

0001 **Running head:** Geodesic Info

0002

0003

0004

0005

0006 Estimating Bayesian Phylogenetic Information

0007 Content Using Geodesic Distances

0008 Analisa Milkey¹ and Paul O. Lewis¹

0009

0010 ¹ *Department of Ecology and Evolutionary Biology, University of Connecticut, 75 N. Ea-*

0011 *gleville Road, Unit 3043, Storrs, Connecticut 06269, U.S.A.*

0012

0013 **Corresponding author:** Paul Lewis, Department of Ecology and Evolutionary Biology,

0014 University of Connecticut, 75 N. Eagleville Road, Unit 3043, Storrs, Connecticut 06269,

0015 U.S.A.; Tel: +01 860 486-2069; FAX: +01 860 486-6364; E-mail: paul.lewis@uconn.edu

0016

0017

0018

0019

0020

0021

0022

0023

0024

0025

ABSTRACT

A new Bayesian measure of phylogenetic information content is introduced based on geodesic distances in treespace. The measure is based on the relative variance of phylogenetic trees sampled from the posterior distribution compared to the prior distribution. This ratio is expected to equal 1 if there is no information in the data about phylogeny and 0 if there is complete information. Trees can be scaled to have the same mean tree length to avoid dominance by edge length information and focus on topological information. The method scales well, requiring only that a valid sample can be obtained from both prior and posterior distributions. We show how dissonance (information conflict) among data sets can also be estimated. Both simulated and empirical examples are provided to illustrate that the new approach produces sensible and intuitive results.

Keywords: geodesic, tree distance, phylogeny, Bayesian, information

INTRODUCTION

The canonical paraphrase of Claude Shannon’s concept of *information* in his 1948 foundational paper on information theory (Shannon, 1948) is that information is the resolution of uncertainty. The amount of information in data is often equated with the amount of data itself, but it is clearly true that a large quantity of data may contain little information (e.g. an entire book filled with randomly-generated text) and a small quantity of data may eliminate all uncertainty (e.g. the six words “I will have the Matar Paneer” uttered at a restaurant with hundreds of menu items). Systematists have long been concerned about how much information relevant to the problem of interest is available in the data they collect. Early work was concerned with determining the degree to which noise from homoplasy was affecting results (consistency index; Kluge and Farris, 1969) and how confident we should feel about different clades in a phylogenetic tree (bootstrapping; Felsenstein, 1985). Archie (1989) and Faith (1991) developed permutation methods designed to test whether there was a significant amount of historical information in data. Hillis and Huelsenbeck (1992) advocated using parsimony tree length skewness to assess the informativeness of data. This early work was followed by many papers addressing various aspects of information content in the data used by systematists to infer phylogenies (Steel et al., 1993, 1995; Lyons-Weiler et al., 1996; Goldman, 1998; Massingham and Goldman, 2000; Shpak and Churchill, 2000; Xia et al., 2003; Geuten et al., 2007; Townsend, 2007; Shi et al., 2008; Fischer and Steel, 2009; Lemey et al., 2009; San Mauro et al., 2009; Xia, 2009; Tippers et al., 2012; Townsend et al., 2012; Brown, 2014; Lewis et al., 2016; Duchêne et al., 2022).

The Bayesian statistical framework provides a reference distribution (the prior) for comparison with the posterior distribution. Information in data transforms the prior distribution into the posterior distribution, increasing the plausibility of some parameter combi-

0076 nations at the expense of other combinations. Comparing the posterior distribution to the
0077 prior distribution provides an explicit way to measure the information contained in the data
0078 (Lindley, 1956).

0079 Lewis et al. (2016) proposed measures of phylogenetic information content and phylo-
0080 genetic dissonance based on the relative entropy of a posterior distribution of tree topologies
0081 compared to the prior topology distribution. Common practice disperses prior probability
0082 mass evenly over every possible tree topology (i.e. maximum entropy). Zero topological **in-**
0083 **formation** occurs when the posterior distribution exactly matches this maximum-entropy
0084 prior distribution. Complete information obtains when the entire posterior is concentrated
0085 over a single tree topology (i.e. minimum entropy).

0086 Because information can be misinformation, it is important to also have ways of mea-
0087 suring conflict information among data sets. Phylogenetic **dissonance** was defined (Lewis
0088 et al., 2016) as the difference between the entropy of the merged posterior tree samples from
0089 two or more data subsets (e.g. loci) and the average posterior entropy within each subset.
0090 If all data subsets concentrate posterior mass in similar ways, then merging the sampled
0091 trees is expected to be simply a larger sample from the same distribution. If, however, some
0092 data sets disagree with others about which tree topologies are best, then the average entropy
0093 within datasets is expected to be much smaller than the entropy of the merged samples.

0094 One of the drawbacks of the method described by Lewis et al. (2016) is scalability.
0095 The number of possible tree topologies quickly becomes too vast to sample adequately as
0096 the number of taxa increases. Consider a problem involving trees of only 12 taxa. Suppose
0097 1 million tree topologies were sampled from the posterior distribution and every sampled
0098 tree topology was distinct from all other sampled topologies. It would seem that there
0099 is not much information about tree topology in the data because even a sample of size 1
0100

0101 million fails to find any tree topology that is more plausible than the others. Given that
0102 a flat prior distribution for unrooted trees of 12 taxa distributes probability evenly over all
0103 654,729,075 possible tree topologies, the fact that the posterior sample is concentrated over
0104 a tiny fraction (0.001527) of the possible tree topologies suggests the opposite: that there
0105 is a tremendous amount of information in the data. In order to confidently conclude that
0106 there is no information in this data set, one would need to have a posterior sample size
0107 several times larger than the number of possible trees. While it is possible to ameliorate the
0108 problem by adjusting the maximum entropy to equal the log of the sample size rather than
0109 the log of the number of possible tree topologies, or by smoothing the empirical posterior
0110 probability distribution of tree topologies using conditional clade distributions (Lewis et al.,
0111 2016; Berling et al., 2025), this overestimation of information content due to the impossibility
0112 of adequately approximating the posterior distribution is a complication. The problem only
0113 gets much, much worse if the number of taxa is in the hundreds or thousands as opposed to
0114 12.

0115 This paper describes a very different approach to measuring phylogenetic information
0116 content that pays attention to edge lengths as well as topology and scales much better
0117 with larger numbers of taxa. In addition to measuring information content, we describe a
0118 corresponding measure of dissonance.

