2014). However, these methods focus almost exclusively on the species tree topology and provide neither accurate information on other parameters, such as divergence times and population sizes, nor confidence intervals for their inferences. The question we address in this paper is: Can the convergence of Bayesian methods be improved in practice without sacrificing the accuracy of the information they provide?

A rich body of literature exists on the development of methods for statistical inference outside phylogenetics, much of which has been adopted by Bayesian phylogenetic methods. The ubiquitous Markov chain Monte Carlo (MCMC) arose from nuclear weapons research (Robert and Casella, 2011), and is the basis for tree and network inference in MrBayes, *BEAST, PhyloNet and other software tools. Metropolis coupling to accelerate MCMC was developed for the inference of spatial statistics (Geyer, 1991), and has been implemented by the above listed phylogenetics software (Ronquist et al., 2012; Bouckaert et al., 2019; Wen et al., 2018). Variational Bayes is a radically different approach which fits parametric distributions to model parameters, unlike MCMC which is non-parametric. Variational Bayes was originally developed for graphical models (Attias, 1999), and has recently been applied to compute posterior distributions and marginal likelihoods of phylogenetic trees (Zhang and Matsen, 2019; Fourment and Darling, 2019).

All of these methods were developed decades before their adoption for phylogenetic inference because tree and network space is far more complex than the typical multidimensional parameter space. The number of unrooted or rooted trees grows superexponentially with the number of taxa (Felsenstein, 1978), and is for all practical purposes infinite when the number of taxa is large. Multilocus MSC inference embeds gene trees within a species tree, with the constraint that between-species coalescent events must take place earlier in time than the most recent common ancestor (MRCA) time of the involved species. This multiplies the complexity of the inference problem by increasing the number of trees to infer, and because the probability distributions of node heights for different trees are not independent.

Rather than trying to adapt an algorithm developed for other fields of natural sciences or mathematics, we have developed a heuristic method that specifically applies to the problem of multilocus MSC inference. The heuristic method constrains the space of gene tree topologies to allow for faster convergence and, consequently, analyses of larger data sets. The idea behind our approach is very simple: A set of constraints in the form of a tree which is usually less than fully resolved is estimated independently for each individual locus, and then MCMC walks in the portion of the tree space that is consistent with these constraints. In other words, the MCMC sampler considers only gene trees that are consistent with the constraints on the individual loci. Using simulated data under a variety of conditions and employing several metrics for assessing performance, we demonstrate that this simple approach results in computational improvements relative to unconstrained Bayesian MCMC without sacrificing accuracy. We then analyze a biological data set and show that the new approach enables analyses that had before necessitated dividing the data set into smaller ones.

Our work presents an approach for improving the computational requirements of Bayesian inference of species phylogenies. The constraints on the individual loci can be obtained in various ways and the proposals that satisfy these constraints can be derived in multiple ways as well. In this work, we implemented one specific method for obtaining the constraints and a standard set of proposals that satisfy them. As this approach can be adopted by any Bayesian species phylogeny inference method, both of these components can be further modified to achieve even further improvement to the computational requirements of Bayesian inference.

Annealed sequential Monte Carlo (SMC) for phylogenetics is another algorithm recently introduced to accelerate Bayesian phylogenetics inference (Wang et al., 2019). Like Metropolis coupling and our approach, it has the advantage of being implementable using existing Metropolis-Hastings moves developed for trees and other parameters. In the future, our approach could be combined with annealed SMC, for even greater improvements in computational performance.

2 New Approach

A Bayesian formulation of the multi-locus species tree inference problem is to estimate the posterior distribution over species trees,

$$f(\Psi) \propto p(\Psi)P(\mathcal{S}|\Psi) = p(\Psi) \int_g P(\mathcal{S}|g)P(g|\Psi)dg. \quad (1)$$

Under this formulation, the data consist of a set of m sequence alignments $\mathcal{S} = \{S_1, S_2, \dots, S_m\}$ obtained from unlinked loci in the genomes of n species. $p(\Psi)$ is the prior on the species tree (topology and all associated parameters, for the sake of brevity), the integration is taken over all possible gene trees for the m loci, and $P(g|\Psi)$ is derived based on the MSC (Degnan and Rosenberg, 2009). Furthermore, under the common assumption of recombination-free, unlinked loci, the posterior becomes proportional to

$$p(\Psi) \int_G \prod_{i=1}^m P(S_i|g_i)P(g_i|\Psi)dG, \quad (2)$$

where each g_i ranges over all possible gene tree topologies and their associated parameters, G .

As the integration in Equation (2) cannot be derived analytically, Markov chain Monte Carlo (MCMC) sampling algorithms are often employed to obtain samples from the posterior distribution and approximate it based on those samples. Beyond a handful of taxa and a handful of loci, MCMC algorithms to approximate Equation (2) become computationally infeasible, getting stuck in the peaks and troughs of the posterior distribution and requiring extremely large numbers of iterations to converge.

Our approach to tackle the computational challenge works as follows. For each sequence alignment S_i , maximum likelihood with bootstrapping is run to obtain a set of gene trees from which a majority-rule consensus tree with a pre-specified support threshold x is built. For example, for $x = 90$, a majority-rule consensus tree is built where only clades that appear in at least 90% of the bootstrap trees are included. We denote this majority-rule consensus tree by C_i (we use the value of x explicitly in the naming only when it is not clear from the context) and call it a constraint gene tree, or CGT. Our approach now samples according to equation (2) with one difference: The integration is taken over g_i 's that are consistent with, i.e., refinements of, their respective C_i constraints. For example, if $C_i = (((A, B), C), (D, E, F))$, the sampler considers gene tree $g_i = (((A, B), C), ((D, E), F))$ as it is a refinement of C_i , but does not consider gene tree $g_i = (((A, B), (C, D)), (E, F))$ as it is not a refinement of C_i . This concept is illustrated in Figure 1.

