# Quantifying MCMC exploration of phylogenetic tree space

Christopher Whidden and Frederick A. Matsen IV

October 20, 2014

**Abstract**

In order to gain an understanding of the effectiveness of phylogenetic Markov chain Monte Carlo (MCMC), it is important to understand how quickly the empirical distribution of the MCMC converges to the posterior distribution. In this paper we investigate this problem on phylogenetic tree topologies with a metric that is especially well suited to the task: the subtree prune-and-regraft (SPR) metric. This metric directly corresponds to the minimum number of MCMC rearrangements required to move between trees in common phylogenetic MCMC implementations. We develop a novel graph-based approach to analyze tree posteriors and find that the SPR metric is much more informative than simpler metrics that are unrelated to MCMC moves. In doing so we show conclusively that topological peaks do occur in Bayesian phylogenetic posteriors from real data sets as sampled with standard MCMC approaches, investigate the efficiency of Metropolis-coupled MCMC (MCMCMC) in traversing the valleys between peaks, and show that conditional clade distribution (CCD) can have systematic problems when there are multiple peaks.

The Bayesian paradigm has been extensively adopted to infer phylogenetic trees and associated parameter values in a consistent probabilistic framework. [We are interested in convergence properties on the discrete structure of unrooted tree topologies, so for the purposes of this paper we will use the word *tree* without further qualification to signify an unrooted leaf-labeled tree topology without branch lengths.] Current Bayesian phylogenetic methods rely on being able to move efficiently through tree hypothesis space with a random walk via Markov chain Monte Carlo (MCMC) (Metropolis et al. 1953; Hastings 1970). These include the widely used BEAST (Drummond and Rambaut 2007; Drummond et al. 2012; Bouckaert et al. 2014) and MrBayes (Ronquist et al. 2012) software packages as well as more recent methods such as BAli-Phy (Suchard and Redelings 2006), RevBayes (`http://github.com/revbayes/revbayes`) and ExaBayes (`http://sco.h-its.org/exelixis/web/software/exabayes/index.html`). The empirical distribution of suitably spaced MCMC samples converges to its true posterior distribution given an infinitely long run of the MCMC (reviewed in Tierney 1994). However, in order to obtain accurate computations of trees and associated confidence levels in practice, it is essential that these Markov chains explore phylogenetic "tree space" efficiently.

Many important questions remain unanswered concerning the practical performance of MCMC for phylogenetics, such as the presence and frequency of multiple peaks (i.e. modes) in phylogenetic tree posteriors, and the ability of chains to move between these posterior peaks. Also, to what extent are "peaky" (i.e. multimodal) posteriors a consequence of the discrete structure of phylogenetic trees, or simply a consequence of simultaneously estimating a large number of real parameters? Do strategies such as Metropolis-coupled MCMC (MCMCMC or (MC)$^3$) (Geyer 1992; Huelsenbeck and Ronquist 2001), which are helpful for multimodal distributions in the real case, effectively solve the problem? To what extent do convergence diagnostics
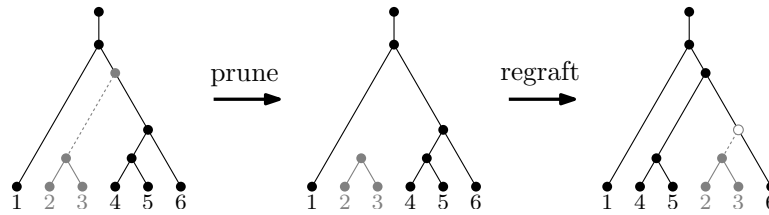
based on tree topologies, such as average standard deviation of split frequencies between independent Markov chains, imply that the empirical distribution of the underlying discrete tree topologies is close to the actual posterior? How many independent chains are required for such convergence diagnostics to adequately assess the level of convergence?

There are continuing (Lakner et al. 2008; Štefankovic and Vigoda 2011; Höhna and Drummond 2012) and sometimes vitriolic (Mossel and Vigoda 2005; Ronquist et al. 2006; Mossel and Vigoda 2006) debates concerning how well MCMC methods explore tree space. Lakner et al. (2008) and Höhna et al. (2008) showed that the random choices of operations used in current methods lead to a low rate of accepted transitions and increase the amount of computation required before MCMC runs achieve a given split frequency distance to golden runs. To address this problem, Höhna and Drummond (2012) introduced improved Metropolized Gibbs samplers—biased operators that use additional computation to select transitions with a higher acceptance rate—and showed that these operators reduced the time to achieve such a given split frequency distance to golden runs using BEAST on 11 empirical data sets. Parsimony-biased tree proposals have been included in MrBayes 3.2 (Ronquist et al. 2012). Mossel and Vigoda (2005) showed mathematically that MCMC methods can give misleading results when the alignments used to construct the trees derive from a site-wise mixture of data generated on two very different trees (note that this usage of "mixture" refers to a means of combining probability distributions, whereas the separate concept of "mixing" as described below refers to a characteristic of Markov chains). On such a site-wise mixture, the Markov chain appears to converge rapidly according to diagnostics but in actuality requires an exponential amount of time to converge due to the large "valleys" of unlikely trees between the two site-wise mixture peaks. Such site-wise mixtures are but one contrived example of a peaky distribution. However, even if we never see the sort of data set they postulate we may still encounter peaky distributions. In such a

situation, the posterior samples from a single peak may appear as though the chain has completely explored the relevant part of tree space, leading to a mistakenly high confidence value for an incomplete sample of trees. Although there has been extensive discussion in the literature about to what extent Metropolis-coupling helps traverse peaks, there have been few conclusions, probably because there hasn't been a clear exploration of peaks and peakiness in phylogenetic posteriors.

Some studies have focused on estimating mixing properties of phylogenetic MCMC using theory (Aldous 2000; Mossel and Vigoda 2005); this is known to be a very hard problem and can only be done in "toy" examples. [As is standard in the field, we will use the word *mixing* to refer to the convergence of the empirical distribution of MCMC samples to their posterior distribution.] Even when we can diagnose the failure of a Markov chain to converge to the posterior distribution, it does not lead to an understanding of why the failure occurred. A more practical approach to understanding movement in discrete tree space is to equip this space with a metric and consider distances traveled by the chain.

Recently, Höhna and Drummond (2012) and Larget (2013) proposed using Conditional Clade Probability (CCP) and Conditional Clade Distribution (CCD) methods, respectively, to approximate tree posterior probabilities. In both methods, the probability of a tree is estimated based on a product of conditional clade probabilities. Larget (2013) uses the approximation that compatible splits, separated by another split, are approximately conditionally independent given the separating split. The approximating equation of CCD is then a product of joint conditional sister clade probabilities, given the parent clade. Conditional probability methods have the potential to estimate the posterior probabilities of many trees using only a small sample of the tree posterior. They have already been productively applied to approximate tree posteriors in phylogenomic analyses (Szöllősi et al. 2013). However, the validity of the assumption of conditional independence of sister clades, given the parent clade, is not clear in practice.

**Figure 1.** An SPR move.

It is thus crucial to determine the accuracy of the CCD approximation on real data sets.

These considerations motivate improved methods to understand the performance of phylogenetic methods and the corresponding "topography" of trees. Hillis et al. (2005) used the Robinson-Foulds (RF) distance (Robinson and Foulds 1981) between phylogenies with multidimensional scaling (MDS) to visualize tree space. However, the RF distance does not correspond to SPR operators, and in fact may be arbitrarily large even for trees separated by a single SPR operation. Matsen (2006) suggested using the nearest-neighbour interchange (NNI) distance with MDS visualization. Höhna and Drummond (2012) used this idea to visualize "islands" among 15 trees from the 27 taxon tree space. Still, the NNI distance does not correspond closely with many rearrangements used in phylogenetic inference and is difficult to compute, limiting the utility of this method.

Subtree prune-and-regraft (SPR) (Hein et al. 1996) moves are the most common rearrangements used by phylogenetic programs (Höhna and Drummond 2012). These involve cutting a subtree off and attaching it somewhere else (Fig. 1). The minimum number of such operations required to transform one tree into another is called the SPR distance. Moreover, SPR operators are closely related to other common rearrangements. NNI operators are a subset of SPR operators. Two other common operators, the subtree swap (SS) and tree-bisection-and-reconnection (TBR) are each equivalent to two SPR operations (Höhna and Drummond 2012).

Thus the SPR distance is especially appropriate to investigate phylogenetic

5

MCMC behaviour in this setting because of the correspondence between SPR operators and most MCMC moves. However, SPR distance is challenging to use due to the computational complexity of its computation (Allen and Steel 2001; Bordewich and Semple 2005; Hickey et al. 2008). Recently, efficient fixed-parameter algorithms for computing the SPR distance have been developed and implemented in the freely available and open source RSPR software package (Whidden and Zeh 2009; Whidden et al. 2010, 2013, 2014). These efficient algorithms require fractions of a second to compute SPR distances between trees with hundreds of taxa and enable, for the first time, tree comparison using the SPR distance on a relevant scale.

In this paper, we use SPR tree space to visualize and analyze Bayesian phylogenetic posterior distributions. Our graph-based method directly shows the difficulty in moving between areas of tree space and can identify topological peaks that are not visible in multidimensional scaling projections. We show that our SPR graphs explain the error rate and time to a given average standard deviation of split frequencies (ASDSF) (Ronquist et al. 2012) of Bayesian phylogenetic methods on various data sets when these statistics do not correlate with the number of taxa alone. Moreover, we show that multiple topological peaks are common in nontrivial posteriors, even with relatively few taxa, and that the graphs can be used to identify bottlenecks in posterior distributions: regions of tree space between peaks that are difficult for MCMC methods to cross. We propose a topological variant of the Gelman-Rubin convergence diagnostic and show that a small ASDSF often implies a small such topological convergence diagnostic. We explore the effect of Metropolis-coupling and show that it greatly improves mixing, particularly between topological peaks, and reduces the number of MCMC iterations required for multiple runs to achieve a given ASDSF threshold. Metropolis-coupling improves overall performance in peaky distributions but may increase computation time in non-peaky distributions, in which case we observe the number of iterations to be reduced by a smaller factor than the number of

Metropolis-coupled chains. For both MCMCMC and single-chain approaches, we find that the current standard of two runs to calculate ASDSF is insufficient to obtain a proper error estimate. Finally, we show that independence of sister clades, conditioned on parent clades, does not hold in some peaky distributions. This causes the CCD distribution to systematically underestimate the probability of trees within alternative peaks and systematically overestimate the probability of trees between peaks.

# METHODS

## *Computing the SPR distance*

We modified RSPR, the open source C++ software package for computing subtree prune-and-regraft distances (Whidden et al. 2010, 2013, 2014). Previous versions of RSPR computed the SPR distance between two input trees or the aggregate SPR distance from a single tree to a set of trees. Our new version 1.3 of RSPR (`https://github.com/cwhidden/rspr/`) adds support for computing pairwise SPR and RF distance matrices. These distance matrices can be used as input to multidimensional scaling methods or to compute tree space graphs.

RSPR computes a maximum agreement forest (MAF) (Hein et al. 1996; Allen and Steel 2001) of two rooted trees with a fixed-parameter algorithm. An agreement forest is a forest of subtrees that can be obtained by cutting edges from both trees. An MAF is obtained by cutting the fewest possible number of edges. This smallest number of cut edges is equivalent to the SPR distance between the trees if they are rooted (Bordewich and Semple 2005). The time required for this fixed-parameter algorithm increases exponentially with the distance computed but only linearly with the size of the trees. In particular, the algorithm can quickly determine whether two rooted trees are separated by a single SPR operation. In practice, RSPR can compute SPR distances between trees

with hundreds of taxa and more than 50 transfers in fractions of a second (Whidden et al. 2014).

Unrooted trees are commonly inferred by phylogenetic methods including MrBayes. However, an MAF of two unrooted trees is equivalent to their tree-bisection-and-reconnection distance (Allen and Steel 2001) and no MAF formulation is known for the SPR distance of unrooted trees. For unrooted trees, we thus consider each possible rooting of the trees and choose the rootings which give the minimal SPR distance. This "best rooting" SPR distance should closely agree with the unrooted SPR distance except in pathological cases where the minimum set of unrooted SPR operations is incompatible with any rooting (e.g. Supplemental Figure 1). In particular, both are guaranteed to agree when the trees are separated by a single SPR operation; much of our work here uses the graph induced by these single SPR moves.

## *SPR tree space graphs*

We used SPR-based graphs, restricted to sets of high probability trees, to model the SPR tree space of Bayesian phylogenetic posterior distributions. We selected these sets of high probability trees as follows. First, we ordered the trees from a posterior sample by descending posterior probability (ties broken by sample order). In cases with a large number of ties (e.g. where every tree is sampled once or twice), breaking ties with sample order may cause bias, so we broke ties randomly in such cases. The 95% credible set is the smallest set of trees at the head of this list with cumulative posterior probability more than 95%. We call the $m$ trees with highest posterior probability the "top $m$ trees," that is, the first $m$ trees in this list. We used $m = 4096$ in our tests unless otherwise noted, and generally used the 95% credible set when it contained fewer than 4096 trees, and the top 4096 trees when it was not. We call these sets of at most 4096 trees the "top trees".

We define the SPR graph for a set of trees $T$ to be the undirected graph $G_T = (V, E)$ such that each tree is represented by a node in $V$ and two trees are

connected by an edge in $E$ if and only if they are separated by an SPR distance of 1. In particular, we constructed a distance matrix $D$ such that an entry $D_{ij} = 1$ if, and only if, the SPR distance between $i$ and $j$ is 1. We constructed such graphs using RSPR version 1.3, then converted these matrices to an edge list format suitable for input to graph visualization software.

## *Clustering high-probability regions of tree space*

We used a simple iterative clustering procedure to aid in the detection of topological peaks. These peaks are intuitively defined as a set of topologies with relatively high probability surrounded by topologies with low probability. Any useful clustering procedure must therefore make use of posterior probabilities in addition to topology, moreover, comparing every pair of trees is computationally expensive even with the simple goal of computing RF distances. We thus employed the following approximate iterative clustering algorithm. First select the most probable topology as the center of our first cluster. Then compare the current cluster center to each unclustered tree, and add each tree within a specified SPR distance radius to the current cluster. This procedure proceeds iteratively, grouping the most probable unclustered topology and the remaining set of unclustered trees until each tree has been clustered or a given number of clusters assigned. For a given cluster center, we used a clustering radius equal to the mean SPR distance from the current cluster center to each unclustered tree, minus the standard deviation of these distances (i.e. $\mu - \sigma$). This radius is recalculated for each new cluster. We stopped this process after 8 clusters had been identified.