0120 BACKGROUND

0121
0122 Brown and Owen (2020) discussed how to estimate the (Fréchet) mean and variance of a set
0123 of phylogenetic trees. In the context of Bayesian phylogenetics, the variance of a posterior
0124 sample of trees relative to the variance of a prior sample is useful as a measure of information
0125 content. Consider two mean trees computed from samples of trees from the posterior and

0126 prior distributions, respectively (Fig. 1). The data used for the posterior sample comprise
0127 1296 DNA sites from the RuBisCO large subunit (*rbcl*) locus from 5 taxa chosen to span
0128 the green plant phylogeny. The prior and posterior samples were obtained using RevBayes
0129 version 1.3.1 (Höhna et al., 2016) (for details of settings, see the *1-fig-rbclmean* directory in
0130 the Supplementary Materials). The trees shown were determined using the Sturm (2003)
0131 algorithm as described in Miller, Owen, and Provan (2015) and Brown and Owen (2020).
0132 The posterior mean tree would make perfect sense to someone familiar with green plant
0133 evolution. For example, it places the two seed plants, *Picea* and *Avena*, as sister taxa in a
0134 tree rooted at the green alga *Chara* as well as placing all vascular plants (*Asplenium*, *Picea*,
0135 and *Avena*) in a clade. On the other hand, the prior mean is completely unresolved, which
0136 is expected because the prior treats every unrooted tree topology as equally probable. Note
0137 the relative scales of these trees. The posterior mean tree length is smaller than the prior
0138 mean tree length by an order of magnitude. Thus, the data have information about both
0139 the topology of the tree as well as the edge lengths.

0140 The mean tree is the tree minimizing the squared distance to all trees in a sample
0141 of phylogenetic trees. The variance of the sample is therefore obtained as a free byproduct
0142 of determining the mean tree. The distance used is the geodesic distance of Owen and
0143 Provan (2010) in the treespace characterized by Billera et al. (2001). For the example in
0144 Fig. 1, the variance of the posterior sample is 0.000728, while the variance of the prior
0145 sample (40.1) is 55 thousand times larger. In this paper, we propose using the difference
0146 between the prior and posterior sample phylogenetic variance as a measure of information
0147 content. A posterior variance equal to zero corresponds to complete information, whereas
0148 a posterior variance equal to the prior variance implies there is no information about the
0149 phylogeny in the observed data (although there may well be considerable information about
0150 other components of the model, such as nucleotide frequencies).

0151
0152
0153
0154
0155
0156
0157
0158
0159
0160
0161
0162
0163
0164
0165
0166
0167
0168
0169
0170
0171
0172
0173
0174
0175

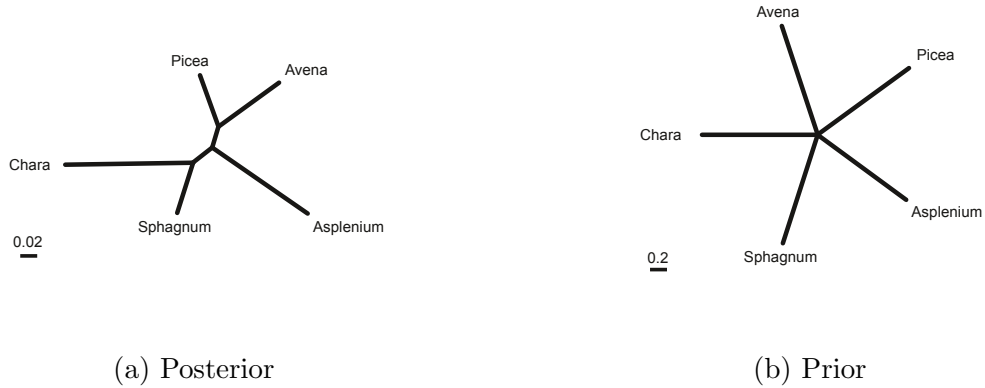


Figure 1: Fréchet mean trees computed from the (a) posterior and (b) prior samples using the *rbcL* dataset.

MATERIALS AND METHODS

INFORMATION MEASURE

Our method requires two samples of phylogenetic trees: a sample of size N_0 from the prior distribution and a sample of size N from the posterior distribution. In this paper, we assume $N_0 = N$ and that both prior and posterior samples are obtained using MCMC. In RevBayes 1.3.1, for example, a sample from the prior may be obtained by calling the `ignoreAllData()` function of the model.

Shi et al. (2022) proposed the information content measure LCR (Log Concentration Ratio):

$$LCR = \log \left\{ \frac{V_0}{V} \right\}, \tag{1}$$

where V_0 and V are (in Shi et al., 2022) the volumes of the 95% credible regions for the prior and posterior distributions, respectively. We equate volume (V and V_0) with the 95%

0176 radius. To compute the 95% radius, a mean tree is obtained from a sample of trees, and the
0177 distance between that mean tree and each sampled tree is calculated. From these distances,
0178 the radius of the smallest hypersphere containing at least 95% of the distances from the
0179 mean tree is used as the “volume”. The danger of measuring volume in a single dimension
0180 is that information will be underestimated. For example, if the radius is reduced from 2 to
0181 1, the volume is halved if assuming 1 dimension but is reduced by a factor of 4 if assuming
0182 even 2 dimensions. Nevertheless, we show that our 1-dimensional *LCR* behaves in a sensible
0183 and intuitive way for simulated data.

0184 In addition to the 95% radius (hereafter denoted the RAD method), we explored two
0185 other ways of measuring dispersion:

- 0186
- 0187 • Letting V equal the standard deviation of distances from the mean tree (the STD
0188 method) and
- 0189 • Letting V equal the radius associated with the 95% highest posterior density credible
0190 set of trees (the HPD method)
- 0191

0192 The HPD method results in a larger radius if the posterior is asymmetric (Fig. S1, Supple-
0193 mentary Materials). We found (Figs. S2 and S3, Supplementary Materials) that, while both
0194 alternative ways of measuring information largely agree with information measured using the
0195 RAD method, we prefer the RAD method because it is both more stable (avoiding the undue
0196 influence of occasional outliers that the STD approach would include) and more intuitive
0197 (avoiding most of the cases of negative information content that are often produced by the
0198 HPD method when information content is low). Ideally, V would equal the actual *volume*
0199 of treespace (as in the original LCR formulation by Shi et al. (2022)) but, unfortunately,
0200 treespace is only partially Euclidean and, lacking any way to accurately measure the volume

of treespace regions, we default to the 95% radius as our measure of dispersion in this paper.

Some may find LCR difficult to interpret given the log scale. Transforming LCR to percent information (I) may be preferable:

$$I = 100 \left\{ \frac{V_0 - V}{V_0} \right\} = 100 \left(1 - e^{-LCR} \right). \quad (2)$$

Whereas LCR ranges from 0 (no information¹) to ∞ (complete information), I ranges from 0 (no information) to 100 (complete information) and thus has a similar interpretation to the measure of Lewis et al. (2016).

SCALING TO A COMMON TREE LENGTH

Under normal circumstances, we expect V to be considerably smaller than V_0 , reflecting the fact that sequence data generally contains information relevant to both tree topology and edge lengths. Information content measured by LCR is often dominated by the difference in tree length between the posterior and prior sample. One way to reduce the influence of tree length on information content, giving the topological information more influence, is to scale trees in both prior and posterior samples such that the prior mean tree and the posterior mean tree both have a common tree length. This rescaling does not mean that the only information measured is topological; the data are expected to induce correlations among edge lengths in the posterior sample that are not present in the prior sample. Such correlations affect the posterior variance even though the prior sample is scaled to match the posterior mean tree length. Unless explicitly stated, information content using equation (1) scales sampled trees so that the prior mean tree and posterior mean tree each have the

¹but can be negative if the posterior has higher variance than the prior

0226 same total length of 1.0.