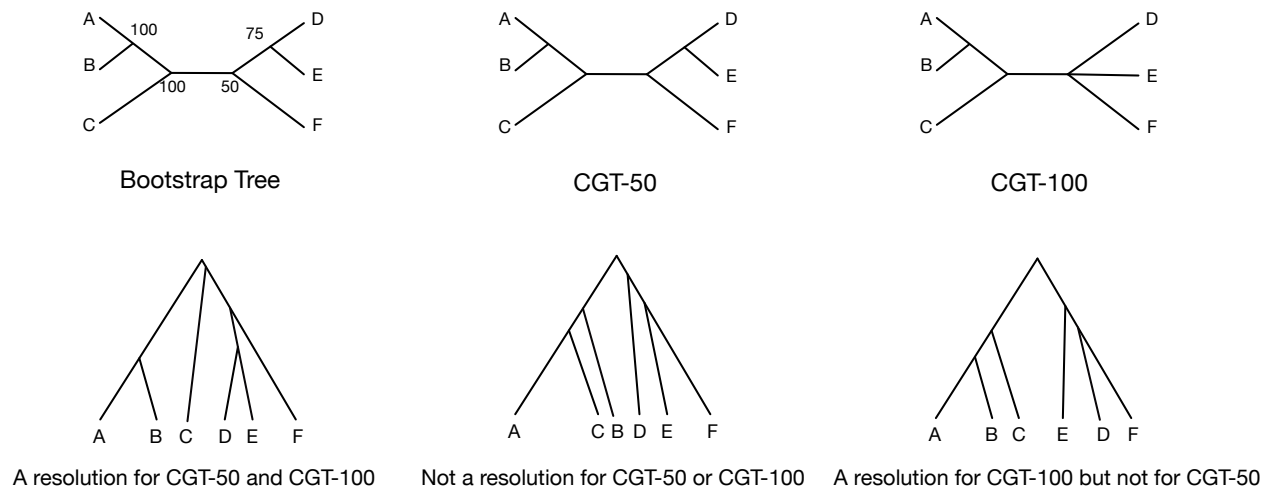


Figure 1: **Constraint trees, their resolutions, and acceptable moves.** Bootstrap tree and CGTs under the consensus threshold of 50 (CGT-50) and 100 (CGT-100) are shown in the first row. In the second row, three possible proposed gene trees are provided. The left tree is acceptable given the constraints CGT-50 and CGT-100. The middle tree is not acceptable given the constraints CGT-50 or CGT-100. The right tree is acceptable given the constraint CGT-100 but not CGT-50.

Observe that if C_i is a star phylogeny (a tree that has no internal branches), then the method is effectively sampling according to Equation (2), whereas if C_i is fully resolved (a binary tree), then the sampler fixes the gene tree topology for locus i and only samples its parameters. It is important to note that C_i imposes only topological constraints; that is, C_i has no branch lengths. Furthermore, we take C_i to be unrooted, so that the sampler is allowed to sample the roots of the gene trees.

The posterior of species trees includes a “long tail” region of model trees where the likelihood of each tree is very small but the number of trees within such region is very large. Unfortunately, as the scale of the data set increases, the autocorrelation time of the MCMC chain increases dramatically (Ogilvie et al., 2016). The motivation of our approach is that by constraining the gene trees, the sampler can avoid the long tail and have better mixing. While utilizing these

constraints necessarily means that the sampler is not sampling from the same posterior distribution as an unconstrained version of the sampler (except for the case where the constraints are star phylogenies), we demonstrate below that this has very little impact on the accuracy of the sampler in practice.

Hereafter we write CMC_{CMC} to denote constrained MCMC according to our new approach and UMC_{CMC} to denote the unconstrained version of MCMC. We also write CMC_{CMC}- x , where x is a value between 50 and 100, to denote the use of CMC_{CMC} with support threshold x in the majority-rule consensus tree, or a value of 0 to denote use with the maximum likelihood tree.

3 Results and Discussion

A simulation study was carried out to comprehensively analyze the performance of CMC_{CMC} and UMC_{CMC}. We varied the simulation data set along three dimensions: evolutionary scenarios of species, the number of loci and the number of taxa. When we focused on one dimension, the other two dimensions were fixed. More details are provided in the Materials and Methods section. To simulate different scenarios of complexity and signal in the data, we varied evolutionary time scales and population sizes in four categories:

- “OH”: old divergence times and high population size;
- “OL”: old divergence times and low population size;
- “YH”: young divergence times and high population size; and,
- “YL”: young divergence times and low population size.

To further examine the performance of each method, we varied the number of loci (10, 20, 40) for the YH condition while fixing the number of taxa as 16. We also varied the numbers of taxa (16, 32, 48) and fixed the number of loci as 10. Unless otherwise stated, there are 10 replicates for each condition.