## *Graph visualization with Cytoscape*

SPR graphs were visualized with the open source Cytoscape platform (Shannon et al. 2003). In addition to the edge list and clusters described above, we computed SPR

distances between the tree with highest posterior probability and the top $m$ trees. We visualized tree space in three ways: (1) distance SPR graphs, (2) cluster SPR graphs, and (3) weighted MCMC graphs. To visualize SPR graphs we used a force-directed graph layout, which essentially means that graph nodes are pushed away from each other, but edges act as "springs" that attempt to maintain a uniform length. We scaled node sizes (area) in proportion to tree posterior probability. The largest node represents the tree with highest posterior probability. We hypothesized that peaks would be visible in such graphs as sets of relatively large (high probability) nodes separated by relatively small (low probability) nodes or in disconnected graph components. In distance SPR graphs, graph nodes are colored on a red-yellow-white scale (dark-light in the print version) with increasing SPR distance from the most probable topology. We further hypothesized that difficult to sample peaks would be visible in distance SPR graphs as large yellow or white nodes. In clustered SPR graphs, graph nodes are colored by cluster. We expected that any significant topological peaks would be grouped in different clusters and therefore receive different colors. Finally, we used another type of graph to visualize Markov chain movement between trees to validate our assumption that SPR tree space corresponds to MCMC movement in practice. These graphs represent movement between MCMC samples (including Metropolis-coupling chain swaps where applicable). We weighted these edges with the number of such transitions and visualized these edge weights using edge thickness and color. Note, however, that posteriors are typically subsampled every given number of iterations, and we followed this practice. Given such subsamples, some of the dependence between sample tree and order may be eliminated and care must be taken when interpreting such graphs.

## *Quantifying tree space mixing*

To quantify mixing behaviour in tree space we computed statistics based on mean access times (MAT)—the mean number of iterations required to transition between

topologies in an MCMC search (Lovász 1993). As with our graph clustering, computing statistics for each pair of trees can be computationally expensive and difficult to visualize. Rather than directly considering access time statistics for each pair of trees, we instead computed the mean commute time (MCT) (Lovász 1993) from the most probable topology to each other high probability tree and back: the sum of pairwise MATs. We also considered a new measure, the mean round trip cover time. This is the mean number of iterations required to cover (visit) each high probability tree, starting from and returning to the highest probability tree. This measure is essentially a round-trip analog of the mean cover time (Lovász 1993). The MAT values (and hence MCT and round trip cover time values) can be computed with a single pass through the tree posterior using a method for updating weighted means (see e.g. West 1979). Formal definitions of these statistics and a description of our dynamic programming method for computing them can be found in the supplementary material.

## *A discrete topological Gelman-Rubin-like convergence diagnostic*

In order to avoid having to project trees down to vectors of split frequencies in order to diagnose convergence, we developed a discrete topological variant of the Gelman-Rubin convergence diagnostic (Gelman and Rubin 1992). The Gelman-Rubin convergence diagnostic for a *real-valued* parameter $x$ requires multiple independent Markov chains and compares the variance within chains and between chains, as we review now. Note that by "chains" here we refer to multiple independent chains, which are equivalent to the MrBayes terminology "runs" rather than Metropolis-coupled chains. Suppose we have $m$ chains, each with $n$ sampled values. The value of chain $i$ at iteration $j$ is denoted $x_{ij}$. The variance between chains, $B$, is estimated by the variance between the $m$ sequence means, $\bar{x}_{i.}$, each based on $n$ values of $x$. That, is,

$$B/n = \frac{1}{m-1} \sum_{i=1}^{m} (\bar{x}_{i.} - \bar{x}_{..})^2,$$

where $\bar{x}_{..} = \frac{1}{m} \sum_{i=1}^{m} \bar{x}_{i.}$. The variance within chains, $W$, is the average of the $m$ within-sequence variances, $s_i^2$, each based on $n-1$ degrees of freedom. That is,

$$W = \frac{1}{m} \sum_{i=1}^{m} s_i^2,$$

where $s_i^2 = \frac{1}{n-1} \sum_{j=1}^{n} (x_{ij} - \bar{x}_i)^2$. The estimated variance is then a weighted average of $W$ and $B$,

$$\hat{V} = \left(1 - \frac{1}{n}\right) W + \frac{1}{n} B.$$

The potential scale reduction factor (PSRF) is defined as $\hat{R} = \sqrt{\hat{V}/W}$. This measures the potential for reducing the difference between $B$ and $W$. $B$ initially overestimates the variance, given multiple chains with overdispersed starting points. $W$ initially underestimates the variance, as it is based on an incomplete sample from a limited region of the parameter space. These values converge as the independent chains converge. As such, the PSRF approaches 1 as the chains converge.

Our topological Gelman-Rubin-like convergence diagnostic estimates the differences within and between Markov chains in terms of topological changes. There is no concept of sample mean for topologies, so we compute an analogous statistic with the mean square deviation instead of variance. In particular, we estimate the SPR distance deviation within and between chains. Again, $x_{ij}$ denotes the tree from chain $i$ at iteration $j$. Let $d(x_{i_1 j_1}, x_{i_2 j_2})$ denote the distance between two such trees.

$W$ is the mean square deviation within a chain:

$$s_i^2 = \frac{1}{n(n-1)} \sum_{j_1=1}^{n} \sum_{j_2=1}^{n} d(x_{ij_1}, x_{ij_2})^2.$$

Similarly, we estimated the between-chain deviation by comparing each chain to the aggregate set of chains :

$$B = \frac{1}{(m-1)mn^2} \sum_{i_1=1}^{m} \sum_{i_2=1}^{m} \sum_{j_1=1}^{n} \sum_{j_2=1}^{n} d(x_{i_1 j_1}, x_{i_2 j_2})^2.$$

12

With this formulation, $\sqrt{V}$ estimates the topology root mean square deviation (RMSD). $\hat{R}$ is computed as before.

As written, these formulas require a great deal of computation. To efficiently compute topological PSRF values, observe that there are many repeated comparisons between identical trees. We thus grouped identical topology comparisons, computed one SPR distance for each and weighted the squared distances accordingly in our calculations. We also limited our comparisons to the top trees, as in our SPR graph construction. We normalized our computations by the number of included distances rather than the total number of samples $n$. Using this method, $B$ is no more complex to compute than $W$.

As with the original Gelman-Rubin convergence diagnostic, the topological PSRF value approaches 1 as the independent chains converge. B initially overestimates the RMSD between topologies, given multiple chains with overdispersed starting points. W initially underestimates the RMSD between topologies, as it is based on an incomplete sample from a limited region of tree space. These values converge as the independent chains converge. We use the name topological Gelman-Rubin-*like* to emphasize that it is inspired by the original but is not the same.

## *Multidimensional scaling*

Multidimensional scaling (MDS) is a method for projecting complex data to a small number of dimensions suitable for visualization (Kruskal 1964a,b). Non-metric MDS is typically applied to create a new two or three dimensional space from a given pairwise distance matrix in a way that preserves the pairwise distances as much as possible. Specifically, it minimizes a *stress function* quantifying the difference between the original distances and Euclidean distances in the projected space. Multidimensional scaling has been used previously to visualize RF distances between trees in a posterior distribution (Hillis et al. 2005) and, on a limited scale, NNI distances (Höhna and

Drummond 2012). We applied MDS to SPR and RF distance matrices using the R `isomds` function from the MASS package (Venables and Ripley 2002).

## *Conditional clade probability*

Recently, the conditional clade probability (CCP) and conditional clade distribution (CCD) concepts have been proposed by Höhna and Drummond (2012) and Larget (2013), respectively. These methods use conditional products of split posterior probabilities on splits to estimate the corresponding phylogenetic posterior probabilities. To test the conditional independence assumption in practice, we applied the CCD software of Larget (2013) to compute conditional clade probabilities and compare the results to posterior probabilities on large posterior samples.

## *Number of runs and chains*

Two MrBayes run parameters are of particular importance to obtain ASDSF estimates that reflect the level of convergence to the posterior distribution: the number of independent runs used for testing ASDSF convergence and the number of Metropolis-coupling chains. The number of independent runs determines the behavior of the average standard deviation of split frequencies (ASDSF) convergence diagnostic, which compares split frequencies between independent runs. As is typical, our ASDSF calculations only consider splits with a frequency exceeding 10% in at least one of the runs. We follow previous researchers by using a 0.01 cutoff for ASDSF as a stopping rule. In the MrBayes version 3.2 manual, Ronquist et al. (2011) suggest that "an average standard deviation below 0.01 is very good indication of convergence, while values between 0.01 and 0.05 may be adequate depending on the purpose of your analysis." Increasing the number of runs increases the stringency of ASDSF convergence at a given limit at the expense of increased computation. Metropolis-coupling (Geyer 1992;

Huelsenbeck and Ronquist 2001) is a commonly applied method to improve MCMC mixing in peaky distributions. In addition to the primary "cold" Markov chain, from which posterior samples are drawn, multiple "hot" chains are maintained. These hot chains typically move more freely through the parameter space. The cold chain is periodically swapped with a hot chain to "jump" through the parameter space.

## *Implementation*

We developed the open source software package sprspace (`https://github.com/cwhidden/sprspace`) to construct SPR graphs. This software package also implements our clustering routine, prepares graph visualizations for Cytoscape, computes access times and commute times and computes our topological Gelman-Rubin-like measure. Our software allows users to specify a fixed clustering radius in case dynamic cluster radius selection provides poor results. Moreover, users may modify the number of top trees considered to change the amount of computation required.

## *Data and run-time parameters*

We investigated MCMC estimation on unrooted trees by applying MrBayes 3.2 (Ronquist et al. 2012) to 17 empirical data sets. The first group of data sets, which we will call DS1-DS11, have become standard data sets for evaluating MCMC methods (Lakner et al. 2008; Höhna and Drummond 2012; Larget 2013). These data sets consist of sequences from 27 to 71 eukaryote species (Table 1), and are fully described elsewhere (Lakner et al. 2008). Note that TreeBASE identifiers for these data sets have changed from those used in some previous publications (Supplemental Table 1). The second group of data sets, which we will call VL1-VL6, consist of alignments with 40 to 63

**Table 1.** The data sets used in this study, DS1-11 (eukaryote) and VL1-6 (bacterial/archaeal). N = number of species; Cols = number of nucleotides; Est error = Estimated maximum standard error of split frequencies in golden runs (in %); rDNA = ribosomal DNA; rRNA = ribosomal RNA; mtDNA = mitochondial DNA; COII = cytochrome oxidase subunit II GARTFase = phosphoribosylglycinamide formyltransferase 2.

| Data | N | Cols | Type of data | Study | Est error |
|------|-----|------|--------------|-------|-----------|
| DS1 | 27 | 1949 | rRNA; 18s | Hedges et al. (1990) | 0.0048 |
| DS2 | 29 | 2520 | rDNA; 18s | Garey et al. (1996) | 0.0002 |
| DS3 | 36 | 1812 | mtDNA; COII (1678); cytb (679-1812) | Yang and Yoder (2003) | 0.0002 |
| DS4 | 41 | 1137 | rDNA; 18s | Henk et al. (2003) | 0.0006 |
| DS5 | 50 | 378 | Nuclear protein coding; wingless | Lakner et al. (2008) | 0.0005 |
| DS6 | 50 | 1133 | rDNA; 18s | Zhang and Blackwell (2001) | 0.0023 |
| DS7 | 59 | 1824 | mtDNA; COII; and cytb | Yoder and Yang (2004) | 0.0011 |
| DS8 | 64 | 1008 | rDNA; 28s | Rossman et al. (2001) | 0.0009 |
| DS9 | 67 | 955 | Plastid ribosomal protein; s16 (rps16) | Ingram and Doyle (2004) | 0.0164 |
| DS10 | 67 | 1098 | rDNA; 18s | Suh and Blackwell (1999) | 0.0164 |
| DS11 | 71 | 1082 | rDNA; internal transcribed spacer | Kroken and Taylor (2000) | 0.0008 |
| VL1 | 40 | 271 | UDP-2,3-diacylglucosamine hydrolase | Beiko et al. (2006) | 0.0019 |
| VL2 | 44 | 472 | coproporphyrinogen III oxidase | Beiko et al. (2006) | 0.0007 |
| VL3 | 50 | 442 | GARTFase | Beiko et al. (2006) | 0.0050 |
| VL4 | 52 | 129 | hypothetical protein | Beiko et al. (2006) | 0.0484 |
| VL5 | 53 | 349 | fructose-1,6-bisphosphatase | Beiko et al. (2006) | 0.0070 |
| VL6 | 63 | 294 | pyridoxine 5'-phosphate synthase | Beiko et al. (2006) | 0.0542 |

bacterial and archaeal sequences (Table 1) of protein-coding genes, and are fully described elsewhere (Beiko et al. 2006).

To analyze the level of convergence to the posterior distribution, we computed large "golden run" posterior samples for each data set, meaning that we repeatedly ran the chains well past the typical number of iterations used for such analyses: for each of our 17 data sets, 10 single-chain MrBayes replicates were run for one billion iterations and sampled every 1000 iterations. These replicates were not Metropolis-coupled. We discarded the first 25% of samples as "burn-in" for a total of 7.5 million posterior samples per data set, and assumed that this long burn-in period implied stationarity, i.e. that after burn-in the chain was sampling from the stationary distribution of the MCMC. Following Höhna and Drummond (2012), we assumed these runs accurately estimated posterior split frequency distributions because of the extreme length of these Markov chains in comparison to our data size. To test this assumption, we estimated the split frequency error between replicated golden runs (maximum standard error of any split) as in Höhna and Drummond (2012) (see Table 1). The estimated split frequency error was below 0.06% for each of our data sets, suggesting that the various golden runs are sampling the same split frequencies. Moreover, commonly applied diagnostics implemented in the MrBayes `sumt` and `sump` tools satisfied common thresholds (Supplemental Table 2), including having a standard error of log likelihoods at most 2.11, maximum standard deviation of split frequencies at most 0.015 (0.007 for all but DS1), maximum Gelman-Rubin split PSRF values of 1.000, and the effective sample size (ESS; a measure of the number of samples correcting for MCMC autocorrelation) for the treelength parameter (the sum of branch lengths) exceeding 650,000.

We cannot similarly assume that these golden runs have accurately estimated the posterior probability of all topologies. We do, however, assume that the golden runs have accurately estimated the posterior probability of high probability topologies, namely the top trees taken from the combined golden runs. To test this assumption, we

17

estimated the topological error between replicated golden runs (maximum standard error of the posterior probability of the top trees) for the eukaryote datasets, analogously to the split frequency error calculation (Supplemental Figure 2 and Supplemental Table 3). The estimated standard error among high probability topologies was generally at least an order of magnitude smaller than the posterior probability, validating this assumption. However, datasets DS9 and DS11 were notable exceptions as each topology was sampled exactly once, with no overlap between runs. As such, we do not assume that the empirical distribution on topologies for DS9 and DS11 are close to their posterior distributions.

We ran MrBayes on each of our data sets with 10 replicates of a varying number of runs (2 through 8) and chains (1 or 4) until the runs had ASDSF less than 0.01 or a maximum of 100 million iterations. We sampled these runs every 100 iterations and again discarded the first 25% of samples. We then compared the effect of these parameters on running time and error in practice, where error was measured by the root mean square deviation (RMSD) of split frequencies as compared to the golden runs.