0227 Both the Fréchet mean tree and tree distances were calculated in this paper using the
0228 open-source software `op`, available from <https://github.com/plewis/op>. For all analyses
0229 in this paper, we instructed the `op` program to stop when the maximum pairwise distance
0230 among the mean trees produced by the $N = 10$ most recent iterations first dropped to less
0231 than $\epsilon = 0.00001$.

0233 DISSONANCE MEASURE

0235 Dissonance between two posterior distributions of trees is measured as a modified effect size.
0236 Let d_{12} be the geodesic distance between the scaled mean trees of the two data sets. Each
0237 mean tree is scaled to have tree length 1.0. Let n_1 and n_2 be the sample sizes and s_1^2 and
0238 s_2^2 the scaled sample variances of the two posterior distributions, respectively. Variances are
0239 scaled using the same scaling factor used to scale mean trees: if the mean tree length for
0240 data set 1 is L , the mean tree is scaled by dividing each edge by L and the variance is scaled
0241 by dividing the unscaled variance by L^2 . The effect size may be computed as Cohen's d :

$$0242 \text{Cohen's } d = \frac{d_{12}}{\sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}}} \quad (3)$$

0245 For compatibility with our information measure, we use a modified effect size that
0246 substitutes the 95% radius for s_1 and s_2 . The formula for dissonance (D) is thus

$$0247 \text{D} = \frac{d_{12}}{\sqrt{\frac{(n_1-1)r_1^2 + (n_2-1)r_2^2}{n_1+n_2-2}}}, \quad (4)$$

0248 where d_{12} is the geodesic distance between the two (scaled) mean trees, and r_1 and r_2 are
0249

0251 the 95% radii for the two (scaled) posterior samples.

0252 Dissonance is expected to be zero if the two posterior distributions compared are
0253 identical (and thus have identical mean trees), and is expected to grow with increasing
0254 distance between the mean trees. Note that two data sets simulated independently from
0255 identical model trees are *not* expected to exhibit zero dissonance. Even though conditionally
0256 independent, the two data sets are not identical and thus have distinct posterior distributions.
0257

0258 SIMULATION EXPERIMENTS

0259 INFORMATION

0260 We conducted a series of simulation experiments to characterize whether *LCR* performs as
0261 expected for differing substitution rates, sequence lengths, proportion missing data, and rate
0262 heterogeneity. Information content is expected to be low if the substitution rate is near zero.
0263 At the extreme of zero substitution rate, all sites would be constant and the data would
0264 therefore provide no evidence relevant for identifying historical groupings of taxa. At the
0265 ideal substitution rate, information content is expected to be 100% given that sequences
0266 are of sufficient length to provide complete evidence for every split. Increasing the substi-
0267 tution rate much beyond this ideal rate should result in a decrease in information because
0268 sequences eventually become saturated with substitutions (newer substitutions overwrite
0269 older substitutions that provided evidence of history). High among-site rate heterogeneity
0270 (ASRV) should also result in decreased phylogenetic information content, even if the mean
0271 substitution rate is ideal, because, with high ASRV, most sites are evolving at a rate less
0272 than the mean rate, and the remaining sites evolve at a high rate. Finally, a high proportion
0273 of missing data (missing-at-random) should decrease information content, and information
0274
0275

0276 should decrease with decreasing number of sites.

0277 We measured information content on simulated DNA sequence data sets with 5 taxa
0278 and variable sequence lengths, mean substitution rates, amount of ASRV, and percent miss-
0279 ing data. The tree used as the model tree for these simulations is that shown in Fig. 1a. DNA
0280 sequences on this tree were generated using Seq-Gen 1.3.4 (Rambaut and Grassly, 1997) and
0281 analyzed using RevBayes 1.3.1 using 50,000 iterations following a 2000-iteration burn-in pe-
0282 riod and a GTR substitution model with Gamma-Dirichlet tree length/edge proportion prior
0283 (Zhang et al., 2012). Trees were sampled every 5 iterations, yielding 10,000 sampled trees.
0284 Scripts for all analyses are provided in the supplementary materials (directory *2-fig-1cr*).
0285

0286 DISSONANCE

0287
0288 To evaluate dissonance, we conducted a simulation experiment in which 20 starting trees,
0289 each of 26 taxa, were drawn from the prior and moved via a random walk through treespace.
0290 Each random walk involved 500 steps and each step consisted of adding a Normal variate
0291 with mean zero and standard deviation 0.005 to each edge length. Pendant (leaf) edges that
0292 dropped below 1×10^{-12} were reflected back into the interval $[1 \times 10^{-12}, \infty)$. Internal edges
0293 that dropped below zero resulted in a change in topology; one of the two possible alternative
0294 splits was chosen uniformly at random and received the residual edge length. Six trees were
0295 sampled at uniformly spaced points (stations) along this path (i.e. stations 1, 2, 3, 4, 5, 6
0296 used trees sampled from steps 0, 100, 200, 300, 400, and 500). For each of the 6 stations
0297 along the path, for each of the 20 replicates, two data sets were simulated using Seq-Gen
0298 1.3.4 and posterior samples were obtained using RevBayes 1.3.1. The first simulated data
0299 set corresponded to the starting tree for that replicate (i.e. station 1); the second data set
0300 corresponded to the tree sampled at that particular station. The first station thus produced

0301 two replicate simulated data sets using the same model tree. Dissonance was measured for
0302 each pair of data sets at each of the 6 stations for all 20 replicates. The expectation is that
0303 dissonance will grow with increasing distance between model trees; however, as mentioned
0304 previously, we do not expect the dissonance between the two data sets simulated from the
0305 same model tree at station 1 to be zero because the data sets, while conditionally indepen-
0306 dent, are not identical. Scripts for all analyses are provided in the supplementary materials
0307 (directory *3-fig-dissonance*)

0309 EMPIRICAL ANALYSES

0311 SATURATION

0312 It is often assumed *a priori* that 3rd codon position sites in protein coding genes are satu-
0313 rated. Lewis et al. (2016) estimated topological information content in 2nd vs. 3rd position
0314 sites in the *psaB* locus in Spheroplealean green algae and found that the 3rd position sites
0315 were in fact not saturated and contained more information than 2nd position sites. We
0316 reanalyzed these data using equation (1). We also applied the saturation test in PhyloMAd
0317 (Duchêne et al., 2018) as an independent check for saturation. Scripts for all analyses are
0318 provided in the supplementary materials (directory *4-fig-saturation*)

0320 DISSONANCE

0322 We re-analyzed an example from Lewis et al. (2016) in which high dissonance is expected.
0323 The 5' half of the mitochondrial *rps11* locus in Bloodroot (*Sanguinaria*, Papaveraceae) is
0324 vertically transferred, whereas the 3' half was evidently horizontally transferred from a
0325 phylogenetically-distant monocot (Bergthorsson et al., 2003). We reanalyzed these data

0326 using equation (3). Scripts for all analyses are provided in the supplementary materials
0327 (directory *5-fig-bloodroot*)

0329 RESULTS

0331 SIMULATION EXPERIMENTS

0333 INFORMATION

0334
0335 Information content was highest when the relative substitution rate was 1 (i.e. the edge
0336 lengths used for simulating data were equal to those in Fig. 1a) and decreased when the
0337 data were simulated with either larger or smaller relative rates (Fig. 2a). Information con-
0338 tent was zero when the number of sites was zero and generally increased as the sequence
0339 length increased, with the exception of the 1-site case, which, on average, resulted in slightly
0340 negative information (Fig. 2b). Information content decreased as the percentage of missing
0341 nucleotides grew from 10% to 95% (Fig. 2c). Finally, information content decreased with
0342 increasing variance in substitution rates across sites (Fig. 2d).