3.1 CGTs improve the convergence of MCMC

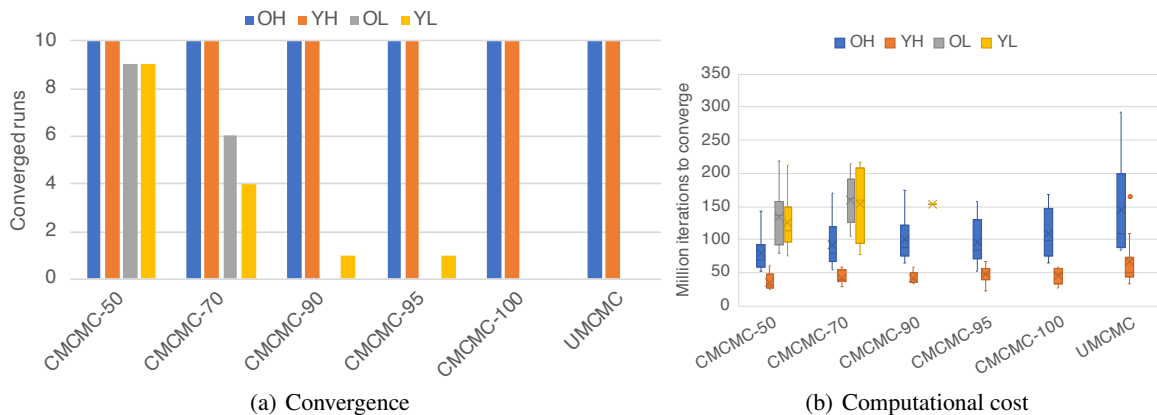


Figure 2: **Convergence and efficiency of CMC and UMC.** (a) Convergence of samplers with or without constraint gene tree. Different samplers are shown on the x axis and the y axis shows the number of data sets (out of 10) on which the sampler converged within 72 hours. (b) The computational costs of the various samplers on the four evolutionary scenarios. The computational cost is quantified by the number of iteration required for convergence.

The ability to converge within a reasonable time is a key metric to evaluate the performance of an MCMC sampler. An effective sample size (ESS) of at least 200 is used as a threshold for convergence in the popular MCMC diagnostic and analysis tool Tracer (Rambaut et al., 2018). In this work, we target the same convergence standard for continuous parameters including the posterior probability, likelihood, prior probability, coalescent likelihood, tree height and population size. We terminated any chain still running after 72 hours.

Figure 2(a) shows that decreasing the consensus threshold enables convergence for low population size conditions, which is impossible for UMC_{CMC} within 72 hours. We also show the improved convergence of CMC_{CMC} as the number of loci increases in Figure S1 and as the number of taxa increases in Figure S2.

When MCMC chains are able to converge, CMCMC reduces the number of iterations required for convergence into less than half of the UCMCMC under various evolutionary parameters as shown in Figure 2(b). Furthermore, CMCMC took fewer iterations than UCMCMC to converge for different numbers of loci and different numbers of taxa (Figure S3 and Figure S4).

3.2 CMCMC and UCMCMC derive similar posterior distributions

The ultimate goal of phylogenetic inference problem with Bayesian sampling is to approach the posterior distribution of the true species tree. One way to verify the posterior distribution is to compare the average standard deviation of split frequencies (ASDSF; Lakner et al., 2008) of the 95% credible set. Note that the 95% credible set or interval of a well calibrated Bayesian method derives the true value in 95% of cases. The smaller the ASDSF is, the more similar the species tree distributions are. A threshold of 0.01 on the ASDSF is commonly used to assess the convergence of two chains. An ASDSF value below 0.01 is taken to indicate that the chains are likely to be sampling from the same underlying distribution (for examples, see Stunžėnas et al., 2011; Mazza et al., 2016; Stensvold et al., 2011).

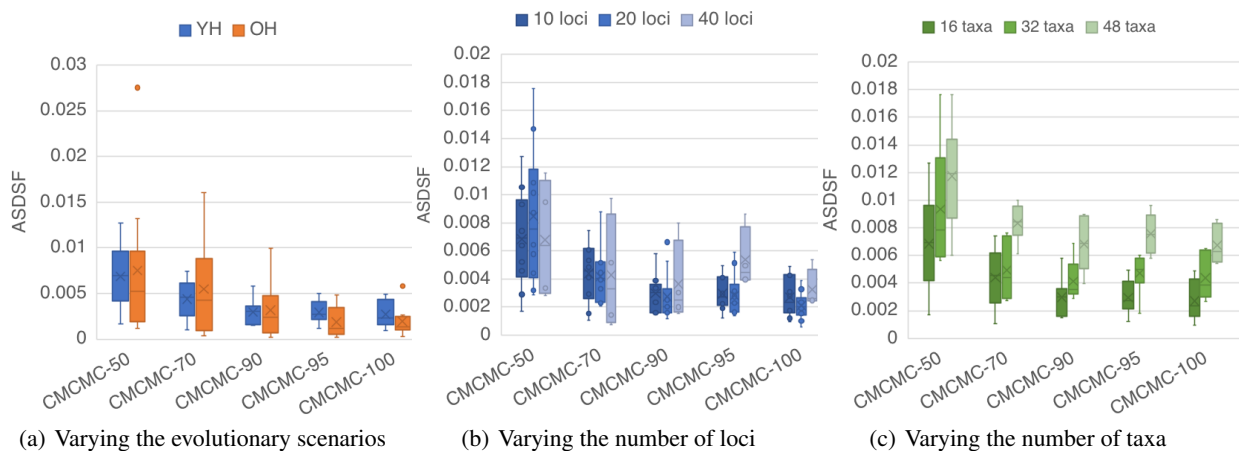


Figure 3: ASDSF between the CMCMC and UCMCMC chains. The x axis lists CMCMC samplers with different support thresholds and the y axis shows the ASDSF between each CMCMC method and UCMCMC. (a) ASDSF values when varying the divergence times fixing the number of taxa and loci as 16 and 10. (b) ASDSF values when varying the number of loci. There are 10 replicates for 10- and 20-locus data sets while for 40-locus data sets there are only 4 replicates because UCMCMC can converge on only 4 out of 10 replicates in this experiment. (c) ASDSF values when varying the number of taxa. There are 10 replicates for 16- and 32-taxon data sets while for 48-taxon data sets there are only 6 replicates because UCMCMC can converge on only 6 out of 10 replicates in this experiment.