# RESULTS

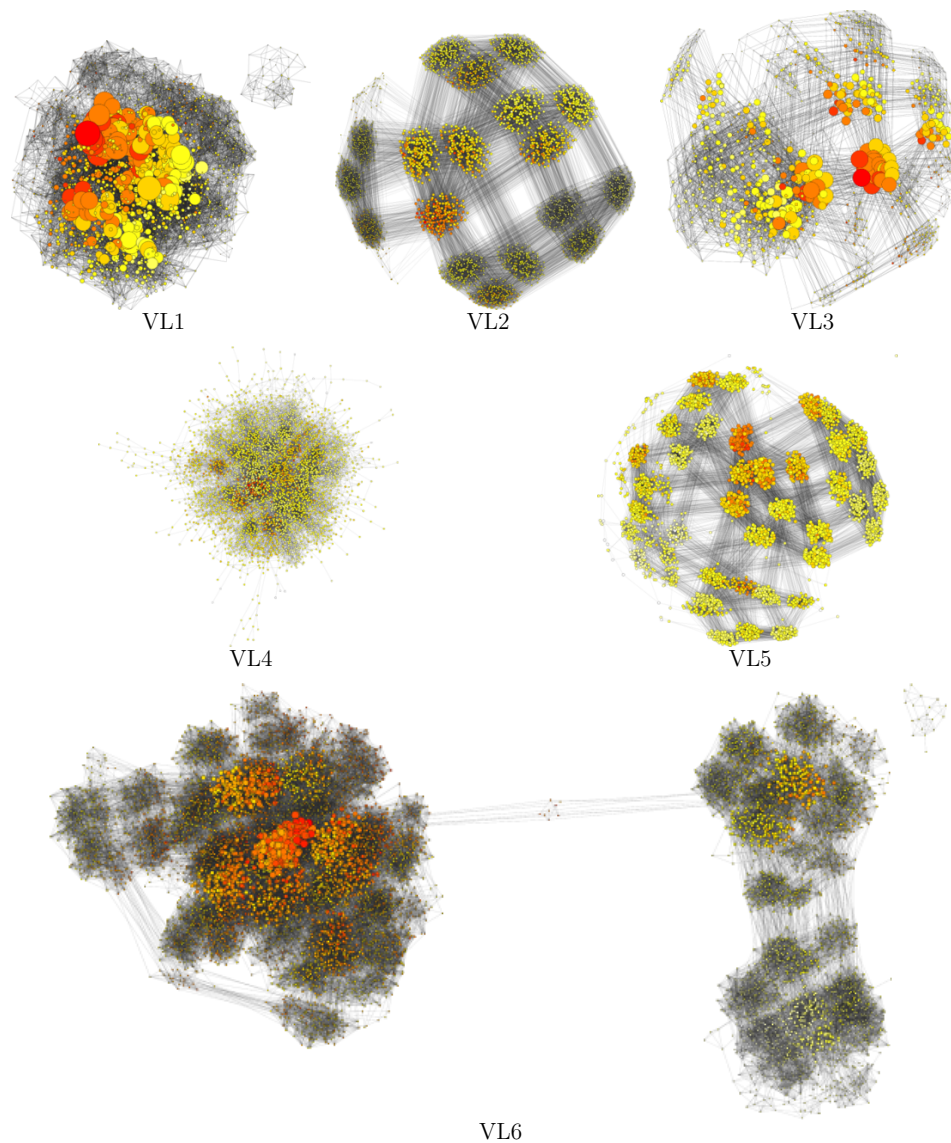*The shape of tree posteriors and identification of peaks*

Distance SPR graphs of the combined golden run tree posteriors from the eukaryote alignments revealed a wide variety of posterior shapes (Supplemental Figure 3). The shapes and complexity of these graphs were clearly not exclusively determined by the number of species or nucleotides in the data set. Topological peaks were evident as large disconnected components (DS1, DS5, DS6) or sets of high probability trees separated by paths of low probability (DS4, DS7). In particular, the trees with highest posterior probability in the two peaks of DS1 were separated by only two

SPR operations but moving between these peaks required leaving the 95% credible set. Large subgraphs of lower probability trees appeared as interesting substructures (e.g. the "tail" on the right hand side of the DS8 graph). No graph could be constructed for DS9 or DS11 as no topology was sampled twice and arbitrary 4096-node subsets were not adjacent in SPR space.
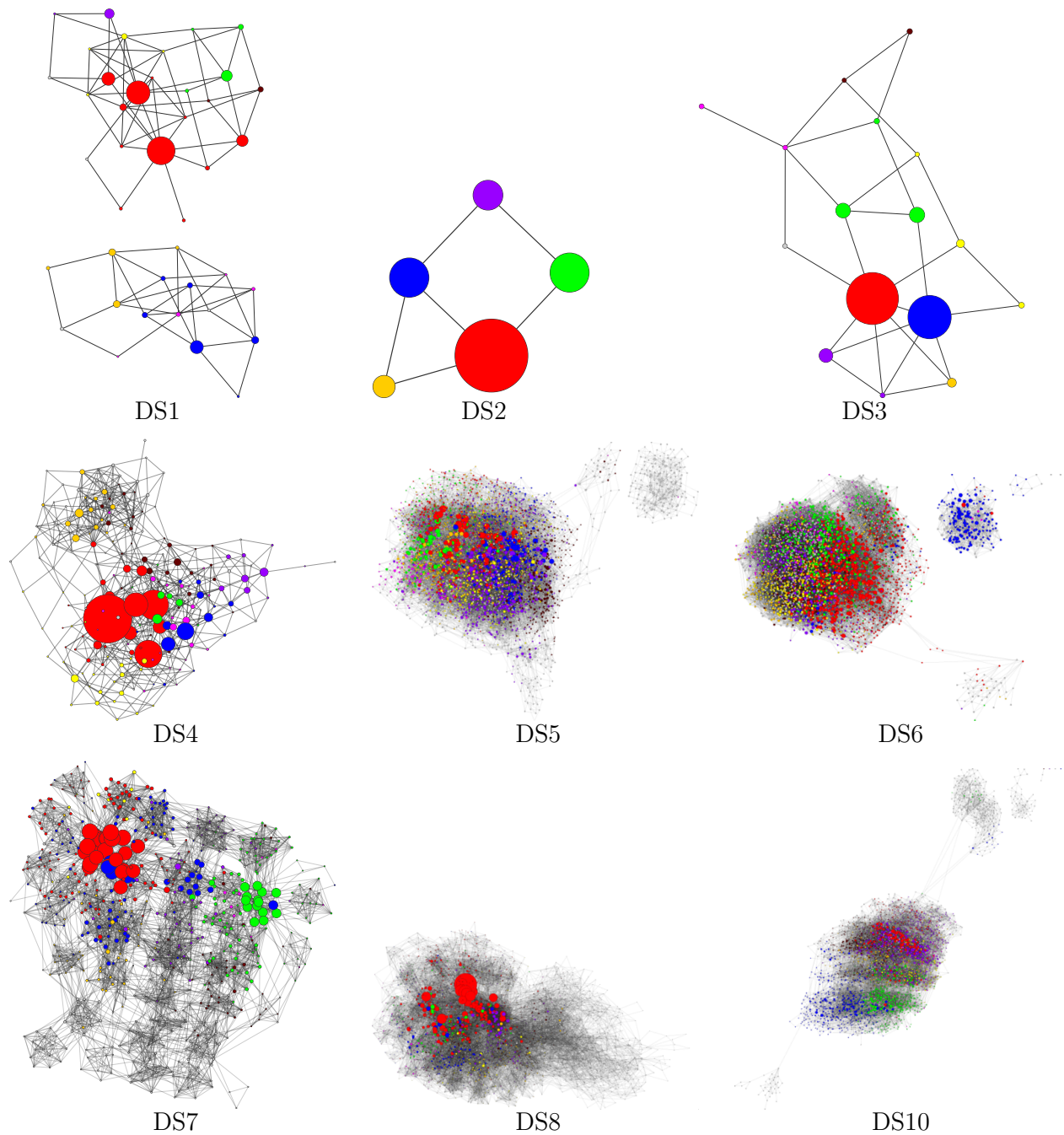
Distance SPR graphs of the combined golden run tree posteriors from the "VL" bacterial and archaeal alignments also showed a wide variety of posterior shapes (Fig. 2). Several posteriors were composed of clumps of trees with similar probability, as in data set DS7, which came from identical or near-identical sequences. These also indicated small changes in uncertain areas of the trees that seldom affect their likelihood but drastically inflate the true 95% credible set of topologies (Supplemental Table 4). We refer to this as the true credible set for brevity. Dataset VL6 provided a striking example of peaks. The 4096 most probable topologies (25.3% credible set) formed 3 disconnected components and the 8192 most probable topologies (31.9% credible set) showed only small paths of connectivity between the 3 peaks. We focused on the eukaryote ("DS") data sets in the remainder of our tests to focus our efforts, unless mentioned otherwise.

Clustering regions of tree space by descending probability (see Methods) highlighted topological peaks and other interesting features (Fig. 3). In addition to the peaky data sets (DS1, DS4, DS5, DS6, DS7) identified with unclustered graphs, DS10 appears to contain at least two peaks. The disconnected sub-peaks of DS1 and DS6 contained the second cluster of both data sets and, thus, the most probable trees outside of the first cluster from each data set. Conversely, the disconnected component of DS5 contained trees of relatively low probability. In non-peaky data sets (e.g. DS3 and DS5) clusters expanded radially from the most probable tree, which indicates relatively easy MCMC mixing.

The number of unique topologies was greatly inflated by ambiguous relationships (Supplemental Table 4). For example, the posterior of data set DS7 had an

**Figure 2.** Distance SPR graphs of the combined bacterial and archaeal golden runs showing at most the 4096 topologies with highest posterior probability (8192 for VL6). Node areas are scaled relative to posterior probability (PP; larger = higher probability) within each graph (but not with respect to the other graphs). Node color indicates SPR distance from the topology with highest posterior probability in each dataset on a red-yellow-white scale (dark-light in the print version), with the highest probability tree colored red.

**Figure 3.** Cluster SPR graphs of the combined golden run eukaryote posteriors. Each graph contains either the 95% credible set or the 4096 topologies with highest PP (DS5, DS6 and DS10). Nodes are scaled relative to posterior probability within each graph (but not with respect to the other graphs). Nodes are colored by SPR-based descending PP clusters (grayscale in the print version).

interesting "grid" structure composed of clumps of 15 trees with similar topology and probability. On closer inspection, trees within a clump differed only in the configuration of a subtree containing four nearly identical *Microcebus rufus* sequences. In fact, these sequences differed in four nucleotides, with one unique mutation per sequence, providing no distinguishing information and inflating the true credible set by a factor of 15 (the number of configurations of four taxa). To verify this effect, we removed three of these four sequences, computed 10 new golden runs, and plotted the resulting tree space (Supplemental Figure 4). As expected we obtained the same structure, but with one tree per 15-node clump and proportional posterior probabilities. The extreme flatness of DS9 and DS11 arose similarly. The majority rule consensus tree for data set DS9 contained two 4-way multifurcations and one 5-way multifurcation. Resolutions of these multifurcations occurred with approximately equal frequency, inflating the true credible set by a factor of $15 * 15 * 105 = 23,625$. Much of the ambiguity was caused by a set of 4 identical sequences and a set of 3 identical sequences. The remaining ambiguity seemed to arise from substantially similar sequences. Similarly, the consensus tree for DS11 contained numerous multifurcations including a multifurcation with 12 edges. The number of samples was insufficient to compare resolutions of this multifurcation and determine if each was equally likely. However, 9 of the taxa involved had the same sequence, which alone inflated the true credible set by at least a factor of $2,027,025$, and this multifurcation likely inflated the credible set by orders of magnitude more. Moreover, the posteriors of data sets DS5, DS6, and DS10 were also inflated by ambiguity. In these cases, none of the sequences involved were identical and resolutions occurred with similar but not equal probability.

The shape of a posterior tree space explains the difficulty of sampling from that distribution (Table 2). Peaky distributions often required a large number of iterations to reach the ASDSF cutoff and/or had high error rates respective to other data sets with a similar number of taxa. In particular, DS1 required the largest number of iterations to

**Table 2.** A comparison of dataset difficulty and posterior shape parameters. The first three columns show the mean number of iterations required to reach ASDSF less than 0.01 ($\mu$Iter) using the MrBayes default parameters (4 runs, 2 chains) as well as the resulting mean maximum split frequency error ($\mu$MaxErr) and mean split frequency RMSD ($\mu$RMSD) as compared to the golden runs. From the golden runs, we considered properties of the *top trees*—the at most 4096 highest probability trees from the 95% credible set. We inferred the SPR radius (Radius) which we define as the maximum SPR distance from each tree in the 95% credible set to the topology with highest posterior probability (radius of the top trees in brackets), the size of the 95% credible set (95CI) , the cumulative posterior probability of the top trees (Cred), and the presence of peaks. Note that our credible set clearly underestimates the true credible set size when it exceeds the number of samples (e.g. DS9 and DS11). "-U" data sets include only one member from each set of identical sequences. Note that each golden run contained 750,000 samples.

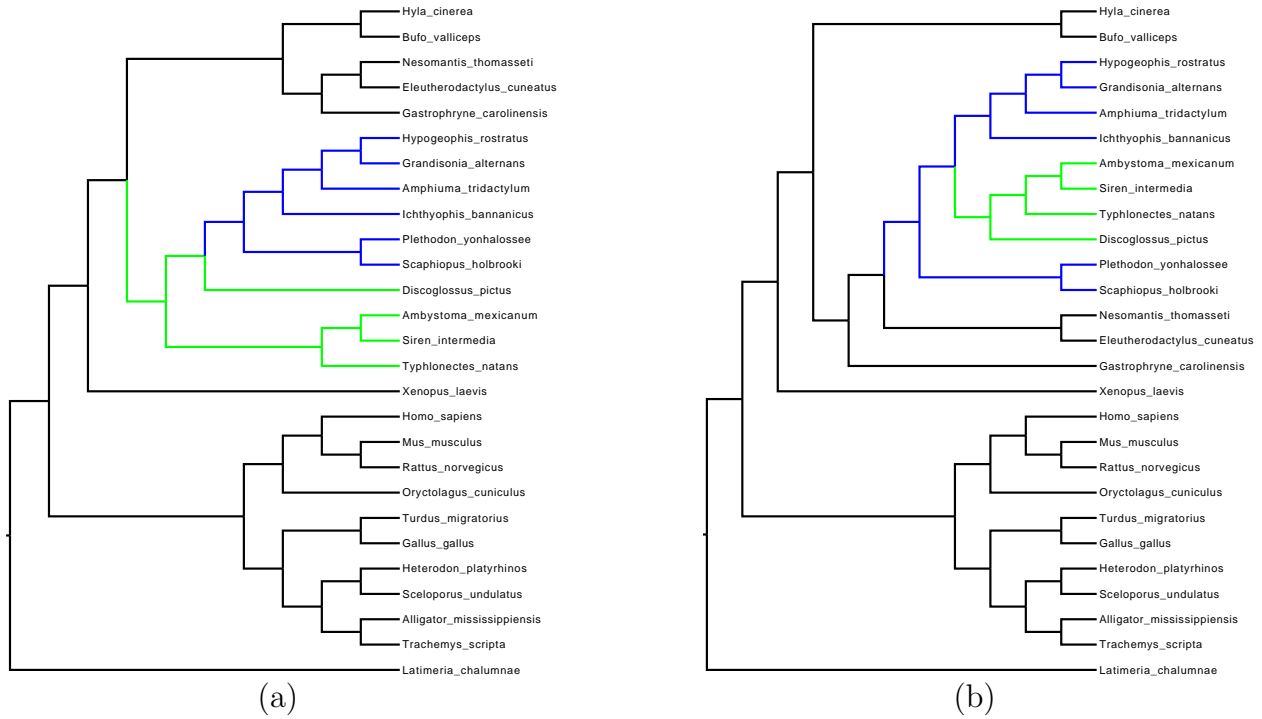| Data | $\mu$Iter | $\mu$MaxErr | $\mu$RMSD | Radius | 95CI | Cred | Peaks |
|---|---|---|---|---|---|---|---|
| DS1 | 850,200 | 0.0819 | 0.0375 | 4 | 41 | 95 | Y |
| DS2 | 8,200 | 0.0976 | 0.0272 | 2 | 5 | 95 | N |
| DS3 | 12,800 | 0.0757 | 0.0225 | 4 | 16 | 95 | N |
| DS4 | 160,800 | 0.1139 | 0.0332 | 6 | 210 | 95 | Y |
| DS5 | 626,000 | 0.0864 | 0.0163 | 16 (8) | 240,311 | 38.9 | Y |
| DS6 | 397,000 | 0.1046 | 0.0244 | 12 (7) | 157,435 | 39.1 | Y |
| DS7 | 62,600 | 0.1616 | 0.0397 | 9 | 735 | 95 | Y |
| DS8 | 283,400 | 0.0882 | 0.0205 | 8 | 3,545 | 95 | N |
| DS9 | 347,200 | 0.1063 | 0.0208 | 23 | 712,502 | 0.6 | ? |
| DS9-U | 255,200 | 0.1019 | 0.0216 | | | | |
| DS10 | 322,400 | 0.1087 | 0.0226 | 15 (12) | 286,604 | 30 | Y |
| DS11 | 338,200 | 0.0503 | 0.0119 | 24 | 712,502 | 0.6 | ? |
| DS11-U | 167,000 | 0.0533 | 0.0143 | | | | |