0344 DISSONANCE

0345
0346 For random walks beginning at trees drawn from the prior, the distance between trees in-
0347 creased as a function of the number of random walk steps taken (Fig. 3a). The mean geodesic
0348 distance between starting and ending trees was 0.606209 (std. dev. 0.066439, $n = 20$). As
0349 expected, the dissonance between two data sets is strongly and positively correlated with the
0350 geodesic distance between the trees used to generate the data (Fig. 3b). Only 4 interstation

0351 segments (of 5 segments \times 20 replicates = 100 total) failed to show increased dissonance.
0352 Dissonance was, on average, 0.124232 (std. dev. 0.011204, $n = 20$) when both data sets
0353 were simulated from the same tree and increased with increasing distance between model
0354 trees. A dissonance of 0.12 means that the distance between means represents 12% of the
0355 pooled radius. The total dissonance between the starting and ending data sets was 0.734085
0356 (std. dev. 0.083814, $n = 20$). Thus, the random walks were not long enough to completely
0357 separate the starting and ending posterior distributions. There is a strong, positive corre-
0358 lation (0.96, $n = 120$, $P < 0.0001$) between the distance separating model trees and the
0359 dissonance of the corresponding posterior distributions.

0361 EMPIRICAL ANALYSES

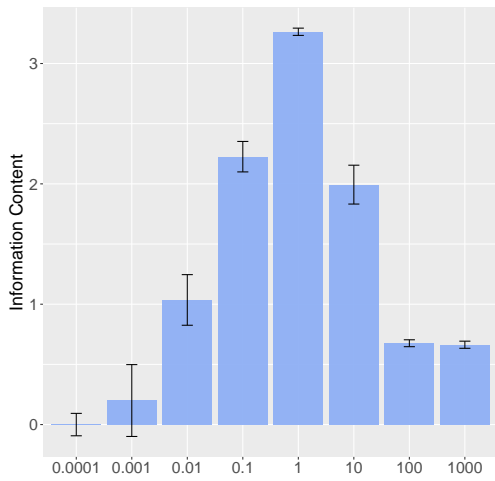
0363 SATURATION

0364 Like Lewis et al. (2016), we found that 3rd position sites at the *psaB* locus in the green
0365 algal data set contained more information ($LCR = 2.73$, $I = 93.5$) than 2nd position sites
0366 ($LCR = 1.75$, $I = 82.6$). The Fréchet mean trees computed from posterior samples from
0367 the two sets of sites clearly show greater resolution in the 3rd-position sites (Fig. 4), which
0368 reflects the greater amount of phylogenetic information in 3rd-position sites.

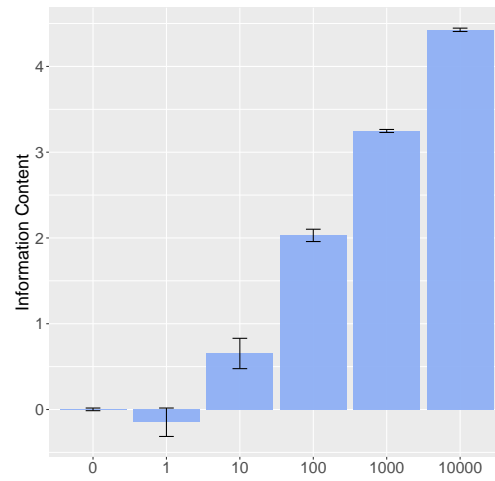
0369
0370 3rd-position sites also had more information than 1st-position sites ($LCR = 2.36$, $I =$
0371 90.5) and even slightly more information than 1st and 2nd positions combined ($LCR = 2.58$,
0372 $I = 92.5$). As expected, combining sites from all three codon positions yielded the most
0373 information ($LCR = 3.48$, $I = 96.9$).

0374 A natural question is “Is the information in 3rd-position sites truthful information
0375 or is it misinformation?” Using either BHV distances (BHV; Owen and Provan, 2010), or

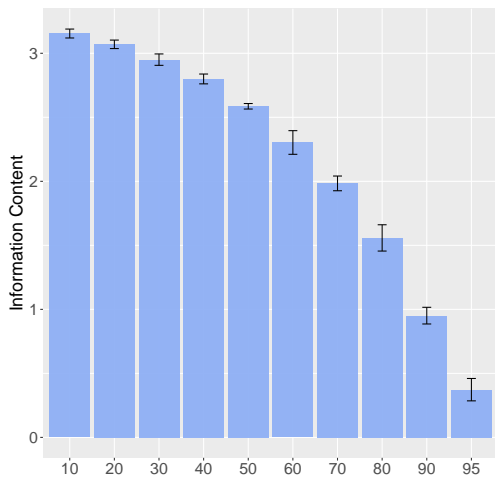
0376
0377
0378
0379
0380
0381
0382
0383
0384
0385
0386
0387
0388
0389
0390
0391
0392
0393
0394
0395
0396
0397
0398
0399
0400



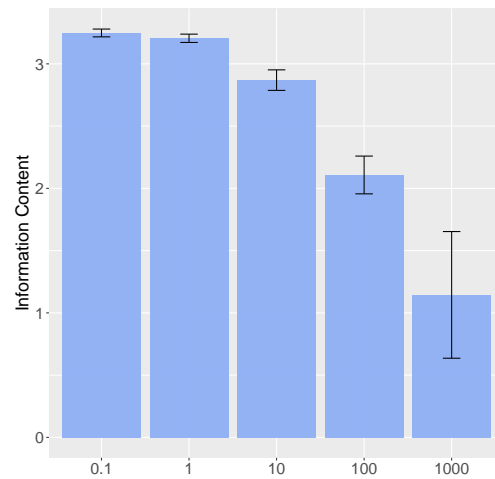
(a) Relative Subst. Rate



(b) No. sites



(c) Missing %



(d) ASRV

Figure 2: Information content (LCR) for 5-taxon datasets. a. Rate is relative to the edge lengths in the model tree. b. No. sites is the sequence length. c. Missing % is the percent of nucleotides missing at random. d. ASRV is measured as the variance in among-site relative rates. In each case, rate=1.0, ASRV=0.0, no. sites=1000, and missing = 0% except for the variable on the x-axis. Error bars show mean \pm standard deviation across $n = 10$ replicates.