Figure 3 shows the ASDSF between the CMCMC and UCMCMC chains for different evolutionary scenarios, different numbers of loci, and different numbers of taxa. For all simulated scenarios in Figure 3(a), the ASDSF values of CMCMC and UCMCMC in most replicates are below 0.01. As we increased the number of loci and taxa, all ASDSF interquartile ranges (IQRs) fell below 0.01, except for CMCMC-50 for which the top of the range rose slightly above higher than 0.01 when the number of loci or the number of taxa increased, as shown in Figure 3(b) and Figure 3(c). In summary, while the underlying distributions sampled by CMCMC and UCMCMC are different by construction, our results show that in practice they are almost the same.

If an internal node in one constraint gene tree (CGT) is binary, we consider such node as resolved because there is only one topology for this CGT. For young divergence time scenarios, there are fewer substitutions and hence less information available to reconstruct the phylogeny because the substitution mutation rate is low if we fix the population times. For a given threshold, CGTs in young divergence time scenarios were less resolved than CGTs in old divergence time scenarios in Figure S5. But for all conditions the proportion of resolved nodes steadily decreased as the threshold was raised. A similar trend was observed when increasing the number of taxa in Figure S6. This shows the role that the support threshold plays as a useful tuning parameter for our heuristic.

While the ASDSF provides a numeric measure reflecting the similarity between the distributions being sampled, we visualize in Figure 4 the distributions sampled by the various samplers as an illustration of this similarity. While decreasing the support threshold increases the difference in the posterior distribution, all CMCMC methods derive posterior distributions similar to UCMCMC, except for CMCMC-0.

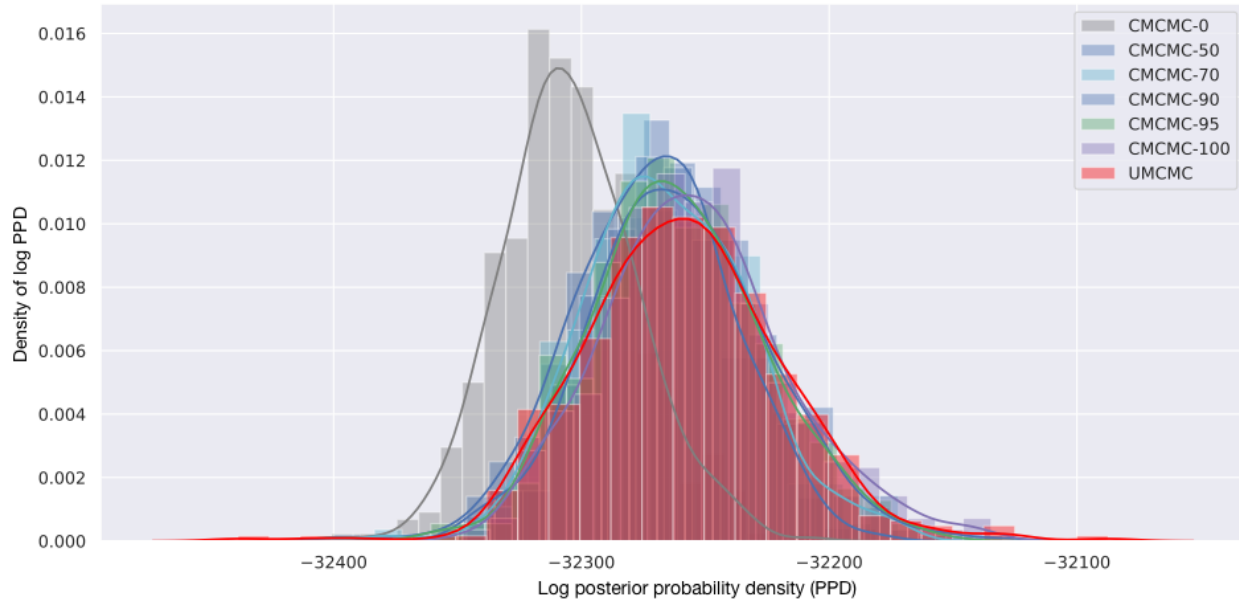


Figure 4: **Kernel Density Estimation of the posterior probability distribution of CMC MC and UMC MC samples.** Different methods are shown in different colors. All methods are run on the example sequence data which contains 16 taxa and 10 independent loci for the young divergence times and high population size scenario.

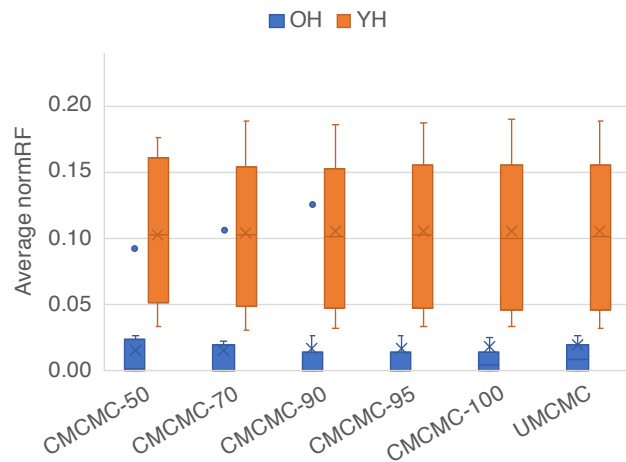


Figure 5: **Topological accuracy of the species trees inferred by CMC MC and UMC MC.** The data pertained to simulations of 10 loci from 16 taxa under the OH and YH scenarios.

3.3 CMC MC and UMC MC infer almost identical species trees

While Bayesian MCMC provides an approximation of the posterior distribution over parameters, the topology of the species tree is often the main quantity of interest. We compared the averaged Robinson-Foulds (RF) distance (Robinson and Foulds, 1981) between the true species tree and the inferred species tree topology in the 95% credible set to assess the topological accuracy. For each topology in the 95% credible set, we calculated the RF distance and divided it by the maximum possible RF distance (twice the number of internal branches in the true species tree) to derive the normalized RF distance (normRF; Kupczok et al., 2010). Then we summed up all normRF scaled by the frequency of each topology in the 95% credible set. More details about how to calculate averaged normRF distance are provided in the Evaluation Metrics section.