reach the ASDSF cutoff and had the second highest RMSD of split frequencies despite having the fewest number of species. The number of credible trees and the radius of the tree space also appears to be a factor. DS5 has a large, wide credible set and required a large number of iterations to reach the ASDSF cutoff. DS7 has a smaller credible set and required relatively few iterations for the split frequencies to converge. The high error rates of DS7, however, may indicate that the sub-peak or posterior shape caused the chain to stop prematurely. Despite the large number of taxa and explored topologies of DS9 and DS11, these flat posteriors had low error rates and average times to achieve an ASDSF of 0.01. To remove the effect of identical sequences, we ran 10 new MrBayes replicates of these two data sets with all but one member of each set of identical sequences removed (DS9-U and DS11-U). Removing duplicate sequences reduced the number of iterations required to reach an ASDSF of 0.01 with little effect on error rates as compared to the DS9 and DS11 golden run splits with the corresponding taxa removed.

## *Identifying bottlenecks in tree space*

We were able to explicitly identify bottlenecks in tree space by examining SPR paths between high probability trees separated by regions of low probability. As mentioned above, the most probable topologies of DS1's two peaks are separated by only two SPR operations. However, these SPR operations have an inverted nested structure (Fig. 4). The intermediate topology in this shortest path was so unlikely that it was never sampled in any of our tests. This induces a severe bottleneck that results in the two peaks of DS1. The peaks of DS6 arise from a different type of bottleneck (Supplemental Fig. 5). Three SPR operations are required that move three subtrees into a common clade. Both types of bottleneck are caused by a dependence between splits.

Topological peaks can lead to incorrect estimation of posterior distributions. In addition to long times to achieve small ASDSF and high error rates, there is a risk of missing a peak entirely. This was particularly evident for data set DS1 where 2 of our 10
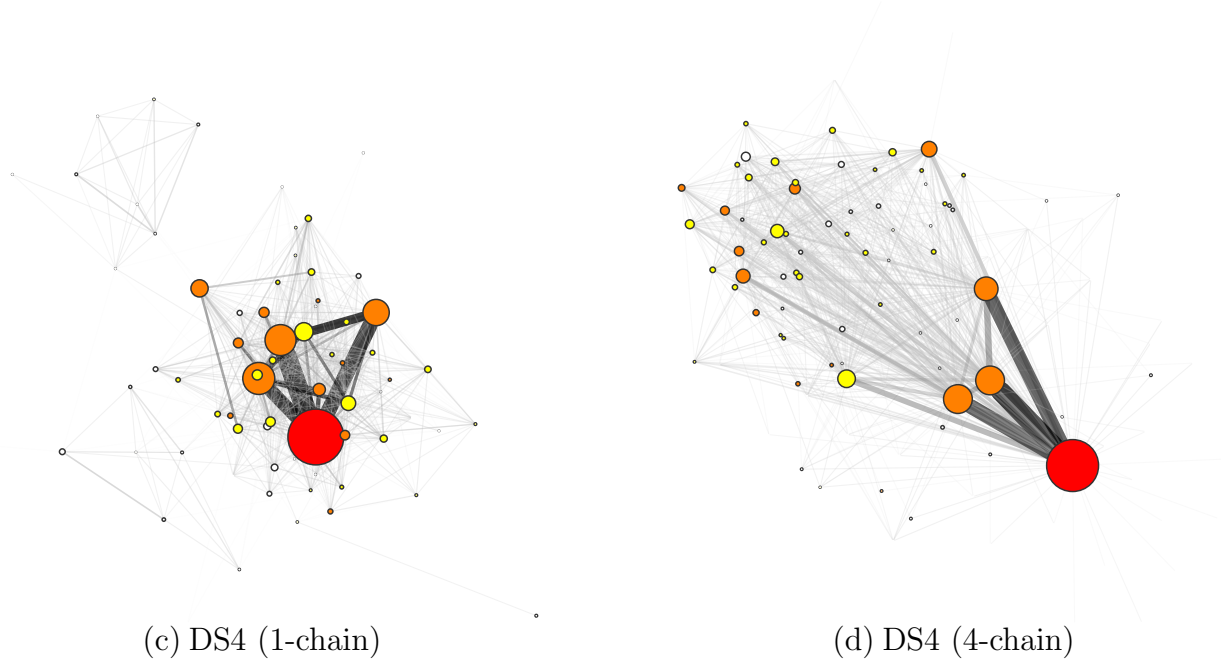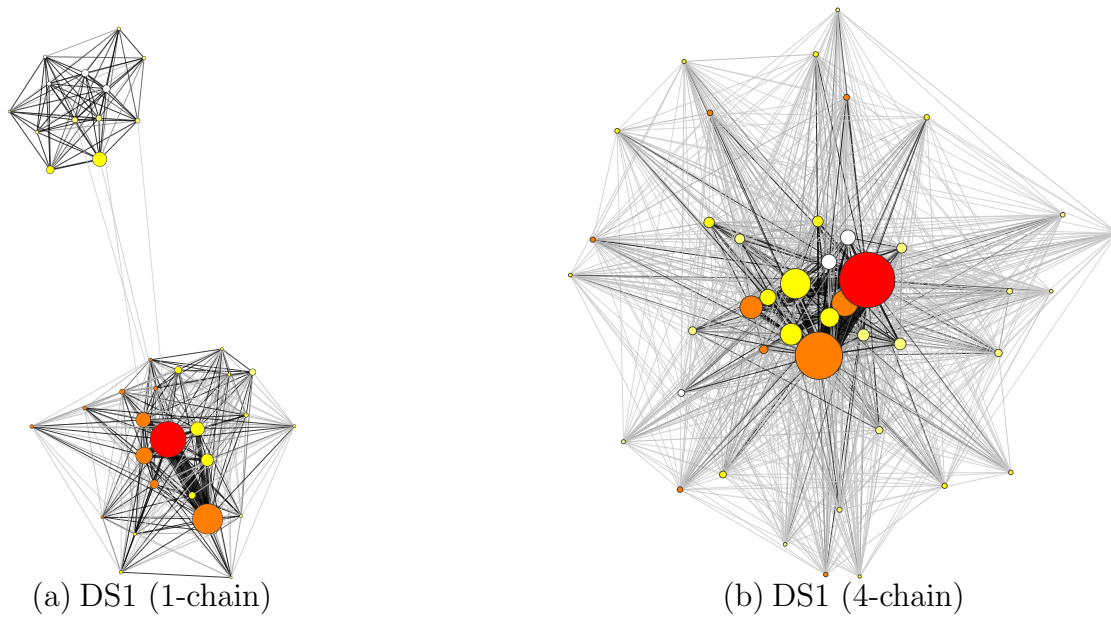
24

**Figure 4.** Central trees of the two topological peaks in dataset DS1. Only two SPR operations separate these trees, moving the blue (gray in the print version) and then green (light gray) clade to traverse from peak 1 to peak 2 and vice versa in the reverse direction. However, the sole intermediate topology is so unlikely that it was never visited in any of our tests, inducing a severe topological bottleneck. Longer paths through multiple trees outside of the 95% confidence interval are taken instead, resulting in long transit times between the peaks.

tests with the MrBayes default settings failed to sample the sub-peak before reaching the ASDSF cutoff. The cumulative posterior probability of this sub-peak (as calculated via golden runs) was approximately 20%. Four splits receive 95-99% support when this sub-peak is missed as opposed to 75-80% support (Supplemental Fig. 6).
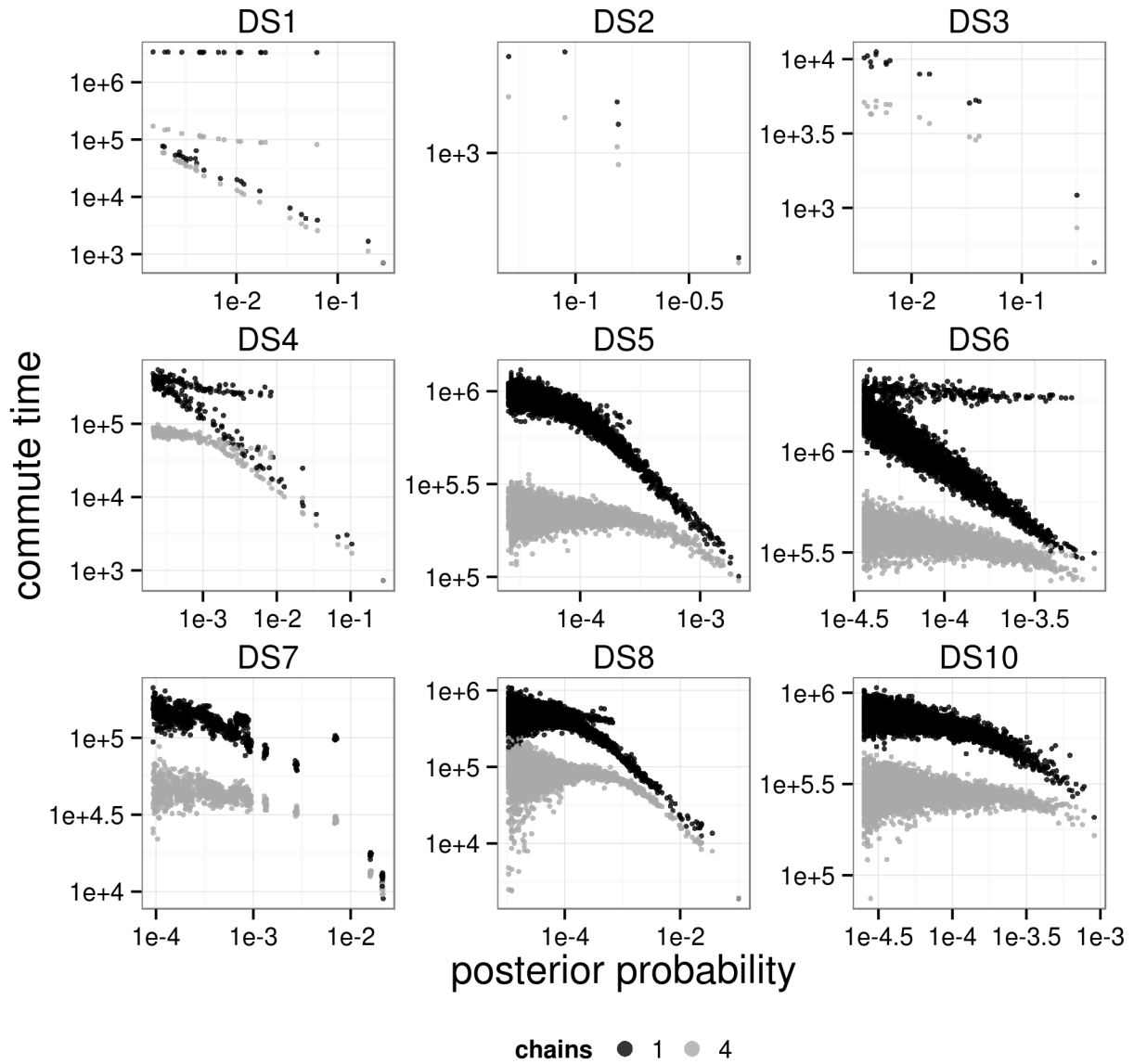
## *Metropolis-coupling improves mixing between peaks*

Metropolis-coupling (Geyer 1992; Huelsenbeck and Ronquist 2001), also known as MCMCMC, connected peaks together for these data sets (Fig. 5). Weighted MCMC graphs of the peaky DS1 for posterior samples without Metropolis-coupling revealed that a single Markov chain rarely transitions between the peaks. For example, there were only 4 observed transitions between peaks in one million-tree sample subsampled from an 100-million iteration MCMC run (Fig. 5(a)). Given the large number of iterations and lack of Metropolis-coupling, it is unlikely that the chain frequently traversed between a peak and returned to the same peak between sampling periods. MCMCMC, however, frequently jumps between the peaks. In one approximately 1.2-million tree sample, subsampled from a 12-million iteration MCMCMC run with 4 chains (Fig. 5(b)), there were more than 4000 observed transitions between the central peak trees. The effect of squashing these graphs together was more pronounced for the deep valley of DS1 as opposed to DS4 (Fig. 5(c)-(d)).

To quantify mixing we computed *commute time* statistics for each topology in the 95% credible set; the commute time here was defined to be the number of Markov chain iterations necessary to move from the highest probability topology to the given tree and back. The *round trip cover time* is the number of iterations necessary to visit every topology in the credible set and return to the highest probability topology. Metropolis-coupling also reduced the mean commute time (Fig. 6) and round trip cover time (Table 3). This effect was particularly pronounced for data set DS1. The round trip cover time decreased by more than a factor of four for DS1, DS4, DS5, DS6, and DS8,

(a) DS1 (1-chain)

(b) DS1 (4-chain)

(c) DS4 (1-chain)

(d) DS4 (4-chain)

**Figure 5.** Weighted MCMC graphs for DS1 and DS4. Node diameters are scaled relative to posterior probability. Nodes are colored on a red-yellow-white scale (dark-light in the print version) with increasing distance from the topology with highest posterior probability. Edges connect trees in successive 100-iteration samples. Edge thickness and color are proportional to the number of MCMC transitions.

**Figure 6.** Comparison of posterior probability and mean commute time with (gray) and without (black) Metropolis-coupling.