Robinson-Foulds distances(RF; Robinson and Foulds, 1981), the 3rd-position mean tree is closer to the all-sites mean tree than is the 2nd-position mean tree (Table 1). This may mean only that the 3rd-position sites, with their greater information content, dominate in

0401
0402
0403
0404
0405
0406
0407
0408
0409
0410
0411
0412
0413
0414
0415
0416
0417
0418
0419
0420
0421
0422
0423
0424
0425

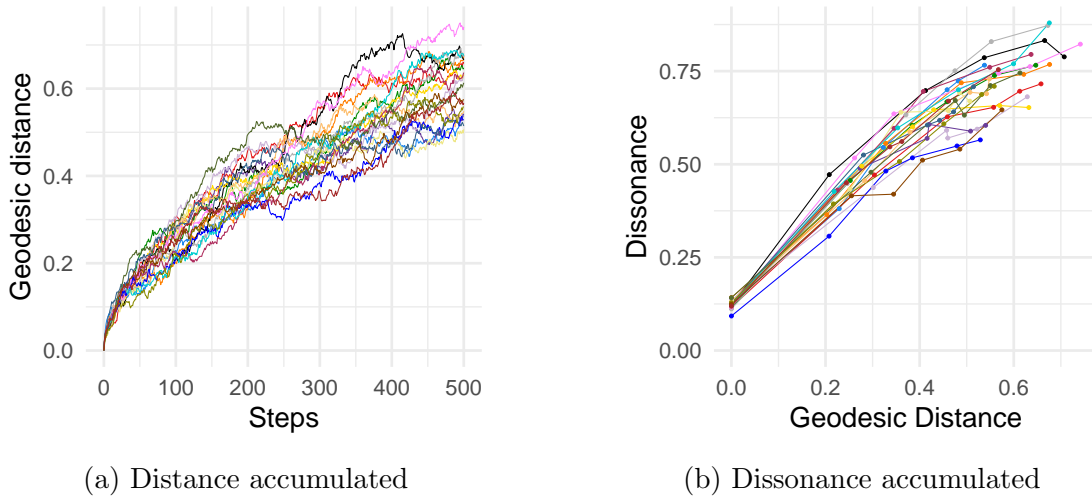


Figure 3: (a) Geodesic distance accumulated during random walk. (b) Dissonance as a function of geodesic distance.

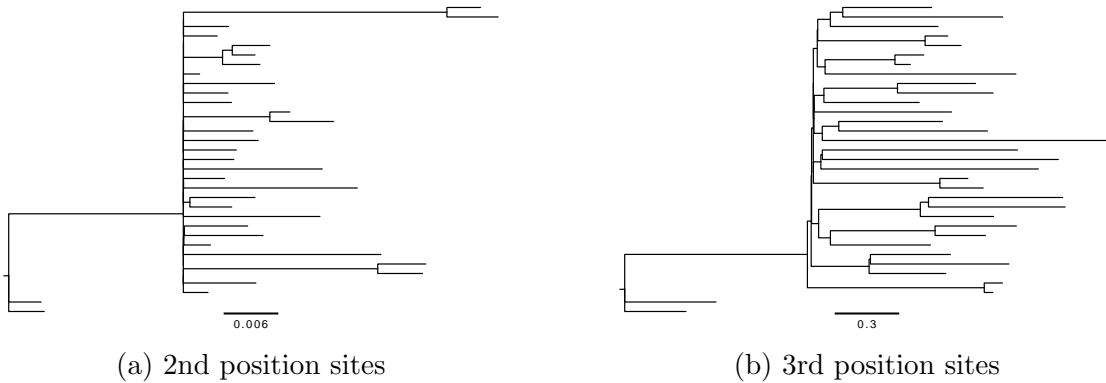


Figure 4: Fréchet mean trees computed from the *psaB* locus: (a) 2nd-codon-position sites only, and (b) 3rd-codon-position sites only.

analyses using all sites. Note, however, that the 3rd-position mean tree is also closer to the 1st-position mean tree than is the 2nd-position mean tree. These results are contradicted by Cluster Information Distance (CID; Smith, 2020), which is arguably a more sensitive topology-only tree distance than the Robinson-Foulds distance. At the very least, given the small differences between these distances, it would be difficult to argue that either 2nd-position or 3rd-position sites are saying something completely different than 1st-position

Table 1: BHV distances (BHV), Cluster Information Distances (CID), and Robinson-Foulds Distances (RFD) between mean trees from Bayesian analyses of codon-based data partitions. Each distance is a simple average of 4 distances, one from each of the two independent runs.

Comparison	BHV	CID	RF
2nd vs 1st	0.09228	0.47349	35.5
3rd vs 1st	0.09080	0.48355	33.5
2nd vs all	0.10893	0.35122	29
3rd vs all	0.05265	0.36291	25

sites.

The saturation test in PhyloMAd concluded that all data subsets (1st-codon, 2nd-codon, and 3rd-codon, as well as the entire locus), for both empirical and all simulated data sets, were ‘low risk’. This test computes a t-statistic that measures how surprised we should be at the site patterns in the data if the data subset was saturated (i.e. observed nucleotides are distributed according to the global empirical nucleotide frequencies). The critical value of the t-statistic is determined not from the t distribution, but, instead, is the value that maximizes the ratio of true positives to false positives from simulations covering a diversity of parameter combinations. The test is thus conservative, not requiring the data to be truly saturated to be declared unsound, yet in all cases the test said that the data were not close to being saturated, confirming our conclusions based on LCR.

DISSONANCE

The mean tree from a sample of 10,000 trees from the posterior distribution of the 5’ subset of *rps11* shows *Sanguinaria* properly associated with other genera (*Bocconia* and *Eschscholzia*) in its eudicot family (Poppy Family, Papaveraceae), whereas the 3’ posterior mean tree shows *Sanguinaria* associated with unrelated monocot genera (*Oryza* and *Disporum*) (Fig. 5).

0451
0452
0453
0454
0455
0456
0457
0458
0459
0460
0461
0462
0463
0464
0465
0466
0467
0468
0469
0470
0471
0472
0473
0474
0475

Two independent MCMC analyses were conducted of both 3' and 5' data subsets, allowing comparison of dissonance measured between independent samples from the same posterior (i.e. 3' vs. 3' or 5' vs. 5') to dissonance measured between samples from different posteriors (e.g. 3' vs 5'). Dissonance was consistently greater than 8 for comparisons between the 3' and 5' subsets but less than 0.2 for comparisons of independent samples from the same posterior distribution (Table 2).

Table 2: Dissonance between posterior samples (D ; lower triangle) and average information (I ; upper triangle). a and b denote independent samples of the same posterior distribution.