Figure 5 shows the topological accuracy of the species trees inferred by the various methods. As the figure shows, under

both evolutionary scenarios, all samplers infer almost the same species trees. This trend is also observed when varying the numbers of loci and taxa (Figure S7 and Figure S8).

3.4 Analysis of a biological data set

Recently, a study that applied exon capture sequencing to Australian rainbow skinks (Bragg *et al.*, 2018) compared the phylogenies inferred using summary MSC methods, a full Bayesian MSC method and concatenation. However due to the computational time required, the full Bayesian species tree method was only applied to 32 locus subsets of the data, despite 304 highly informative loci being available. This data set contains 46 taxa from 43 recognized species.

CMCMC enabled us to double the number of loci to 64, so inspired by Bragg *et al.*, we studied the effect of increasing the number of loci in a subset on the inferred phylogenies. Once the MCMC chain converges in one subset, we summarize the samples to generate the maximum clade credibility (MCC) tree. To quantify the variation between species trees inferred from different subsets, we calculated the normRF distance between the MCC tree from one MCMC chain or the inferred tree from ASTRAL (Mirarab *et al.*, 2014). As shown in Figure 6, CMCMC derived more consistent results compared with UMCML, as the highest normRF distance between CMCMC subsets was 0.23, but the highest pairwise distance between the smaller UMCML subsets was 0.3.

The precision of the 64 loci CMCMC posterior distributions was higher, as expected given the larger number of loci employed. For both normRF and branch score, the average distances between the maximum clade credibility (MCC) tree and individual samples in the posterior distribution were smaller for CMCMC (Figure S9 and Figure S10).

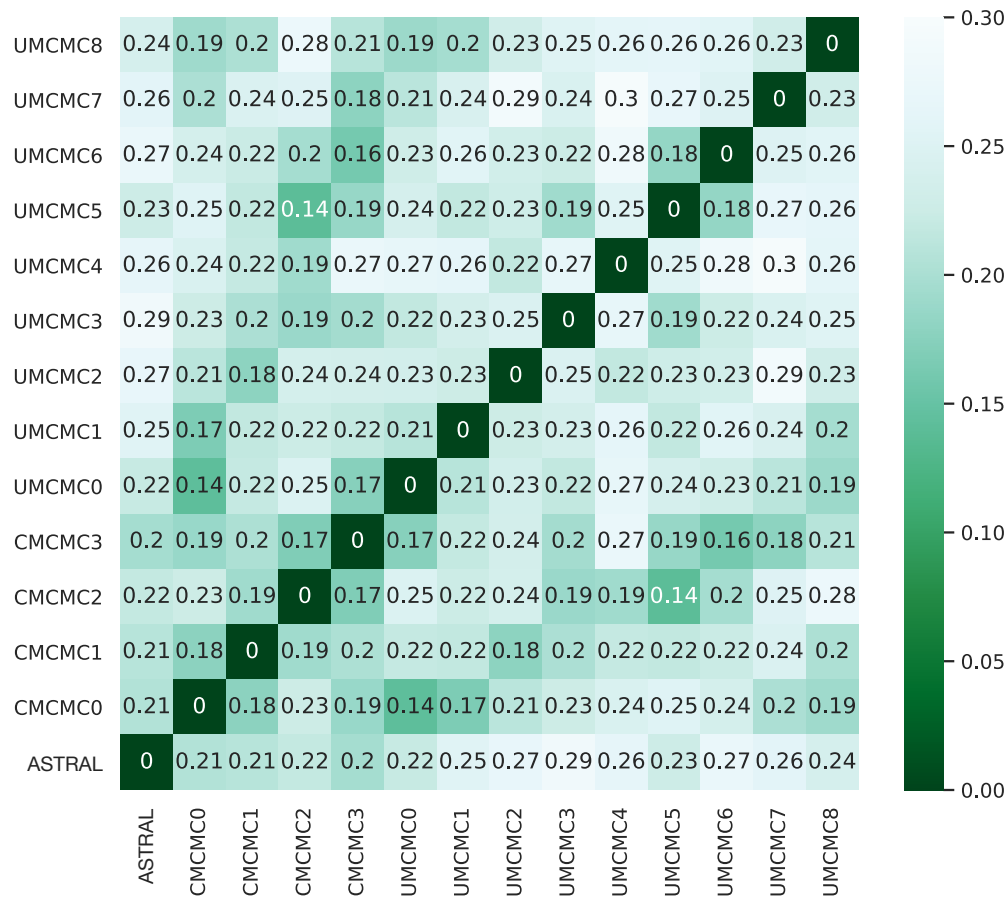


Figure 6: **Discordance among phylogenies estimated by ASTRAL, CMCMC and UMCML.** CMCMC was applied to four non-overlapping subsets of 64 loci each, and UMCML was applied to nine non-overlapping subsets of 32 loci each. The color and the number in each entry of the matrix indicates the normalized Robinson-Foulds distance between trees estimated by two different methods.

The increased precision of the CMCMC analyses enables taxonomic refinement of rainbow skinks. When 32 loci are used with UCMCMC, the trio *Carlia inconnexa*, *C. pectoralis* and *C. rubigo* form a clade but the relationships within that clade are unclear, as the best supported topology for this trio has *C. rubigo* as sister with an average posterior probability of 69% across subsets. When 64 loci are used with CMCMC, this average rises to 98% (Figure S11 and Figure S12).