**Table 3.** The mean round trip cover time (MRT) for each dataset with and without MCMCMC.

| Data | 1-chain MRT | 4-chain MRT | Ratio |
|------|------------:|------------:|------:|
| DS1 | 3,388,506 | 171,339 | 19.8 |
| DS2 | 2,471 | 1,653 | 1.5 |
| DS3 | 11,182 | 5,246 | 2.1 |
| DS4 | 545,726 | 97,442 | 5.6 |
| DS5 | 2,913,041 | 540,336 | 5.4 |
| DS6 | 9,028,010 | 1,094,217 | 8.3 |
| DS7 | 211,873 | 87,779 | 2.4 |
| DS8 | 1,083,837 | 245,794 | 4.4 |
| DS10 | 2,141,752 | 789,460 | 2.7 |
| All | 19,326,398 | 3,033,266 | 6.4 |

outweighing the factor of four increase in computation time, whereas on data sets DS2, DS3, DS7, and DS10 the improved mixing rate of Metropolis-coupling did not outweigh the increased computation. However, Metropolis-coupling reduced total computation time substantially, as the data sets where it did not reduce total computational effort to achieve a fixed ASDSF mixed relatively quickly compared to the ones for which it did. Commute and cover time statistics could not be estimated for the flat DS9 and DS11 posteriors. These results suggest that Metropolis-coupling does improve mixing between peaks and reduce total computational effort on average, but may not be beneficial for all posterior shapes.

Trees within sub-peaks were observed to have much larger commute times than other trees with a similar posterior probability. This effect was particularly prominent in data sets DS1, DS4, DS6, and DS7 (Fig. 6). For example, the commute time of the central tree in the sub-peak of DS1 was 2.6 million iterations as opposed to 3,300-5,500 iterations

for trees with similar probability. This reduced to 80,200 and 2,100-3,700 iterations, respectively with Metropolis-coupling. Similarly, the most probable trees within the three sub-peaks of DS4 (Fig. 3) had commute times between 200,000-307,000 iterations (37,000-47,000 with Metropolis-coupling). Other trees of similar probability had commute times between 16,000-25,000 iterations. Generally, our commute time analysis further demonstrates the difficulty of sampling sub-peaks and allows quantification of this difficulty.

*A small ASDSF when calculated with two runs is not always sufficient to ensure that empirical split frequencies are close to their posterior distribution; Metropolis-coupling aids in split frequency mixing*

As described in Methods, ASDSF compares split frequencies between runs (we emphasize that these runs are completely distinct and not coupled as for MCMCMC). Increasing the number of simultaneous Markov chain runs greatly increased the stringency of a given ASDSF cutoff (Supplemental Figure 7). We found that ASDSF calculated using two runs is not sufficient for estimating the convergence of split frequencies. Adding additional runs both increased the number of iterations required to reach the ASDSF cutoff and decreased the amount of error. This effect varied by data set and peaky distributions saw the greatest decrease in error with additional runs.

In most cases, a small ASDSF implied that other convergence diagnostics were satisfied, regardless of the number of runs. The mean potential scale reduction factor (PSRF; see Methods) for branch lengths was less than 1.01 in all but the 2-run DS2 and DS3 cases and 4 2-run DS7 cases, where the mean PSRF was less than 1.042. Similarly, the ESS for the treelength parameter was greater than 200 except for data sets DS2 and DS3 and 8 of the 4-chain 2-run DS7 cases.

The ASDSF and split frequency error varied considerably over Markov chains of peaky data sets as runs transitioned between peaks. These statistics often dipped below

commonly applied thresholds only to increase rapidly after one run began exploring an alternative peak. The subsequent rise and fall of these statistics decreased in magnitude as we gathered a sufficient sample. However, current convergence diagnostics assume that these statistics decrease smoothly and, in particular, do not rise sharply. The first time that an ASDSF cutoff is reached may not indicate that split frequency estimates are close to their posterior probabilities in peaky posteriors. Moreover, the stopping time for Markov chains is often determined by the first occurrence of a sufficiently small split frequency deviation from golden runs (Höhna and Drummond 2012). This approach may underestimate the time needed to run these chains in the presence of topological peaks because the running observation of the split frequency may get close to the golden run split frequency because of stochasticity.

Metropolis-coupling decreased error when peaky distributions were sampled with a small number of runs. Dataset DS1, in particular, required 7 or more runs to achieve a mean RMSD below 0.02 without Metropolis-coupling but only 3 runs with Metropolis-coupling (Table 4). Much of this error occurred when runs prematurely stopped using a split frequency criterion on the larger peak. Even with 8 runs, one replicate without Metropolis-coupling reached an ASDSF of 0.01 after only 740,000 iterations, compared to the mean 82 million iterations. None of the 8 runs visited any tree in the sub-peak, resulting in an RMSD of 0.08 and similar effects to those detailed above (Supplementary Fig. 6). The common diagnostics were satisfied for this replicate, including an ASDSF value less than 0.01, tree length ESS value of 3054, tree length PSRF of 1.000 and a maximum split frequency PSRF of 1.001. Even with Metropolis-coupling, the MrBayes default of two runs was insufficient to adequately sample data sets DS1, DS4, and DS7 at the 0.01 ASDSF threshold.

*Topological Gelman-Rubin-like statistic*

Because split frequency is a projection of the actual posterior on phylogenetic

31

**Table 4.** A detailed look at performance on dataset DS1 with and without MCMCMC using varying number of runs. The number of replicates (out of 10) with a given number of chains (Ch) are shown which found both peaks (Peak), converged to an RMSD at most 0.02 (Conv), or exceeded the iteration limit (Lim). The mean number of iterations, running time (iterations*chains*runs), maximum split frequency error (MaxErr), and RMSD are also shown.

| Ch | Runs | Peak | Conv | Lim | Iterations | Run Time | MaxErr | RMSD |
|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 0 | 0 | 0 | 2.42E+6 | 4.84E+6 | 0.266 | 0.106 |
| 1 | 3 | 6 | 3 | 0 | 2.06E+7 | 6.18E+7 | 0.108 | 0.045 |
| 1 | 4 | 8 | 3 | 1 | 3.39E+7 | 1.36E+8 | 0.081 | 0.036 |
| 1 | 5 | 9 | 7 | 3 | 5.88E+7 | 2.94E+8 | 0.052 | 0.024 |
| 1 | 6 | 9 | 6 | 3 | 6.64E+7 | 3.98E+8 | 0.053 | 0.024 |
| 1 | 7 | 9 | 8 | 4 | 7.52E+7 | 5.26E+8 | 0.041 | 0.018 |
| 1 | 8 | 9 | 8 | 6 | 8.26E+7 | 6.61E+8 | 0.037 | 0.016 |
| 4 | 2 | 8 | 5 | 0 | 8.50E+5 | 6.80E+6 | 0.082 | 0.038 |
| 4 | 3 | 10 | 6 | 0 | 4.05E+6 | 4.86E+7 | 0.038 | 0.018 |
| 4 | 4 | 10 | 6 | 0 | 4.07E+6 | 6.51E+7 | 0.030 | 0.015 |
| 4 | 5 | 10 | 8 | 0 | 6.52E+6 | 1.30E+8 | 0.026 | 0.013 |
| 4 | 6 | 10 | 9 | 0 | 1.20E+7 | 2.88E+8 | 0.021 | 0.010 |
| 4 | 7 | 10 | 10 | 0 | 1.36E+7 | 3.81E+8 | 0.013 | 0.006 |
| 4 | 8 | 10 | 10 | 0 | 1.38E+7 | 4.42E+8 | 0.011 | 0.005 |

**Table 5.** Estimated topology deviation (RMSD) and potential scale reduction factor (PSRF) using our topological Gelman-Rubin-like measure (TGR). Ch = number of chains.

| Data | Ch | TGR-RMSD | | TGR-PSRF | |
|------|-----|--------|--------|--------|--------|
| | | 2-runs | 8-runs | 2-runs | 8-runs |
| DS1 | 1 | 1.9 | 2.2 | 1.001 | 1.002 |
| DS1 | 4 | 1.9 | 2.2 | 1.001 | 1.001 |
| DS2 | 1 | 0.9 | 1.1 | 1.005 | 1.005 |
| DS2 | 4 | 0.9 | 1.1 | 1.007 | 1.005 |
| DS3 | 1 | 0.9 | 1.2 | 1.004 | 1.007 |
| DS3 | 4 | 1.0 | 1.2 | 1.004 | 1.004 |
| DS4 | 1 | 1.8 | 2.3 | 1.002 | 1.002 |
| DS4 | 4 | 2.0 | 2.3 | 1.001 | 1.001 |
| DS5 | 1 | 5.7 | 6.4 | 1.000 | 1.001 |
| DS5 | 4 | 6.3 | 6.3 | 1.000 | 1.001 |
| DS6 | 1 | 5.6 | 6.0 | 1.001 | 1.001 |
| DS6 | 4 | 6.0 | 6.0 | 1.001 | 1.001 |
| DS7 | 1 | 2.5 | 3.3 | 1.002 | 1.005 |
| DS7 | 4 | 3.1 | 3.3 | 1.002 | 1.004 |
| DS8 | 1 | 3.1 | 3.3 | 1.001 | 1.002 |
| DS8 | 4 | 3.3 | 3.4 | 1.002 | 1.001 |
| DS9 | 1 | 16.3 | 16.4 | 1.000 | 1.001 |
| DS9 | 4 | 16.2 | 16.4 | 1.000 | 1.001 |
| DS10 | 1 | 6.8 | 5.8 | 1.001 | 1.001 |
| DS10 | 4 | 7.6 | 6.2 | 1.001 | 1.001 |
| DS11 | 1 | 18.5 | 18.6 | 1.001 | 1.001 |
| DS11 | 4 | 18.5 | 18.6 | 1.000 | 1.001 |

trees rather than the posterior itself, we wondered to what extent split-based measures being small implies that the empirical frequency on phylogenetic tree topologies is close to the posterior. To explore this question, we developed a variant of the Gelman-Rubin statistic that used SPR distances (Methods). This measure compares the mean square topology deviation within independent Markov chains to that between the chains. The corresponding PSRF will approach 1 as the independent runs converge in topology distribution.

On our data sets, a small ASDSF generally implied that the topological measure was small (Table 5). PSRF estimates with our topological Gelman-Rubin-like measure approached 1, regardless of the number of runs. Surprisingly, this also held for the flat posteriors of DS9 and DS11. This suggests that similar trees were explored between runs of these posteriors, even if no two trees were identical. There was little difference in topology deviation or PSRF with or without Metropolis-coupling. Moreover, topological PSRF and ASDSF showed similar trends over time (Supplemental Fig. 9), although the scale of this relationship appears to vary between different data sets and even different replicated tests on the same data set.

## *Multidimensional scaling*

In general, MDS projections were insufficient to diagnose peaks (Fig. 7) and extra information is required such as commute time, posterior density, and connectivity. For flat posteriors, however, where extra information is unavailable, multidimensional scaling remains the only method of visualizing tree space (Supplemental Fig. 8).

Specifically, MDS plots often highlighted topological differences that did not impede mixing and missed sub-peaks. For data set DS1, MDS displayed 4 clusters. One axis separated the two peaks of DS1, but the other axis separated trees according to a common difference that did not impede mixing. MDS plots of DS4 identified only one of three difficult to reach areas of tree space. MDS plots were generally similar between RF

(a) DS1 (SPR)

(b) DS1 (RF)

(c) DS4 (SPR)

(d) DS4 (RF)

(e) DS6 (SPR)

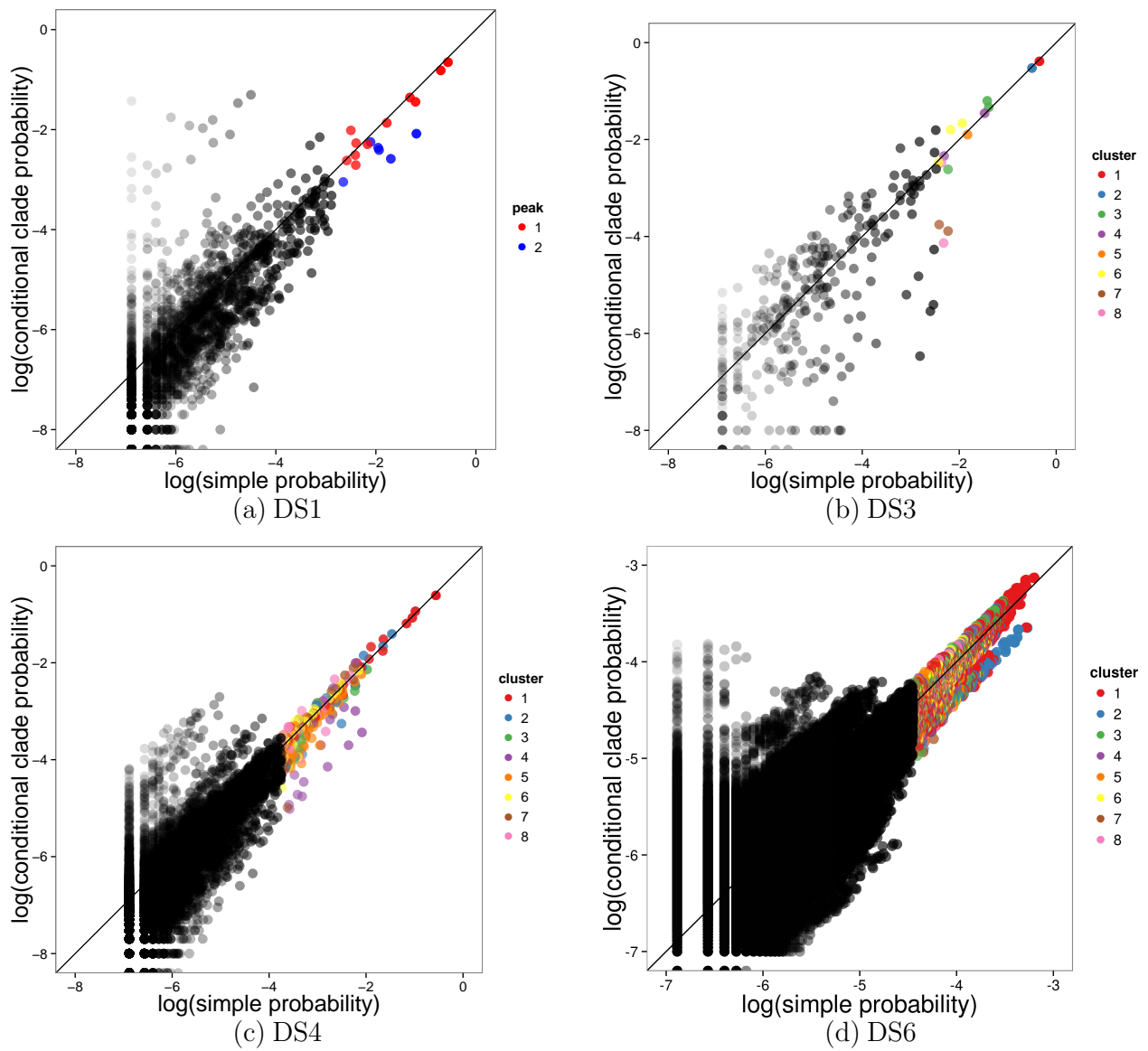(f) DS6 (RF)

(g) VL6 (SPR)

(h) VL6 (RF)

**Figure 7.** Comparison of multidimensional scaling representations with SPR and RF distances. Nodes are colored by identified peaks to match the primary cluster of the peak (grayscale in the print version).

and SPR. In DS6, however, the plots differed significantly. This highlighted the fact that the peaks of DS6 are quite close in SPR terms (despite their separation by a valley of improbable trees) but had very different splits. We also compared MDS plots on the peaky microbial data set VL6. Only one of the peaks was separated from the others in the SPR plot and the RF plot broke the peaks into multiple clumps.