	3'a	3'b	5'a	5'b
3'a	—	81.51472	79.58368	79.09154
3'b	0.06490	—	79.01913	78.52700
5'a	8.31250	8.30348	—	76.59596
5'b	8.18291	8.17435	0.15654	—

0476
0477
0478
0479
0480
0481
0482
0483
0484
0485
0486
0487
0488
0489
0490
0491
0492
0493
0494
0495
0496
0497
0498
0499
0500

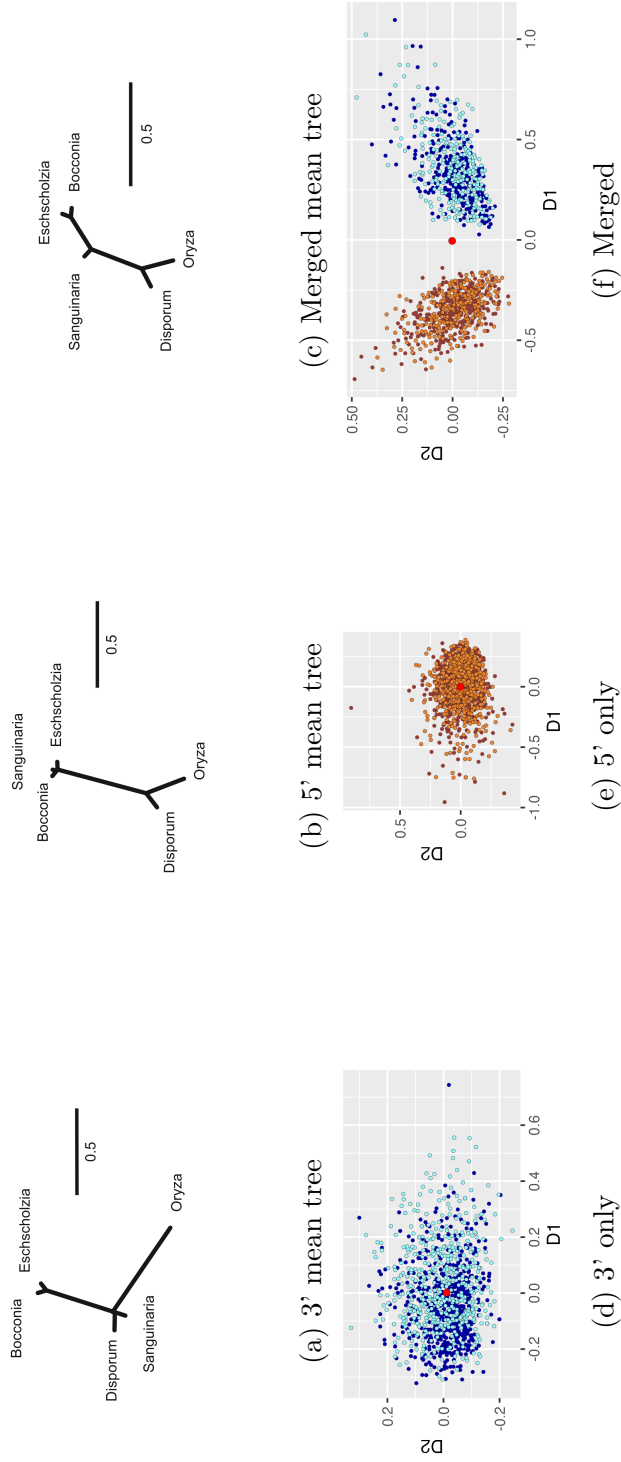


Figure 5: Comparisons of posterior samples from subsets of the *rps11* locus. (a-c) Mean trees from posterior samples conditioned on the 3' subset, the 5' subset, and the merged sample, respectively. (d-e) Multidimensional scaling (MDS) comparisons of 2 random subsamples, each of size 200, from independent analyses of the same 3' and 5' posterior distributions, respectively. (f) MDS comparison of merger of the 4 subsamples in (d) and (e). Larger red point in each plot is the mean tree.

DISCUSSION

The methods presented in this paper use the difference in phylogenetic variance between prior and posterior tree samples to make inferences about the phylogenetic information content of data. This builds on previous work by Lewis et al. (2016), who measured information content using the difference between prior and posterior entropy in discrete topological probability distributions. The primary drawback of the Lewis et al. (2016) method was scalability. The methods we present here are very different, making use of recent algorithms for efficiently determining geodesic distances between trees to estimate the mean and variance of phylogenetic trees (Owen and Provan, 2010; Brown and Owen, 2020; Miller, Owen, and Provan, 2015). The new measure yields similar inferences of information content as the entropy-based measures if the prior sample is scaled to have the same mean tree length as the posterior sample. They have the advantage of being at least as scalable as the Bayesian sampling itself; that is, if it is possible to obtain a good sample from the posterior and prior distributions, then it is possible to apply our method.

In addition to measuring information content, we have also introduced a geodesic-distance-based measure of dissonance that is large and positive if data subsets conflict in their choice of tree and small if data subsets do not conflict.

One nice feature of the approach described by Lewis et al. (2016) is that information can be partitioned. For example, it is possible to say that 90% of the total information is attributable to a single well-supported clade. This additivity of information content is not currently possible with the method described here, although future work may reveal ways to assess the relative contribution of different subsets of taxa to the total information content.

More work is also needed to determine the best definition of "volume" to use in the

0526 LCR measure. Here, we've used as volume the 95% radius: that is, the distance from the
0527 mean tree to the most distant sampled tree that nevertheless lies within the 95% set of trees
0528 closest to the mean tree. It would be more appropriate to use a measure of the true volume
0529 of treespace enclosed by the 95% set of sampled trees closest to the mean, but the non-
0530 Euclidean component of tree space as defined by Billera et al. (2001) places this definition
0531 of volume currently out of reach.

0532 One clear application of our information content measure is in phylogenomics. When
0533 data from many loci are available, it is not always advisable to include all loci (Duchêne et
0534 al., 2022). In particular, species tree methods such as BEAST2 (Bouckaert et al., 2019) that
0535 jointly estimate gene trees and the species tree could be made more computationally efficient
0536 by filtering out loci that have very little information to contribute. Species tree inference
0537 using ASTRAL, which treats gene trees as observed data, could be improved by using the
0538 mean tree from each locus as input. The mean tree is expected to be poorly resolved if
0539 information content is low, but the clades that are present in the mean tree are those that
0540 are supported by the information in the data for that locus. Thus, using the mean tree as
0541 input to a species tree method would be less likely to contribute false information than using
0542 the maximum likelihood or maximum *a posteriori* (MAP) tree and provides an alternative
0543 to collapsing nodes in gene trees using arbitrary rules.

0544 With extremely fast tests for saturation available in software such as DAMBE (Xia,
0545 2009) and PhyloMAd (Duchêne et al., 2018), why use the method described here, which is
0546 itself fast but requires a valid Bayesian posterior sample that may take hours or days to
0547 obtain? One reason is that saturation tests test only one end of the information spectrum:
0548 they do not reveal loci that exhibit low information due to a paucity of substitutions. An-
0549 other reason is that our approach makes use of the exact Bayesian model being used for
0550

0551 inference, whereas entropy t-tests such as those in PhyloMAd use critical values predicted
0552 from simulation experiments that, while exploring many relevant parameter combinations,
0553 cannot anticipate the behavior of the specific model being used. In particular, PhyloMAd
0554 may have difficulty accurately identifying problematic data sets when the model used is com-
0555 plex (Duchêne et al., 2022), for example, the Bayesian CAT model (Lartillot and Philippe,
0556 2004). Finally, our approach directly measures the information content of data by comparing
0557 the posterior distribution, which is informed by the data, to the prior distribution, which
0558 serves as a reference point not influenced by the observed data.