4 Conclusions

In this paper we reported on a simple heuristic method for speeding the convergence of Bayesian MCMC under the multispecies coalescent. The heuristic works by restricting the space of gene trees that can be sampled. The constraints can be obtained in various ways including bootstrap trees contracted according to some support threshold, majority-rule consensus trees of posterior samples, or even constraints provided based on biological knowledge. As the approach restricts the explored space by design, we evaluated the method's performance in terms of convergence and, when converged, the distribution it samples. The evaluation was done on simulated data sets as well as a biological data set, and for evaluation metrics, we focused mainly on the time to convergence, the ASDSF between the constrained version of the sampler (CMCMC) and the unconstrained one (UCMCMC), as well as the topological accuracy of the inferred species trees.

Our findings indicate the CMCMC is a promising direction to pursue. CMCMC improves the number of iterations it takes for Bayesian MCMC to converge. CMCMC does not sacrifice the accuracy as long as the appropriate support threshold is selected. For the biological analysis, CMCMC enables greater precision, and therefore better inference and downstream analyses. Last but not least, CMCMC is easy to implement and combine with other techniques like annealed SMC.

5 Materials and Methods

5.1 CMCMC settings

CMCMC can be easily implemented in commonly used Bayesian based phylogenetics software packages such as PhyloNet (Wen et al., 2016) and BEAST2.5 (Bouckaert et al., 2019). In this paper, we implemented CMCMC in PhyloNet. In all our experiments, we generated 50 bootstrap trees for each locus and obtained the majority-rule consensus trees from those. Firstly, we generated 50 bootstrap trees given an alignment using RAxML (Stamatakis, 2014). Then, we estimated the constraint tree given a specific support threshold. MCMC chains for were run on a compute node with four Intel Xeon CPU E5-2650 v2 CPUs (with 2GB of memory). More implementation details are provided in Supplementary Materials.

6 Simulating data

For all simulated data sets, we used DendroPy (Sukumaran and Holder, 2010) to generate random species trees and ms (Hudson, 2002) to generate gene trees on these species trees under the multispecies coalescent. Sequence data were generated by Seq-Gen (Rambaut and Grass, 1997) under the Jukes-Cantor model (Jukes and Cantor, 1969). We derived the CGT for each locus by bootstrapping from the sequences by RAxML (Stamatakis, 2014).

Because species tree inference methods are employed over a range of evolutionary time scales and to clades with different population sizes, we varied both parameters for our simulation study which are shown in Table 1. For each simulated species tree we scaled its root to different heights; an "old" height of 50 million years ago (mya) akin to the rice-Pooideae split (Sandve et al., 2008), and a "young" height of 10 mya akin to the split of gorillas with humans and chimpanzees (Langergraber et al., 2012).

For both the old and young species tree scales, we simulated gene trees under large and small population sizes of 500,000 and 100,000 respectively, with annual generation times. Chimpanzees, gorillas and ancient humans all have effective population sizes (N_e) of around 20,000 individuals (Huff et al., 2010). Assuming a great ape generation time of 25 years, this population will have the same distribution of coalescent times as a clade of species with annual generation times and an N_e of 500,000, the same as our large population size condition. Given that human effective population sizes are often considered low, the small population size condition therefore corresponds to species with very low effective population sizes.

The evolutionary parameters also affect the proportion of resolved internal nodes of the CGT as shown in Figure S5. The "OH" and "OL" scenarios have higher proportion of resolved internal nodes than the "YH" and "YL" which means that the evolutionary mutation rate more effectively restrict the gene tree search space. The proportion of resolved

Table 1: **The evolutionary parameters varied to control the complexity and signal in the data.**

	Low population sizes	High population sizes
Old divergence times	OL: 50 mya, 100k individuals	OH: 50 mya, 500k individuals
Young divergence times	YL: 10 mya, 100k individuals	YH: 10 mya, 500k individuals

internal nodes in consensus trees decreases exponentially as the number of taxa increases as shown in Figure S6. In contrast, the population size does not have such obvious effect as the mutation rate or number of taxa.

More details on the simulations are provided in Supplementary Materials.

6.1 Biological data

We analyzed the Australian skinks data set which is provided in Bragg et al. (2018). We randomly selected one sample from each species. Note that the species names in the data set and in the paper are not consistent. More details about how to map the species in the data set and in the paper are shown in Table S1.

The Australian skinks data set contains three in-group genera: *Carlia*, *Lygisaurus* and *Liburnascincus*. There are 46 taxa from 43 recognized species. All details of the biological data including genus, species, tissue, collection, sample library and focal clade are provided in supplementary Table S2.

To obtain informative gene trees, we included 304 complete informative loci whose length ranges from 240 to 6,534 sites. Figure S13 shows the proportion of resolved internal nodes of constraint gene trees for different ranges of sequence length. In general, as the length of sequence increases the number of resolved internal nodes gets larger. This is because longer sequences are likely to contain more mutations to inform the resolution of nodes.

6.2 Evaluation metrics

6.2.1 Effective sample size

The Effective Sample Size or ESS is the number of effectively independent draws from some distributions sampled by the MCMC chain. Adequate ESS is a sign of good mixing of the MCMC chain and Kuhner (2009) argued that the ESS should be more than 200, a value that has been adopted in the Bayesian phylogenetics community. Therefore, an MCMC chain where the ESS of all selected parameters were higher than 200 was considered to have converged. These parameters were posterior, likelihood, prior, coalescent, tree height and population size.