## *Effect of peaks on conditional clade probabilities*

Recent work uses a product of conditional posterior probabilities on splits as a proxy for the corresponding phylogenetic posterior probability (Höhna and Drummond 2012; Larget 2013; Szöllősi et al. 2013). This assumes an independence between the split probabilities of sister clades conditioned on their parent clade. Larget (2013) found several examples of trees where CCD probabilities differ from well-sampled empirical frequencies in the eukaryote datasets and where the simple estimates were well above the sampling threshold. Larget (2013) conducted two tests per eukaryote dataset using MrBayes 3.2 with the GTR model for 5,500,000 iterations, subsampling 100,000 trees with a 500,000-tree burn-in period. He used these long runs as being representative of the posterior distribution, noting that "A second set of runs under the same conditions, but with different random numbers, shows very similar results, indicating that these MCMC samples are likely not to suffer from poor convergence (data not shown)." We extended his investigation of differences between CCD probabilities and well-sampled empirical frequencies with our substantially larger 1-billion iteration MCMC golden runs, subsampling 750,000 trees with a 250,000,000-tree burn-in period, replicated 10 times.

We found that conditional independence clearly did not hold in peaky distributions (Fig. 8). Specifically, conditional clade probabilities systematically underestimated the probability of trees within sub-peaks and overestimated the probability of trees between the peaks. This effect was exemplified in data set DS1 where highly unlikely trees between peaks had conditional clade probabilities exceeding one

**Figure 8.** A comparison of posterior probability and CCD estimates for the aggregated golden runs on datasets DS1, DS3, DS4, and DS6. Probability is shown on a log-log scale in base 10. The top trees for each dataset are colored by peak in DS1 and cluster for the other datasets. Transparency of points increases as posterior probability decreases.

percentage point (points significantly above the line in Figure 8(a)). We observed similar effects in DS4 and DS6. Surprisingly, even in the relatively simple posterior of DS3, CCD underestimated the posterior probability of three trees in the 95% credible set by an order of magnitude. However, CCD performed well overall on non-peaky data sets and is currently the only way to estimate probabilities below the sampling threshold (e.g. DS9 and DS11).

# DISCUSSION

We developed the first practical method for examining the subtree prune-and-regraft tree space of Bayesian phylogenetic posteriors. Our novel graph-based approach uses size and color to visualize connectivity, posterior probability, and relative distance. Our simple clustering procedure identified topological peaks in several real data sets. Additionally, we investigated the impact of Metropolis-coupling, the number of runs used for ASDSF calculation, and developed a convergence diagnostic that uses phylogenetic tree topologies directly.

We find that multimodal or "peaky" posteriors are common in data sets with 30 or more taxa, confirming the suggestion by Beiko et al. (2006). Markov chains on peaky posteriors often required a large number of iterations to obtain small ASDSF values and had high error rates relative to the number of taxa. We used dynamic programming to compare tree commute times and found that trees within sub-peaks were difficult to sample. The "height" of a peak compared to the "depth" of the corresponding valley influenced sampling difficulty. Dataset DS1, despite its relatively small number of taxa, has a large sub-peak separated by a particularly deep valley. In many cases, this led to premature termination of chains by the ASDSF measure and erroneously assigning greater than 95% confidence to some relationships with an actual frequency less than 80%.

We explicitly identified tree space bottlenecks in two data sets with tall sub-peaks and found that they were caused by a dependence between splits. These peaks were only isolated by a handful of SPR operations. However, the intermediate valley topologies were exceedingly unlikely and the SPR operations modified a large number of splits.

Dependence between sister clades caused systematic errors in CCD probability estimates. Specifically, CCD overestimated the probability of trees between peaks and underestimated the probability of trees in sub-peaks. These observations suggest that CCD-guided proposal operators could hide sub-peaks and further aggravate the difficulty of sampling peaky phylogenetic posteriors. On the other hand, CCD-guided proposal operators may sample valley trees more frequently and therefore provide more chances to cross valleys and sample sub-peaks. Tree space sampling methods that penalize or even prevent SPR and TBR operators that change a large number of splits could also hide sub-peaks, such as the "pruning distance" of Höhna and Drummond (2012) and similar suggestions (Huelsenbeck et al. 2008; Lakner et al. 2008). Moreover, an anti-peak bias would be undetectable and, perversely, decrease running times using an ASDSF rule or other common convergence diagnostic. One strategy to alleviate bias, while still retaining the benefit of CCD, might be to use CCD or other biased proposal distributions in a subset of Markov chains (Metropolis-coupled or otherwise) along with chains using a general proposal distribution. CCD has also begun to see use in phylogenomic methods such as amalgamated likelihood estimation (Szöllősi et al. 2013), which uses CCD directly as a proxy for posterior probability in order to infer a species tree joint with a set of gene trees. Biases in CCD will bias the results of this approach.

Identical and closely related sequences cause ambiguity which greatly inflated the tree space sampled by MCMC methods. Such data sets had large and flat posteriors, which were difficult to quantify and visualize. Ignoring duplicate sequences reduced mean run time (determined by the ASDSF stopping rule) by 26% and 50% in our two flattest posteriors. Thus, we suggest that users of MCMC should identify and ignore

duplicate sequences, maintaining only a single representative from each set of identical sequences. The post-processing of such an analysis could either expand the representative to a multifurcating clade containing each ignored sequence from a set, or spread the probability of each sampled tree uniformly across each full tree with monophyletic clades for the expanded sets. Ignoring duplicate sequences may make branch length priors harder to interpret, due to a consequent ascertainment bias. This may make little difference in practice, however, as commonly used priors do not allow for large numbers of simultaneous or nearly simultaneous branching events. Intelligently handling duplicate sequences may be a useful feature of future MCMC software. Moreover, future work should explore methods for handling closely related sequences without inflating tree space. This could be done with reversible jump Markov chain Monte Carlo (Lewis et al. 2005). Tree space inflation will be of particular importance when estimating trees for a large number of closely related sequences as in personalized medicine and metagenomics.

Metropolis-coupling was effective in reducing commute times and decreased the mean cover time of the 95% credible set by more than a factor of 4 in peaky distributions. Metropolis-coupling increased the number of transitions between peaks by three orders of magnitude in our peakiest data set. However, Metropolis-coupling may not be effective for all posterior shapes. The observed cover time decrease in non-peaky data sets did not outweigh the increased computation of Metropolis-coupling, however because Metropolis-coupling significantly reduced computational load for the most difficult and time-consuming posteriors, it appears to be a useful default option on average. We tested the effect of Metropolis-coupling with 4 chains, the default number of chains in MrBayes, and future work should investigate the optimal number of Metropolis-coupled chains for peaky and non-peaky posteriors. Moreover, further research could investigate whether it is heating, multiple chains, or both that improves mixing in peaky posteriors.

The magnitude of the ASDSF convergence diagnostic depends heavily on the number of Markov chains used for comparison. We found that using 2 independent runs with an ASDSF cutoff of 0.01 resulted in insufficient chain lengths for peaky posterior distributions. Indeed, MCMC runs often stopped using a 2-run ASDSF stopping rule before sampling a sub-peak. This is a serious consideration, as MrBayes uses 2 runs by default, and MrBayes uses a default ASDSF termination threshold of 0.05 when ASDSF termination is enabled but no threshold provided. Moreover, MrBayes does not provide a warning message unless the ASDSF exceeds 0.1 when run for a fixed number of iterations. ExaBayes uses ASDSF termination by default with a threshold of 0.05. We did not test similar single-chain convergence diagnostics (e.g the methods of Raftery and Lewis (1992) or Geweke (1991)) but they may experience similar problems. MCMC analyses should use at least 3 independent runs and an ASDSF threshold of at least 0.01 in any MCMC analysis for which accurate topological posterior estimation is an important concern. Moreover, multiple independent MCMC replicates should be compared—using even 8 runs was not enough to prevent one of our MCMC tests from stopping with an ASDSF stopping rule before sampling the sub-peak of DS1. A wide variance in chain lengths using split frequency stopping rules on independent replicates may be a sign of topological sub-peaks.

We developed a topological Gelman-Rubin-like convergence diagnostic which works directly on tree topologies. This diagnostic can be applied with any distance metric on tree topologies. Tests with this topological Gelman-Rubin-like measure suggest that small ASDSF often implies a small topological Gelman-Rubin-like diagnostic for high-probability topologies, although neither measure can detect unsampled topological peaks.

A major and natural difficulty of peak detection is that the peaks must be sampled in order to be detected. Similarly, it is difficult to accurately estimate time to satisfy some convergence criterion. Convergence time estimates using golden runs

(Höhna and Drummond 2012) are based on the first time that split frequencies of a Markov chain approach the golden run split distribution. However, this approach may underestimate the running time of MCMC methods in practice because sampled split frequencies can approximately hit the golden run split distribution before they have stabilized. It may be worth checking that split distributions have stabilized in addition to requiring them to hit the golden run split distribution.

Our methods could be expanded in several ways. We limited many of our comparisons to subsets of at most 4096 trees due to the computational overhead of pairwise comparisons. Our approach would benefit greatly from faster methods for unrooted SPR comparisons or a way to construct SPR graphs without comparing every pair of trees. There also are special challenges in moving through the space of rooted trees with a time component (as estimated by BEAST), which would be interesting to investigate; our methods would also be much more efficient on posteriors of these rooted trees. We developed a very simple method for highlighting topological peaks that was designed to dynamically select cluster radii with few SPR comparisons. Our clustering procedure worked well in our tests, but in multiple situations could select unreasonably small cluster sizes (e.g. if the standard deviation approached or exceeded the mean). Improved methods for identifying such peaks and analyzing tree space graphs should be explored. In particular, methods are needed to rapidly scan posteriors for common bottlenecks in order to develop new phylogenetic operators that cross those bottlenecks. Moreover, future work should determine the cause of such bottlenecks in terms of sequence features (for example mixtures of tree topologies). Nontrivial methods will be required to do so in the likelihood-based framework. Finally, our observations need to be confirmed on other data sets. This work is but a first step in quantifying MCMC exploration of phylogenetic tree space using topological methods.

# FUNDING

## ACKNOWLEDGEMENTS

## REFERENCES

Aldous, D. J. 2000. Mixing time for a Markov chain on cladograms. Combin. Probab. Comput. 9:191–204.

Allen, B. L. and M. Steel. 2001. Subtree transfer operations and their induced metrics on evolutionary trees. Ann. Comb. 5:1–15.

Beiko, R. G., J. M. Keith, T. J. Harlow, and M. A. Ragan. 2006. Searching for convergence in phylogenetic markov chain monte carlo. Syst. Biol. 55:553–565.

Bordewich, M. and C. Semple. 2005. On the computational complexity of the rooted subtree prune and regraft distance. Ann. Comb. 8:409–423.

Bouckaert, R., J. Heled, D. Kühnert, T. Vaughan, C.-H. Wu, D. Xie, M. A. Suchard, A. Rambaut, and A. J. Drummond. 2014. Beast 2: a software platform for bayesian evolutionary analysis. PLoS computational biology 10:e1003537.

Drummond, A. J. and A. Rambaut. 2007. BEAST: Bayesian evolutionary analysis by sampling trees. BMC Evol. Biol. 7:214.

Drummond, A. J., M. A. Suchard, D. Xie, and A. Rambaut. 2012. Bayesian phylogenetics with beauti and the beast 1.7. Molecular biology and evolution 29:1969–1973.

Garey, J. R., T. J. Near, M. R. Nonnemacher, and S. A. Nadler. 1996. Molecular evidence for acanthocephala as a subtaxon of rotifera. Journal of Molecular Evolution 43:287–292.

Gelman, A. and D. B. Rubin. 1992. Inference from iterative simulation using multiple sequences. Statist. Sci. 7:457–472.

Geweke, J. 1991. Evaluating the accuracy of sampling-based approaches to the calculation of posterior moments. Pages 169–193 *in* Bayesian Statistics 4 Oxford University Press.

Geyer, C. J. 1992. Practical Markov chain Monte Carlo. Statist. Sci. 7:473–483.

Hastings, W. K. 1970. Monte Carlo sampling methods using Markov chains and their applications. Biometrika 57:97–109.

Hedges, S. B., K. D. Moberg, and L. R. Maxson. 1990. Tetrapod phylogeny inferred from 18s and 28s ribosomal rna sequences and a review of the evidence for amniote relationships. Molecular Biology and Evolution 7:607–633.

Hein, J., T. Jiang, L. Wang, and K. Zhang. 1996. On the complexity of comparing evolutionary trees. Discrete Appl. Math. 71:153–169.

Henk, D. A., A. Weir, and M. Blackwell. 2003. Laboulbeniopsis termitarius, an ectoparasite of termites newly recognized as a member of the laboulbeniomycetes. Mycologia 95:561–564.

Hickey, G., F. Dehne, A. Rau-Chaplin, and C. Blouin. 2008. SPR distance computation for unrooted trees. Evolutionary Bioinformatics 4:17–27.

Hillis, D. M., T. A. Heath, and K. S. John. 2005. Analysis and visualization of tree space. Syst. Biol. 54:471–482.

Höhna, S., M. Defoin-Platel, and A. J. Drummond. 2008. Clock-constrained tree proposal operators in Bayesian phylogenetic inference. Pages 1–7 *in* BioInformatics and BioEngineering, 2008. BIBE 2008. 8th IEEE International Conference on IEEE.

Höhna, S. and A. J. Drummond. 2012. Guided tree topology proposals for bayesian phylogenetic inference. Syst. Biol. 61:1–11.

Huelsenbeck, J. P., C. Ané, B. Larget, and F. Ronquist. 2008. A bayesian perspective on a non-parsimonious parsimony model. Syst. Biol. 57:406–419.

Huelsenbeck, J. P. and F. Ronquist. 2001. Mrbayes: Bayesian inference of phylogenetic trees. Bioinformatics 17:754–755.

Ingram, A. L. and J. J. Doyle. 2004. Is eragrostis (poaceae) monophyletic? insights from nuclear and plastid sequence data. Systematic botany 29:545–552.

Kroken, S. and J. W. Taylor. 2000. Phylogenetic species, reproductive mode, and specificity of the green alga trebouxia forming lichens with the fungal genus letharia. The Bryologist 103:645–660.

Kruskal, J. B. 1964a. Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. Psychometrika 29:1–27.

Kruskal, J. B. 1964b. Nonmetric multidimensional scaling: a numerical method. Psychometrika 29:115–129.

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. Syst. Biol. 57:86–103.

Larget, B. 2013. The estimation of tree posterior probabilities using conditional clade probability distributions. Syst. Biol. 62:501–511.

Lewis, P. O., M. T. Holder, and K. E. Holsinger. 2005. Polytomies and Bayesian phylogenetic inference. Syst. Biol. 54:241–253.

Lovász, L. 1993. Random walks on graphs: a survey. Combinatorics, Paul Erdös is Eighty 2:1–46.

Matsen, F. A. 2006. A geometric approach to tree shape statistics. Syst. Biol. 55:652–661.

Metropolis, N., A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller. 1953. Equation of state calculations by fast computing machines. The Journal of Chemical Physics 21:1087.

Mossel, E. and E. Vigoda. 2005. Phylogenetic MCMC algorithms are misleading on mixtures of trees. Science 309:2207–2209.