0560 FUNDING

0561
0562 This work was supported by the National Science Foundation Graduate Research Fellowship
0563 Program (Grant No. DGE 2136520 to AAM). Any opinions, findings, and conclusions or
0564 recommendations expressed in this material are those of the author(s) and do not necessarily
0565 reflect the views of the National Science Foundation.

0568 ACKNOWLEDGMENTS

0569
0570 We would like to thank xxxxx for their constructive comments.

0571 The computational work for this project was conducted using resources provided by
0572 the Storrs High-Performance Computing (HPC) cluster. We extend our gratitude to the
0573 UConn Storrs HPC and its team for their resources and support, which aided in achieving
0574 these results.

SUPPLEMENTARY MATERIAL

Data available from the Dryad Digital Repository: <http://datadryad.org/xxxxx>

REFERENCES

- Archie, J. W. 1989. A randomization test for phylogenetic information in systematic data. *Systematic Zoology* 38:239–252. <https://doi.org/10.2307/2992285>
- Bergthorsson, U., K. L. Adams, B. Thomason, and J. D. Palmer. 2003. Widespread horizontal transfer of mitochondrial genes in flowering plants. *Nature* 424:197–201. <https://doi.org/10.1038/nature01743>
- Berling, L., J. Klawitter, R. Bouckaert, D. Xie, A. Gavryushkin, and A. J. Drummond. 2025. Accurate Bayesian phylogenetic point estimation using a tree distribution parameterized by clade probabilities. *PLoS Computational Biology* 21:e1012789. <https://doi.org/10.1371/journal.pcbi.1012789>
- Billera, L. J., S. P. Holmes, and K. Vogtmann. 2001. Geometry of the space of phylogenetic trees. *Advances in Applied Mathematics* 27(4):733–767. <https://doi.org/10.1006/aama.2001.0759>
- Remco Bouckaert, Timothy G. Vaughan, Joëlle Barido-Sottani, Sebastián Duchêne, Mathieu Fourment, Alexandra Gavryushkina, Joseph Heled, Graham Jones, Denise Kühnert, Nicola De Maio, Michael Matschiner, Fábio K. Mendes, Nicola F. Müller, HewvA. Ogilvie, Louis du Plessis, Alex Poppinga, Andrew Rambaut, David Rasmussen, Igor Siveroni, Marc A. Suchard, Chieh-Hsi Wu, Dong Xie, Chi Zhang, Tanja Stadler, and Alexei J. Drummond. 2019. BEAST 2.5: a(n) advanced software platform for Bayesian evolutionary analysis.

0601 PLoS Computational Biology 15:e1006650. <https://doi.org/10.1371/journal.pcbi.>
0602 1006650

0603 Brown, J. M. 2014. Detection of implausible phylogenetic inferences using posterior predic-
0604 tive assessment of model fit. *Systematic Biology* 63:334–348. <https://doi.org/10.1093/>
0605 [sysbio/syu002](https://doi.org/10.1093/sysbio/syu002)

0606

0607 Brown, D. G., and M. Owen. 2020. Mean and variance of phylogenetic trees. *Systematic*
0608 *Biology* 69:139–154. <https://doi.org/10.1093/sysbio/syz041>

0609

0610 Duchêne, D. A., S. Duchêne, and S. Y. W. Ho. 2022. PhyloMAd: efficient assessment of
0611 phylogenomic model adequacy. *Bioinformatics* 34:2300–2301 <https://doi.org/10.1093/>
0612 [bioinformatics/bty103](https://doi.org/10.1093/bioinformatics/bty103)

0613

0614 Duchêne, D. A., N. Mather, C Van Der Wal, and S. Y. W. Ho. 2022. Excluding loci with
0615 substitution saturation improves inferences from phylogenomic data. *Systematic Biology*
0616 71:676–689. <https://doi.org/10.1093/sysbio/syab075>

0617

0618 Faith, D. P. 1991. Cladistic permutation tests for monophyly and nonmonophyly. *Systematic*
0619 *Zoology* 48:366–375. <https://doi.org/10.2307/2992329>

0620

0621 Felsenstein, J. 1985. Confidence limits on phylogenies: an approach using the bootstrap.
0622 *Evolution* 39: 783–791. <https://doi.org/10.1111/j.1558-5646.1985.tb00420.x>

0623

0624 Fischer, M., and M. Steel. 2009. Sequence length bounds for resolving a deep phylogenetic
0625 divergence. *Journal of Theoretical Biology* 256:247–252. <https://doi.org/10.1016/j.jtbi.2008.09.031>

0626 design criteria in phylogenetics: where to add taxa. *Systematic Biology* 56:609–622. <https://doi.org/10.1080/10635150701499563>

0627

0628 Goldman, N. 1998. Phylogenetic information and experimental design in molecular sys-
0629 tematics. *Proceedings of the Royal Society of London Series B* 265:1779–1786. <https://doi.org/10.1098/rspb.1998.0502>

0630

0631 Hillis, D. M., and Huelsenbeck, John P. 1992. Signal, noise, and reliability in molecu-
0632 lar phylogenetic analyses. *Journal of Heredity* 83:189–195. <https://doi.org/10.1093/oxfordjournals.jhered.a111190>

0633

0634

0635 Höhna, S., M. J. Landis, T. A. Heath, B. Boussau, N. Lartillot, B. R. Moore, J. P. Huelsen-
0636 beck, and F. Ronquist. 2016. RevBayes: Bayesian phylogenetic inference using graphical
0637 models and an interactive model-specification language. *Systematic Biology* 65:726–736.
0638 <https://doi.org/10.1093/sysbio/syw021>

0639

0640 Kluge, A., and S. Farris. 1969. Quantitative phyletics and the evolution of Anurans. *System-*
0641 *atic Zoology* 18:1–32. <https://doi.org/10.1093/sysbio/18.1.1>

0642

0643 Lartillot, N., and H. Philippe. 2004. A Bayesian mixture model for across-site heterogeneities
0644 in the amino-acid replacement process. *Molecular Biology and Evolution* 21:1095–1109.
<https://doi.org/10.1093/molbev/msh112>

0645

0646 Lemey, P., A. Rambout, A. J. Drummond, and M. A. Suchard. 2009. Bayesian phylogeog-
0647 raphy finds its roots. *PLoS Computational Biology* 5:e1000520. <https://doi.org/10.1371/journal.pcbi.1000520>

0648

0649 Lewis, P. O., M.-H. Chen, L. Kuo, L. A. Lewis, K. Fučíková, S. Neupane, Y.-B. Wang, and
0650 D. Shi. 2016. Estimating Bayesian phylogenetic information content. *Systematic Biology*
65:1009–1023. <https://doi.org/10.1093/sysbio/syw042>