6.2.2 Average Robinson-Foulds distance

When calculating the total Robinson-Foulds distance, we only considered posterior samples where the species tree topology was within the 95% credible set. We call this 95% credible set of posterior samples T^* and the full set of posterior samples T . To quantify differences between true tree t and the 95% credible set T^* , we calculate the averaged normRF distance (Kupczok et al., 2010) as

$$\frac{1}{|T|} \sum_{t^* \in T^*} \text{normRF}(t, t^*). \quad (3)$$

6.2.3 Average standard deviation of split frequencies

Average standard deviation of split frequencies (ASDSF) is a measure of convergence that has been used in tools such as ExaBayes (Aberer et al., 2014) and MrBayes (Ronquist et al., 2012). ASDSF can be calculated by comparing split or clades frequencies between two MCMC chains. Given two posterior distributions T_1 and T_2 from two MCMC chains and their corresponding 95% credible sets T_1^* and T_2^* , C is all unique, non-trivial clades in $T_1^* \cup T_2^*$. Set C^* is defined as

$$C^* = \{c \in C | \max(f(c, T_1), f(c, T_2)) \geq \epsilon\}, \quad (4)$$

where $f(c, T)$ is the frequency of clade c in distribution T , and ϵ is a threshold. We used $\epsilon = 0.1$, the same as the default setting in MrBayes 3.2 (Ronquist et al., 2012). Finally, the ASDSF between T_1 and T_2 is defined as

$$\text{ASDSF}(T_1, T_2) = \frac{1}{|C^*|} \sum_{c \in C^*} \frac{|f(c, T_1) - f(c, T_2)|}{2}, \quad (5)$$

because the standard deviation of two numbers is half of the absolute difference.

7 Supplementary Material

Supplementary Figures S1-S13, Tables S1-S2 and external tool commands are available in supplementary.pdf. All simulation data are available online: <https://rice.box.com/s/4ibr0b562fplbh4c1sndrpe49dwjs09g>.

8 Acknowledgments

This work was supported in part by NSF grants DBI-1355998, CCF-1514177, CCF-1800723, and DMS-1547433. Ana C. Afonso Silva assisted us in interpreting, and mapping the species names, of the biological data set.

References

- Aberer, A. J., K. Kobert, and A. Stamatakis (2014). Exabayes: massively parallel bayesian tree inference for the whole-genome era. *Molecular Biology and Evolution* 31(10), 2553–2556.
- Attias, H. (1999). Inferring parameters and structure of latent variable models by variational bayes. In *Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence, UAI'99*, San Francisco, CA, USA, pp. 21–30. Morgan Kaufmann Publishers Inc.
- Bouckaert, R., T. G. Vaughan, J. Barido-Sottani, S. Duchêne, M. Fourment, A. Gavryushkina, J. Heled, G. Jones, D. Kühnert, N. De Maio, M. Matschiner, F. K. Mendes, N. F. Müller, H. A. Ogilvie, L. du Plessis, A. Poppinga, A. Rambaut, D. Rasmussen, I. Siveroni, M. A. Suchard, C.-H. Wu, D. Xie, C. Zhang, T. Stadler, and A. J. Drummond (2019, 04). Beast 2.5: An advanced software platform for bayesian evolutionary analysis. *PLOS Computational Biology* 15(4), 1–28.
- Bragg, J. G., S. Potter, A. C. A. Silva, C. J. Hoskin, B. Y. Bai, and C. Moritz (2018). Phylogenomics of a rapid radiation: the australian rainbow skinks. *BMC Evolutionary Biology* 18(1), 15.
- Chifman, J. and L. Kubatko (2014). Quartet inference from snp data under the coalescent model. *Bioinformatics* 30(23), 3317–3324.
- Degnan, J. H. and N. A. Rosenberg (2009). Gene tree discordance, phylogenetic inference and the multispecies coalescent. *Trends in Ecology & Evolution* 24(6), 332–340.
- Felsenstein, J. (1978, 03). The number of evolutionary trees. *Systematic Biology* 27(1), 27–33.
- Flouri, T., X. Jiao, B. Rannala, and Z. Yang (2018). Species tree inference with bpp using genomic sequences and the multispecies coalescent. *Molecular Biology and Evolution* 35(10), 2585–2593.
- Fourment, M. and A. E. Darling (2019). Evaluating probabilistic programming and fast variational bayesian inference in phylogenetics. *bioRxiv*.
- Geyer, C. J. (1991). Markov chain Monte Carlo maximum likelihood. In E. M. Keramidas (Ed.), *Computing Science and Statistics: Proceedings of the 23rd Symposium on the Interface*, pp. 156–163.
- Hudson, R. R. (2002). Generating samples under a wright–fisher neutral model of genetic variation. *Bioinformatics* 18(2), 337–338.
- Huelsenbeck, J. P., F. Ronquist, R. Nielsen, and J. P. Bollback (2001). Bayesian inference of phylogeny and its impact on evolutionary biology. *Science* 294(5550), 2310–2314.
- Huff, C. D., J. Xing, A. R. Rogers, D. Witherspoon, and L. B. Jorde (2010). Mobile elements reveal small population size in the ancient ancestors of *Homo sapiens*. *Proceedings of the National Academy of Sciences* 107(5), 2147–2152.
- Jukes, T. H. and C. R. Cantor (1969). Evolution of protein molecules. In *Mammalian Protein Metabolism*, pp. 21–132. Academic Press.
- Kuhner, M. K. (2009). Coalescent genealogy samplers: windows into population history. *Trends in Ecology & Evolution* 24(2), 86–93.
- Kupczok, A., H. A. Schmidt, and A. von Haeseler (2010). Accuracy of phylogeny reconstruction methods combining overlapping gene data sets. *Algorithms for Molecular Biology* 5(1), 37.
- Lakner, C., P. Van Der Mark, J. P. Huelsenbeck, B. Larget, and F. Ronquist (2008). Efficiency of markov chain monte carlo tree proposals in bayesian phylogenetics. *Systematic Biology* 57(1), 86–103.
- Langergraber, K. E., K. Prüfer, C. Rowney, C. Boesch, C. Crockford, K. Fawcett, E. Inoue, M. Inoue-Muruyama, J. C. Mitani, M. N. Muller, M. M. Robbins, G. Schubert, T. S. Stoinski, B. Viola, D. Watts, R. M. Wittig, R. W. Wrangham, K. Zuberbühler, S. Pääbo, and L. Vigilant (2012). Generation times in wild chimpanzees and gorillas suggest earlier