Mossel, E. and E. Vigoda. 2006. Limitations of Markov chain Monte Carlo algorithms for Bayesian inference of phylogeny. Ann. Appl. Probab. 16:2215–2234.

Raftery, A. E. and S. Lewis. 1992. How many iterations in the gibbs sampler. Bayesian Statistics 4:763–773.

Robinson, D. F. and L. R. Foulds. 1981. Comparison of phylogenetic trees. Math. Biosci. 53:131–147.

Ronquist, F., J. Huelsenbeck, and M. Teslenko. 2011. Draft MrBayes version 3.2 manual: tutorials and model summaries. Distributed with the software from `http://brahms.biology.rochester.edu/software.html`.

Ronquist, F., B. Larget, J. P. Huelsenbeck, J. B. Kadane, D. Simon, and P. van der Mark. 2006. Comment on "phylogenetic MCMC algorithms are misleading on mixtures of trees". Science 312:367–367.

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. MrBayes 3.2: efficient bayesian phylogenetic inference and model choice across a large model space. Syst. Biol. 61:539–542.

Rossman, A. Y., J. M. McKemy, R. A. Pardo-Schultheiss, and H.-J. Schroers. 2001. Molecular studies of the bionectriaceae using large subunit rDNA sequences. Mycologia 93:100–110.

Shannon, P., A. Markiel, O. Ozier, N. S. Baliga, J. T. Wang, D. Ramage, N. Amin, B. Schwikowski, and T. Ideker. 2003. Cytoscape: a software environment for integrated models of biomolecular interaction networks. Genome Res. 13:2498–2504.

Štefankovic, D. and E. Vigoda. 2011. Fast convergence of Markov chain Monte Carlo algorithms for phylogenetic reconstruction with homogeneous data on closely related species. SIAM J. Discrete Math. 25:1194–1211.

Suchard, M. A. and B. D. Redelings. 2006. Bali-phy: simultaneous bayesian inference of alignment and phylogeny. Bioinformatics 22:2047–2048.

Suh, S.-O. and M. Blackwell. 1999. Molecular phylogeny of the cleistothecial fungi placed in cephalothecaceae and pseudeurotiaceae. Mycologia 91:836–848.

Szöllősi, G. J., W. Rosikiewicz, B. Boussau, E. Tannier, and V. Daubin. 2013. Efficient exploration of the space of reconciled gene trees. Syst. Biol. 62:901–912.

Tierney, L. 1994. Markov chains for exploring posterior distributions. Ann. Statist. 22:1701–1728.

Venables, W. N. and B. D. Ripley. 2002. Modern Applied Statistics with S. Fourth ed. Springer, New York.

West, D. 1979. Updating mean and variance estimates: An improved method. Communications of the ACM 22:532–535.

Whidden, C., R. G. Beiko, and N. Zeh. 2010. Fast FPT algorithms for computing rooted agreement forests: Theory and experiments. Pages 141–153 *in* Experimental Algorithms (P. Festa, ed.) vol. 6049 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg.

Whidden, C., R. G. Beiko, and N. Zeh. 2013. Fixed-Parameter algorithms for maximum agreement forests. SIAM J. Comput. 42:1431–1466.

Whidden, C. and N. Zeh. 2009. A unifying view on approximation and FPT ofagreement forests. Pages 390–401 *in* Proceedings of the 9th International Workshop, WABI 2009 vol. 5724 of *Lecture Notes in Bioinformatics* Springer-Verlag.

Whidden, C., N. Zeh, and R. G. Beiko. 2014. Supertrees based on the subtree prune-and-regraft distance. Syst. Biol. 63:566–581.

Yang, Z. and A. D. Yoder. 2003. Comparison of likelihood and bayesian methods for estimating divergence times using multiple gene loci and calibration points, with application to a radiation of cute-looking mouse lemur species. Systematic Biology 52:705–716.

Yoder, A. D. and Z. Yang. 2004. Divergence dates for malagasy lemurs estimated from multiple gene loci: geological and evolutionary context. Molecular Ecology 13:757–773.

Zhang, N. and M. Blackwell. 2001. Molecular phylogeny of dogwood anthracnose fungus (discula destructiva) and the diaporthales. Mycologia 93:355–365.

# SUPPLEMENTARY MATERIAL

## COMPUTING ACCESS TIME STATISTICS

The MAT between two trees $i$ and $j$, $\mathrm{MAT}(i,j)$, is the mean number of iterations before node $j$ is visited after node $i$ is visited. The mean commute time $\mathrm{MCT}(i,j)$ is the mean number of iterations for a random walk to visit tree $i$, then tree $j$, and finally return to tree $i$. The MCT for two trees can be computed as $\mathrm{MCT}(i,j) = \mathrm{MAT}(i,j) + \mathrm{MAT}(j,i)$. The $\mathrm{MRT}(i,T)$ is the mean number of iterations required to cover (visit) each tree in set $T$ starting from tree $i$ and then return to tree $i$. The MRT of a graph $T$ is the maximum of $\mathrm{MRT}(i,T)$ across nodes $i$. The MRT can be computed as the maximum MCT from node $i$ to a tree $t$ in $T$.

The MAT values (and hence MCT and MRT values) involving the highest probability tree, $t_0$, can be computed with a single pass through the tree posterior using a method for updating weighted means. To do so, we use dynamic programming and store three values: (1) $c_{ij}$, the number of times a topology $j$ has been seen since the last visit to topology $i$, (2) $m_{ij}$, the mean iteration number of each such visit, and (3) the current $\mathrm{MAT}(i,j)$ estimate. We perform updates when one of $i$ and $j$ is $t_0$ as follows. For each posterior sample $j$ with $j = t_0$, we update our values for each topology $i$. We update the access time estimates $\mathrm{MAT}(i,j)$ with weight $c_{ij}$ and value $m_{ij}$ and then reset the weight and value. We then update the stored values $c_{ji}$ and $m_{ji}$. If $j \neq t_0$, we apply the same update procedure but only for $i = t_0$. This requires linear storage with respect to the number of distinct compared topologies.

**Supplemental table 1.** TreeBASE identifiers and legacy identifiers for the eukaryotic data sets used in this study.

| Data | ID | Legacy ID |
|------|------|-----------|
| DS1 | M2017 | M336 |
| DS2 | M2131 | M501 |
| DS3 | M127 | M1510 |
| DS4 | M487 | M1366 |
| DS5 | M2907 | M3475 |
| DS6 | M220 | M1044 |
| DS7 | M2449 | M1809 |
| DS8 | M2261 | M755 |
| DS9 | M2389 | M1748 |
| DS10 | M2152 | M520 |
| DS11 | M2274 | M767 |

**Supplemental table 2.** Convergence diagnostics for the golden runs on eukaryotic datasets as reported by the MrBayes `sumt` and `sump` tools. We report the mean log likelihood ($\mu$LL), standard error of log likelihoods (Est LL error), maximum standard deviation of split frequencies (maxSDSF), maximum topological Gelman-Rubin potential scale reduction factor for splits (maxPSRF), and the minimum estimated sample size for the treelength parameter (ESS).

| Data | $\mu$LL | Est LL error | maxSDSF | maxPSRF | ESS |
|------|--------|-------------|---------|---------|-----|
| DS1 | -6,901.25 | 0.22 | 0.015 | 1.000 | 712,555 |
| DS2 | -26,166.93 | 0.68 | 0.001 | 1.000 | 739,035 |
| DS3 | -33,466.94 | 1.43 | 0.001 | 1.000 | 734,698 |
| DS4 | -13,034.24 | 0.27 | 0.002 | 1.000 | 724,599 |
| DS5 | -7,914.11 | 0.35 | 0.002 | 1.000 | 718,327 |
| DS6 | -6,298.55 | 0.45 | 0.007 | 1.000 | 738,077 |
| DS7 | -36,823.46 | 2.11 | 0.003 | 1.000 | 724,628 |
| DS8 | -8,123.66 | 0.61 | 0.003 | 1.000 | 653,977 |
| DS9 | -3,599.49 | 0.46 | 0.002 | 1.000 | 725,545 |
| DS10 | -9,537.43 | 0.82 | 0.003 | 1.000 | 729,304 |
| DS11 | -5,725.33 | 1.20 | 0.003 | 1.000 | 728,818 |

**Supplemental table 3.** The standard error of topology posterior probabilities for the top trees between the golden run eukaryote posteriors. Columns labeled by posterior probabilities of the top trees. Dashes from the left indicate that all posterior probabilities satisfied the given threshold. Dashes from the right indicate that all posterior probabilities fit within the next smaller threshold.
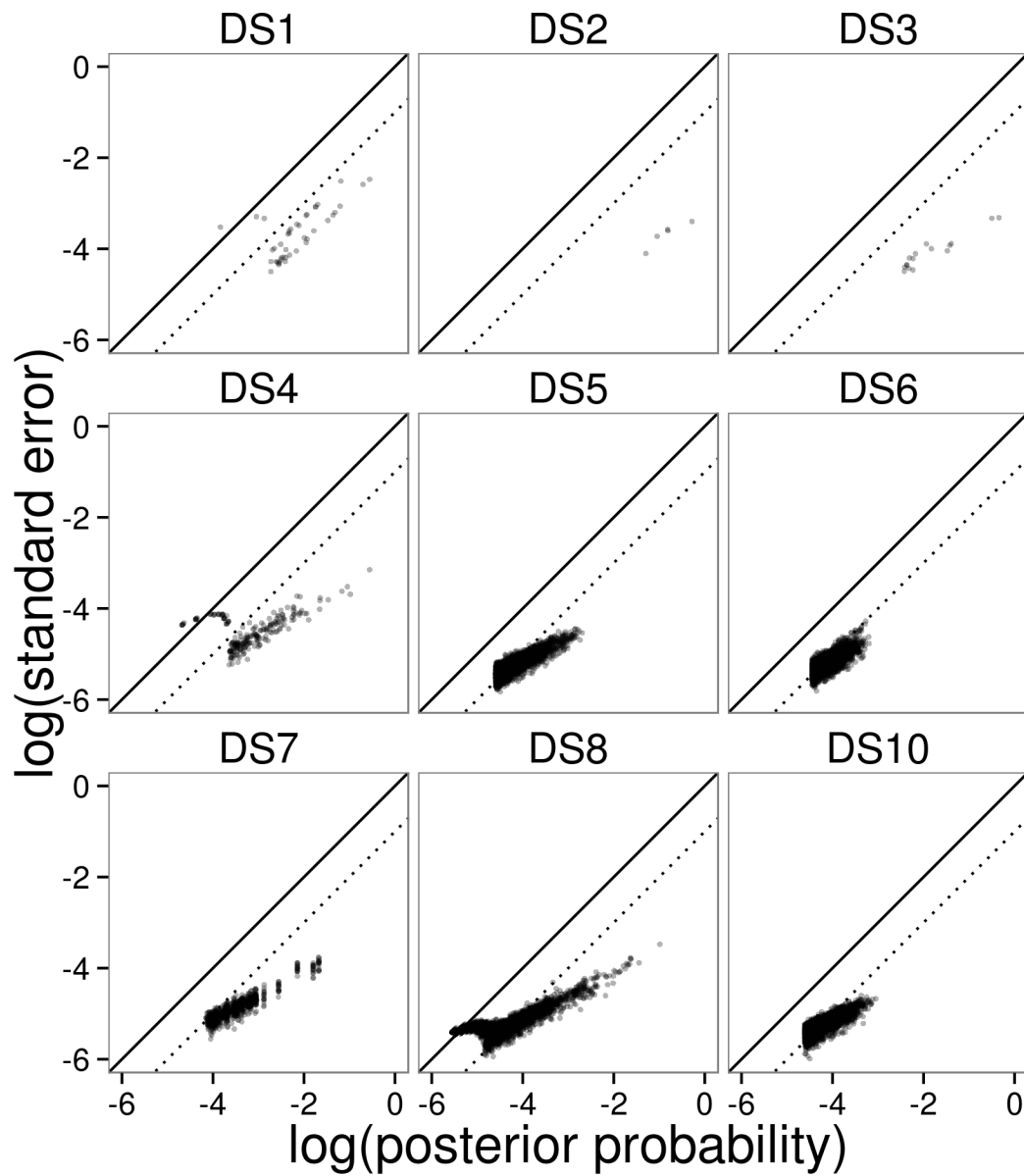
| Data | < 1e-05 | < 1e-04 | < 1e-03 | < 1e-02 | < 1e-01 | < 1e+00 |
|------|---------|---------|---------|---------|---------|---------|
| DS1 | - | - | 5.07e-04 | 5.07e-04 | 3.07e-03 | 3.35e-03 |
| DS2 | - | - | - | - | 1.88e-04 | 3.98e-04 |
| DS3 | - | - | - | 7.84e-05 | 1.30e-04 | 4.82e-04 |
| DS4 | - | 7.57e-05 | 7.58e-05 | 1.87e-04 | 3.00e-04 | 7.09e-04 |
| DS5 | - | 1.19e-05 | 3.13e-05 | 3.57e-05 | - | - |
| DS6 | - | 1.26e-05 | 5.29e-05 | - | - | - |
| DS7 | - | 1.11e-05 | 3.55e-05 | 1.31e-04 | 1.73e-04 | - |
| DS8 | 7.10e-06 | 1.19e-05 | 4.91e-05 | 8.23e-05 | 1.69e-04 | 3.34e-04 |
| DS10 | - | 1.28e-05 | 2.50e-05 | - | - | - |