- 0651 Lindley, D. V. 1956. On a measure of the information provided by an experiment. *The Annals*
0652 *of Mathematical Statistics* 27:986–1005. <https://doi.org/10.1214/aoms/1177728069>
- 0653 Lyons-Weiler, J., G. A. Hoelzer, and R. J. Tusch. 1996. Relative apparent synapomorphy anal-
0654 ysis (RASA). I: The statistical measurement of phylogenetic signal. *Molecular Biology and*
0655 *Evolution* 13:749–757. <https://doi.org/10.1093/oxfordjournals.molbev.a025635>
- 0656 Massingham, T., and N. Goldman. 2000. EDIBLE: experimental design and information
0657 calculations in phylogenetics. *Bioinformatics* 16:294–295. [https://doi.org/10.1093/](https://doi.org/10.1093/bioinformatics/16.3.294)
0658 [bioinformatics/16.3.294](https://doi.org/10.1093/bioinformatics/16.3.294)
- 0659 Miller, E., M. Owen, and J. S. Provan 2015. Polyhedral computational geometry for averaging
0660 metric phylogenetic trees. *Advances in Applied Mathematics* 68:51–91. [http://dx.doi.](http://dx.doi.org/10.1016/j.aam.2015.04.002)
0661 [org/10.1016/j.aam.2015.04.002](http://dx.doi.org/10.1016/j.aam.2015.04.002)
- 0662 Owen, M., and J. S. Provan. 2010. A fast algorithm for computing geodesic distances in tree
0663 spaces. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 8(1):2–13.
0664 [10.1109/TCBB.2010.3](https://doi.org/10.1109/TCBB.2010.3)
- 0665 Rambaut, A. and N. C. Grassly 1997. Seq-Gen: an application for the Monte Carlo simula-
0666 tion of DNA sequence evolution along phylogenetic trees. *Computer Applications in the*
0667 *Biosciences* 13:235–238. <https://pubmed.ncbi.nlm.nih.gov/9183526/>
- 0668 Robinson, D. R., and L. R. Foulds. 1981. Comparison of phylogenetic trees. *Mathematical*
0669 *Biosciences* 53:131–147. [https://doi.org/10.1016/0025-5564\(81\)90043-2.](https://doi.org/10.1016/0025-5564(81)90043-2)
- 0670 San Mauro, D., D. J. Gower, T. Massingham, M. Wilkinson, R. Zardoya, J. A. Cotton. 2009.
0671 Experimental Design in Caecilian Systematics: Phylogenetic Information of Mitochondrial
0672 Genomes and Nuclear rag1. *Systematic Biology* 58:425–438. [https://doi.org/10.1093/](https://doi.org/10.1093/sysbio/syp043)
0673 [sysbio/syp043](https://doi.org/10.1093/sysbio/syp043)

- 0676 Shpak, M., and G. A. Churchill. 2008. The information content of a character under a
0677 Markov model of evolution. *Molecular Phylogenetics and Evolution* 17:231–243. <https://doi.org/10.1006/mpev.2000.0846>
0678
- 0679 Shannon C.E. 1948. A mathematical theory of communication. *Bell System Technical Journal*
0680 27:379–423. <https://doi.org/10.1002/j.1538-7305.1948.tb01338.x>
0681
- 0682 Shi, X., H. Gu, and C. Field. 2008. Pattern Classification of Phylogeny Signals. *Statistical Applications in Genetics and Molecular Biology* 7:1. <https://doi.org/10.2202/1544-6115.1289>
0683
0684
- 0685 Shi, W., M.-H. Chen, L. Kuo, and P. O. Lewis. 2022. Bayesian concentration ratio and
0686 dissonance. *Bayesian Analysis* 17:817–847. <https://doi.org/10.1214/21-BA1277>
0687
- 0688 Smith, M. R. 2020. Information theoretic generalized Robinson–Foulds metrics for com-
0689 paring phylogenetic trees. *Bioinformatics* 36:5007–5013. <https://doi.org/10.1093/bioinformatics/btaa614>
0690
- 0691 Steel, M., P. J. Lockhart, and D. Penny. 1993. A frequency-dependent significance test for
0692 parsimony. *Molecular Phylogenetics and Evolution* 4:64–71. <https://doi.org/10.1006/mpev.1995.1006>
0693
- 0694 Steel, M., P. J. Lockhart, and D. Penny. 1995. Confidence in evolutionary trees from biological
0695 sequence data. *Nature* 364:440–442. <https://doi.org/10.1038/364440a0>
0696
- 0697 Sturm, Karl-Theodor. 2003. Probability measures on metric spaces of nonpositive curvature.
0698 In: P. Auscher, T. Coulhon, A. Grigoryan (Eds.), *Heat kernels and analysis on manifolds,*
0699 *graphs, and metric spaces: lecture notes from a quarter program on heat kernels, random*
0700 *walks, and analysis on manifolds and graphs.* In: *Contemp. Math.*, vol. 338, American
Mathematics Society, USA, 2003, pp. 357–390. <https://doi.org/10.1090/conm/338>

- 0701 Tippery, N. P., K. Fučíková, P. O. Lewis, and L. A. Lewis. 2012. Probing the monophyly
0702 of the Sphaeropleales (Chlorophyceae) using data from five genes. *Journal Of Phycology*
0703 48:1482–1493. <https://doi.org/10.1111/jpy.12003>
- 0704 Townsend, J. P. 2007. phylogenetic informativeness. *Systematic Biology* 56:222–231. <https://doi.org/10.1080/10635150701311362>
- 0705
0706
- 0707 Townsend, J. P., Z. Su, and Y. I. Tekle. 2012. Phylogenetic signal and noise: predicting
0708 the power of a data set to resolve phylogeny. *Systematic Biology* 61:835–849. <https://doi.org/10.1093/sysbio/sys036>
- 0709
- 0710 Xia, X., Z. Xie, M. Salemi, L. Chen, and Y. Wang. 2003. An index of substitution saturation
0711 and its application. *Molecular Phylogenetics and Evolution* 26:1–7. [https://doi.org/10.](https://doi.org/10.1016/S1055-7903(02)00326-3)
0712 [1016/S1055-7903\(02\)00326-3](https://doi.org/10.1016/S1055-7903(02)00326-3)
- 0713
- 0714 Xia, X. 2009. Assessing substitution saturation with DAMBE. Pages 613–629 in: Lemey,
0715 P., M. Salemi, and A.-M. Vandamme. (eds.). *The phylogenetic handbook: a practical*
0716 *approach to phylogenetic analysis and hypothesis testing*. Cambridge University Press.
0717 2nd edition. <https://doi.org/10.1017/CB09780511819049>
- 0718 Zhang, C., B. Rannala, and Z. Yang 2012. Robustness of compound Dirichlet priors for
0719 Bayesian inference of branch lengths. *Systematic Biology* 61:779–784. [https://doi.org/](https://doi.org/10.1093/sysbio/sys030)
0720 [10.1093/sysbio/sys030](https://doi.org/10.1093/sysbio/sys030)
- 0721
0722
0723
0724
0725