- divergence times in great ape and human evolution. *Proceedings of the National Academy of Sciences* 109(39), 15716–15721.
- Liu, L. and L. Yu (2011). Estimating species trees from unrooted gene trees. *Systematic Biology* 60(5), 661–667.
- Liu, L., L. Yu, and S. V. Edwards (2010). A maximum pseudo-likelihood approach for estimating species trees under the coalescent model. *BMC Evolutionary Biology* 10(1), 302.
- Mazza, G., M. Menchetti, R. Sluys, E. Sola, M. Riutort, E. Tricarico, J.-L. Justine, L. Cavigioli, and E. Mori (2016). First report of the land planarian *Diversibipalium multilineatum* (makino & shirasawa, 1983)(platyhelminthes, tricladida, continenticola) in europe. *Zootaxa* 4067(5), 577–580.
- Mirarab, S., R. Reaz, M. S. Bayzid, T. Zimmermann, M. S. Swenson, and T. Warnow (2014). Astral: genome-scale coalescent-based species tree estimation. *Bioinformatics* 30(17), i541–i548.
- Ogilvie, H. A., R. R. Bouckaert, and A. J. Drummond (2017). Starbeast2 brings faster species tree inference and accurate estimates of substitution rates. *Molecular Biology and Evolution* 34(8), 2101–2114.
- Ogilvie, H. A., J. Heled, D. Xie, and A. J. Drummond (2016, 01). Computational performance and statistical accuracy of *BEAST and comparisons with other methods. *Systematic Biology* 65(3), 381–396.
- Rambaut, A., A. J. Drummond, D. Xie, G. Baele, and M. A. Suchard (2018, 04). Posterior summarization in Bayesian phylogenetics using tracer 1.7. *Systematic Biology* 67(5), 901–904.
- Rambaut, A. and N. C. Grass (1997). Seq-gen: an application for the monte carlo simulation of dna sequence evolution along phylogenetic trees. *Bioinformatics* 13(3), 235–238.
- Robert, C. and G. Casella (2011, 02). A short history of markov chain monte carlo: Subjective recollections from incomplete data. *Statistical Science* 26(1), 102–115.
- Robinson, D. F. and L. R. Foulds (1981). Comparison of phylogenetic trees. *Mathematical Biosciences* 53(1-2), 131–147.
- Ronquist, F., M. Teslenko, P. van der Mark, D. L. Ayres, A. Darling, S. Höhna, B. Larget, L. Liu, M. A. Suchard, and J. P. Huelsenbeck (2012, 02). MrBayes 3.2: efficient Bayesian phylogenetic inference and model choice across a large model space. *Systematic Biology* 61(3), 539–542.
- Sandve, S. R., H. Rudi, T. Asp, and O. A. Rognli (2008, Sep). Tracking the evolution of a cold stress associated gene family in cold tolerant grasses. *BMC Evolutionary Biology* 8(1), 245.
- Stamatakis, A. (2014). Raxml version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* 30(9), 1312–1313.
- Stensvold, C. R., M. Lebbad, and C. G. Clark (2011). Last of the human protists: the phylogeny and genetic diversity of iodamoeba. *Molecular Biology and Evolution* 29(1), 39–42.
- Stunžėnas, V., R. Petkevičiūtė, and G. Stanevičiūtė (2011). Phylogeny of *Sphaerium solidum* (bivalvia) based on karyotype and sequences of 16s and its 1 rdna. *Central European Journal of Biology* 6(1), 105–117.
- Sukumaran, J. and M. T. Holder (2010). Dendropy: a python library for phylogenetic computing. *Bioinformatics* 26(12), 1569–1571.
- Wang, L., S. Wang, and A. Bouchard-Côté (2019, 06). An Annealed Sequential Monte Carlo Method for Bayesian Phylogenetics. *Systematic Biology*. Published online in advance of print.
- Wang, Y. and L. K. Nakhleh (2018). Towards an accurate and efficient heuristic for species/gene tree co-estimation. *Bioinformatics* 34 17, i697–i705.
- Wen, D. and L. Nakhleh (2017). Coestimating reticulate phylogenies and gene trees from multilocus sequence data. *Systematic Biology* 67(3), 439–457.
- Wen, D., Y. Yu, and L. Nakhleh (2016). Bayesian inference of reticulate phylogenies under the multispecies network coalescent. *PLOS Genetics* 12(5), e1006006.
- Wen, D., Y. Yu, J. Zhu, and L. Nakhleh (2018). Inferring phylogenetic networks using PhyloNet. *Systematic Biology* 67(4), 735–740.
- Yu, Y., J. Dong, K. J. Liu, and L. Nakhleh (2014). Maximum likelihood inference of reticulate evolutionary histories. *Proceedings of the National Academy of Sciences* 111(46), 16448–16453.
- Zhang, C. and F. A. Matsen (2019). Variational bayesian phylogenetic inference. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=SJVmjjR9FX> accessed September 15 2019.
- Zhang, C., H. A. Ogilvie, A. J. Drummond, and T. Stadler (2017). Bayesian inference of species networks from multilocus sequence data. *Molecular Biology and Evolution* 35(2), 504–517.

Zhu, J., D. Wen, Y. Yu, H. M. Meudt, and L. Nakhleh (2018). Bayesian inference of phylogenetic networks from bi-allelic genetic markers. *PLOS Computational Biology* 14(1), e1005932.