**Supplemental table 4.** Sets of identical sequences and multifurcations for each dataset with at least one multifurcation in the majority rule consensus tree. Multifurcations increase the size of the true credible set by the factor noted.
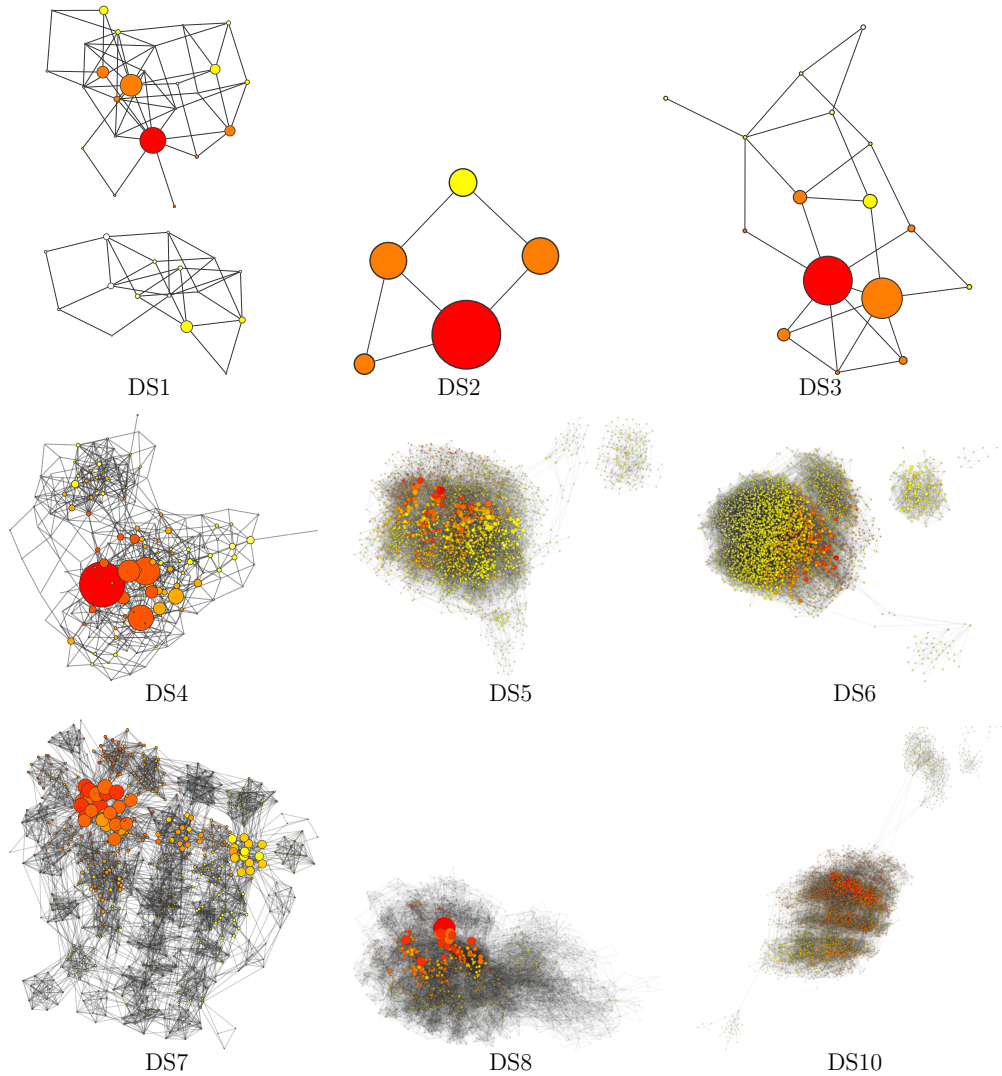
| Data | Identical Sequences | Multifurcations | Inflation Factor |
|------|---------------------|-----------------|------------------|
| DS5  | 0                   | 3,3,4           | 375              |
| DS6  | 0                   | 3,3,3           | 27               |
| DS7  | 2,2                 | 4               | 15               |
| DS9  | 2,2,3,3,4           | 4,4,5           | 23,625           |
| DS10 | 0                   | 5               | 105              |
| DS11 | 2,3,3,4,5,9         | 3,4,4,4,6,12    | more than 1e+13  |
| VL1  | 2,2,2,2,2           | 3               | 3                |
| VL2  | 2,2,2,3,4           | 4,4             | 225              |
| VL3  | 2,2,2,2,2           | 3,3,3           | 27               |
| VL4  | 2,2,2,3,3,3,5       | 3,3,3,6         | 25,515           |
| VL5  | 2,2,2,2,3,4         | 3,6             | 315              |
| VL6  | 2,2,2,2,2,2,3       | 3,3,4           | 135              |

**Supplemental figure 1.** Two trees such that the best rooting SPR distance overestimates the unrooted SPR distance. This occurs when every leaf is part of some moved subtree in every minimal unrooted set of SPR operations. In this example, the unrooted distance remains 2, while the best rooting SPR distance is $\lfloor \frac{m+3}{2} \rfloor$ for $m \geq 4$. Dashed lines indicate the edges modified by these minimal unrooted SPR operations.
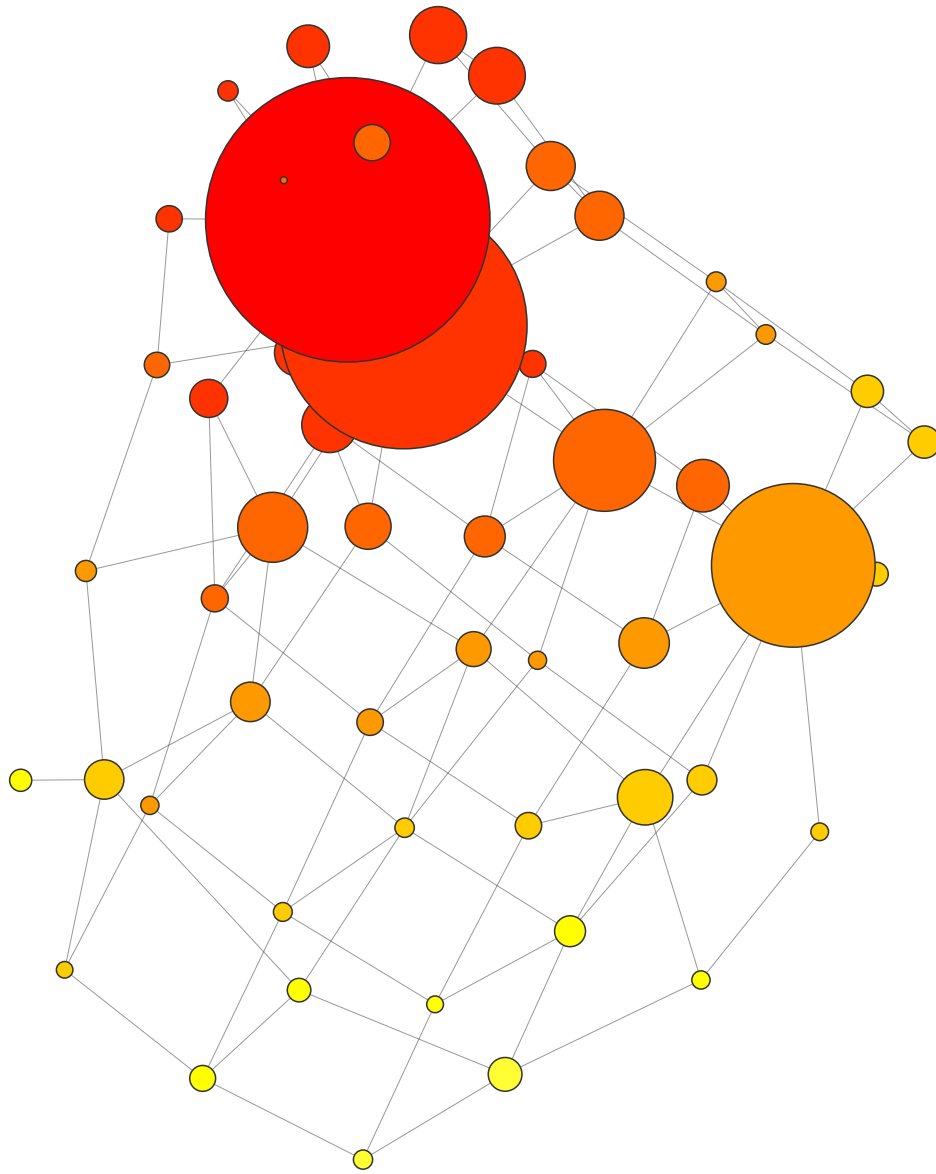
**Supplemental figure 2.** The standard error of topology posterior probabilities for the top trees between the golden run eukaryote posteriors. The standard error is smaller than the posterior probability for estimates below the solid line and an order of magnitude smaller for estimates below the dotted line.

**Supplemental figure 3.** Distance SPR graphs of the combined golden run eukaryote posteriors. Each graph contains either the 95% credible set or the 4096 topologies with highest estimated posterior probabilities (DS5, DS6 and DS10). Node areas are scaled relative to posterior probability (PP; larger = higher probability) within each graph (but not with respect to the other graphs). Node color indicates SPR distance from the topology with highest posterior probability in each dataset on a red-yellow-white scale (dark-light in the print version), with the highest probability tree colored red.

**Supplemental figure 4.** Distance SPR graph of the 95% credible set of the combined golden run eukaryote posteriors for DS7 with 3 of the 4 nearly identical *Microcebus rufus* sequences removed. Area indicates posterior probability and color indicates SPR distance from the topology with highest posterior probability on a red-yellow-white scale (dark-light in the print version), with the highest probability tree colored red.
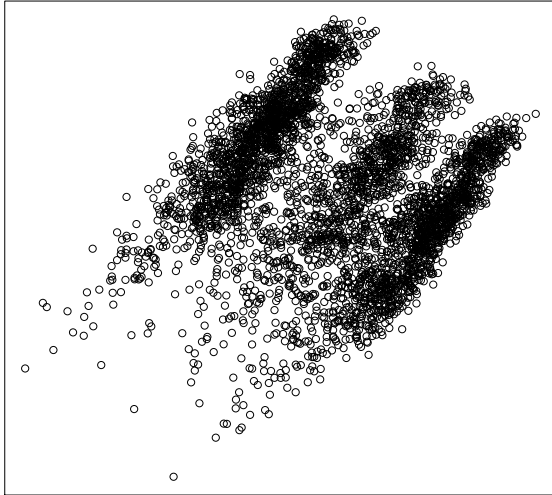
**Supplemental figure 5.** Central trees of the two topological peaks in dataset DS6. Only three SPR operations separate these trees, moving the colored subtrees. Intermediate groupings are unsupported and intermediate trees unlikely.

**Supplemental figure 6.** Extended Majority Rule Consensus tree of the DS1 golden runs. Node labels indicate the percentage of trees within the 95% credible set containing that split. Four splits, indicated in red, receive erroneously high support when the second peak is not sampled (numbers in brackets)
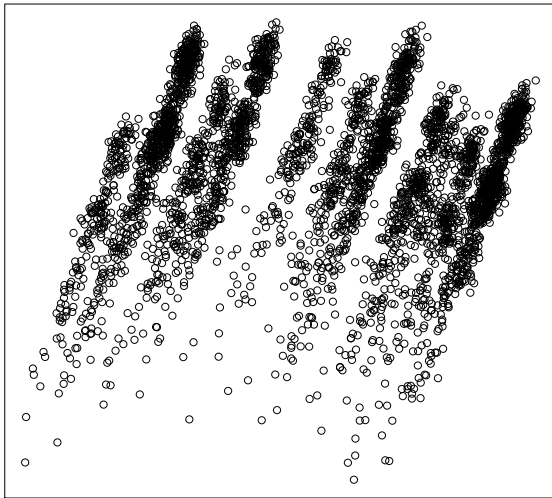
**Supplemental figure 7.** Comparison of mean running time and split frequency error with and without Metropolis-coupling using varying numbers of runs.
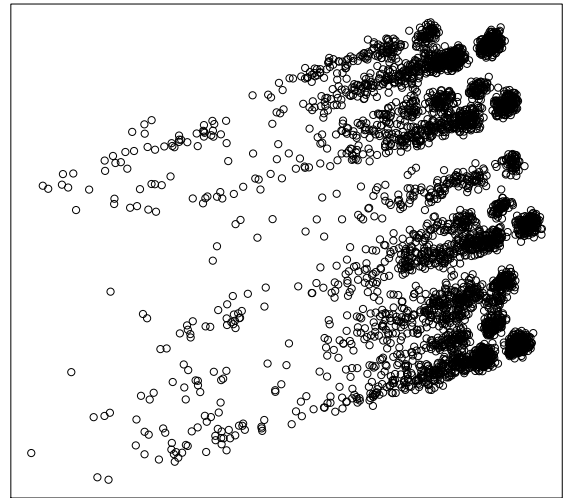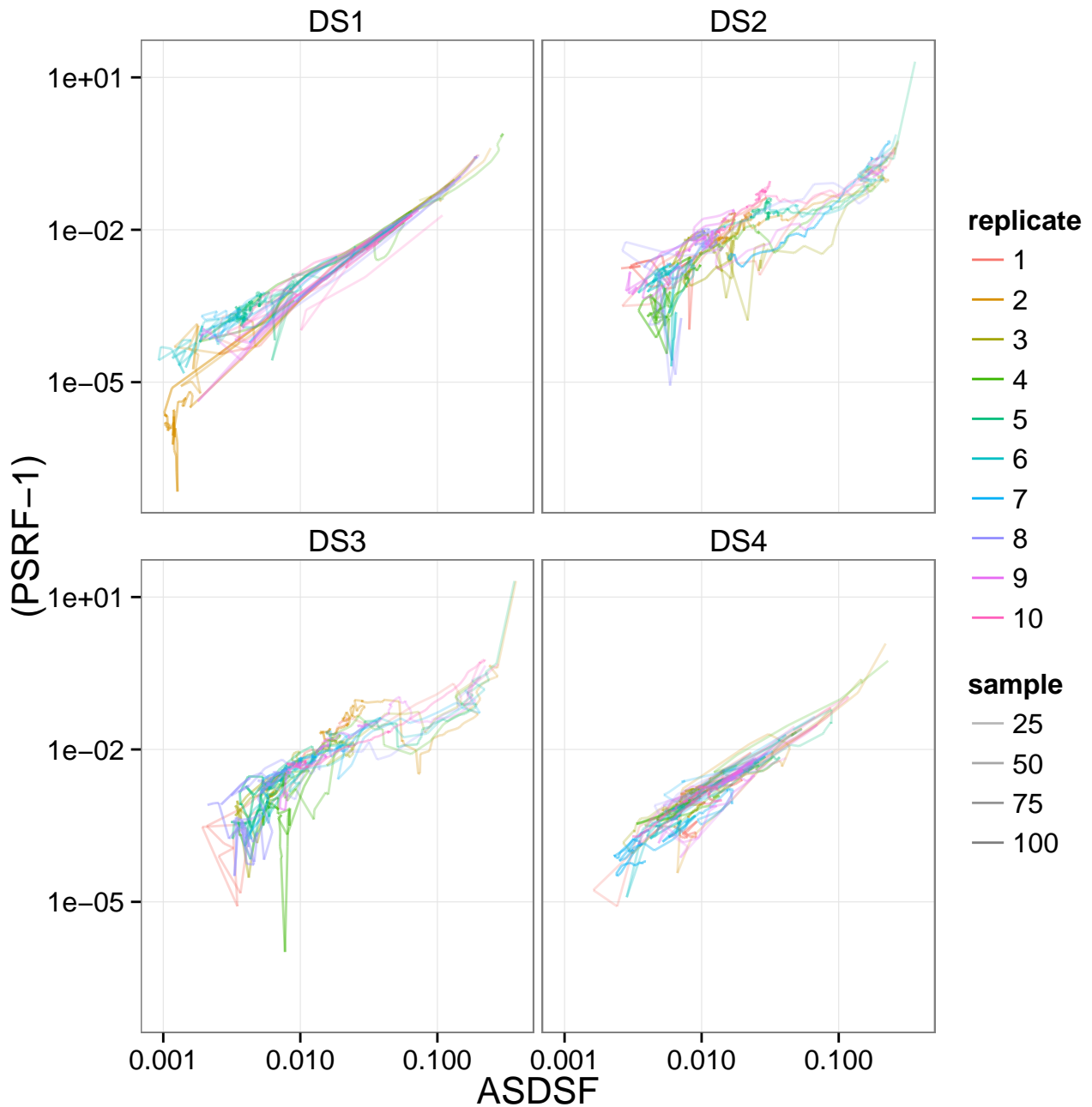
(a) DS9 (SPR)

(b) DS9 (RF)

(c) DS11 (SPR)

(d) DS11 (RF)

**Supplemental figure 8.** Comparison of multidimensional scaling representations with SPR and RF distances for flat posteriors DS9 and DS11.

**Supplemental figure 9.** A comparison of the average standard deviation of split frequencies and SPR topological Gelman-Rubin-like convergence diagnostics for datasets DS1, DS2, DS3, and DS4 using 2 independent runs. Values are shown on a log-log scale in base 10. 100 evenly-spaced samples were taken from the first 2 runs of each 8-run replicate that had achieved ASDSF of less than 0.01. Transparency decreases with time.