# 1 Robust analysis of phylogenetic tree space

Martin R. Smith
Department of Earth Sciences, Durham University, Lower Mountjoy, Durham, DH1 3LE, UK;
martin.smith@durham.ac.uk
Running head: Analysis of phylogenetic tree space
6

7 Abstract.—Phylogenetic analyses often produce large numbers of trees. Mapping trees' 8 distribution in 'tree space' can illuminate the behaviour and performance of search strategies, 9 reveal distinct clusters of optimal trees, and expose differences between different data sources or 10 phylogenetic methods – but the high-dimensional spaces defined by metric distances are 11 necessarily distorted when represented in fewer dimensions. Here, I explore the consequences of 12 this transformation in phylogenetic search results from 128 morphological datasets, using 13 stratigraphic congruence – a complementary aspect of tree similarity – to evaluate the utility of 14 low-dimensional mappings.

I find that phylogenetic similarities between cladograms are most accurately depicted in tree spaces derived from information-theoretic tree distances or the quartet distance. Robinson– Foulds tree spaces exhibit prominent distortions and often fail to group trees according to phylogenetic similarity, whereas the strong influence of tree shape on the Kendall–Colijn distance makes its tree space unsuitable for many purposes.

20 Distances mapped into two or even three dimensions often display little correspondence
21 with true distances, which can lead to profound misrepresentation of clustering structure.

22	Without explicit testing, one cannot be confident that a tree space mapping faithfully represents
23	the true distribution of trees, nor that visually evident structure is valid.
24	My recommendations for tree space validation and visualization are implemented in a
25	new graphical user interface in the 'TreeDist' R package.
26	Key words
27	tree distance metrics; multidimensional scaling; treespace projections; phylogenetic software.
28	
29	Phylogenetic analysis seeks to reconstruct historical relationships between evolving lineages,
30	such as species, languages or cell lines. Such analyses often identify many candidate trees,
31	making it difficult to encapsulate the underlying phylogenetic signal. Single summary trees
32	generated through consensus, compromise, or centroid methods (Wilkinson 1994; Nixon and
33	Carpenter 1996) cannot communicate information about the 'landscape' (Bastert et al. 2002) that
34	trees occupy, such as the existence of tightly defined but potentially dissimilar 'islands' or
35	'terraces' (Maddison 1991).
36	The structure of 'tree space' – formally, the metric space defined by the distances
37	between each pair of trees in a set – can help to establish the progress of tree searches; to
38	produce more informative summary trees; to reveal relationships within a set of optimal trees
39	obtained from different datasets or methods; and to illuminate the posterior distribution of trees
40	resulting from Bayesian analysis (Amenta and Klingner 2002; Stockham et al. 2002; Hillis et al.
41	2005; Holmes 2006; Chakerian and Holmes 2012; Whidden and Matsen 2015; Willis and Bell
42	2018; Wright and Lloyd 2020).
43	To appreciate this structure, a tree space that may have many intrinsic dimensions must

44 be mapped into fewer: ideally two or three. However, dimensionality reduction discards

45	information: mapping into too few dimensions will misrepresent spatial relationships. Few
46	published studies evaluate whether a mapped tree space meaningfully depicts true tree-to-tree
47	distances – perhaps because such distortion is deemed a theoretical rather than practical concern.
48	Alongside the dimensionality of a mapping, other factors known to influence the nature
49	and utility of tree space include the method of dimensionality reduction; the means of calculating
50	distances between trees; the specific trees used to generate the tree space; and how clusters
51	('islands') of trees are identified (Hillis et al. 2005; Huang and Li 2013; Wilgenbusch et al.
52	2017). The methods implemented in the popular 'TreeSetVis' and 'treespace' software packages
53	(Amenta and Klingner 2002; Jombart et al. 2017) are frequently used, but there otherwise seems
54	to be little consensus as to how a method should be selected.
55	Here I evaluate the behaviour of eight distance metrics, four clustering approaches and
56	six mapping methods in the construction and interrogation of tree spaces from 128 sets of
57	stratigraphically-calibrated cladograms (Lloyd and Wright 2020). I explore the degree to which
58	methodological decisions can materially impact the analysis and interpretation of tree space, and
59	identify recommendations for best practice.

60 Methods

Wright and Lloyd (2020) used a selection of 128 morphological datasets to demonstrate how tree space analysis can facilitate the interpretation of phylogenetic results. They estimated Bayesian trees under the Mk model of morphological evolution (Lewis 2001), partitioning datasets according to the number of observed tokens per character, and using four rate categories to describe the speed of morphological change, with each category's mean rate drawn from the quartiles of a gamma distribution. A single MCMC run was executed in 'RevBayes' (Höhna et al.

67	2016) for 300 000 generations. To minimize the risk of artefacts due to non-convergence of
68	chains, I conservatively discard the first 50% of Bayesian trees as burn-in, and sample 2 500 of
69	the remaining trees at uniform intervals to represent the posterior distribution.
70	Wright and Lloyd (2020) identified most-parsimonious trees using TNT (Goloboff and
71	Catalano 2016) under equal-weights parsimony, using exhaustive searches for datasets with < 25
72	leaves, and heuristic searches for larger datasets. I include all most parsimonious trees reported,
73	with an upper limit of 1 000 trees for each dataset.
74	I treat all trees as cladograms, discarding branch length information in order to focus
75	exclusively on the evolutionary relationships contained within each tree.
76	The underlying palaeontological datasets contain 4-88 (median: 15) terminal taxa and 8-
77	540 (median: 57) morphological characters, address a broad range of vertebrate and invertebrate
78	taxa, and are each associated with stratigraphic occurrence data from the fossil record (Lloyd and
79	Wright 2020). This broad suite of tree sets with disparate properties helps to illuminate, if
80	incompletely, the nature of tree spaces constructed from typical morphological data sets.
81	Molecular datasets are not added to this sample because they cannot be directly
82	integrated with stratigraphic information from the fossil record. Besides data type, the character
83	of tree space may also depend on factors such as the method of inference, the signal:noise ratio,
84	or the number of sites per taxon. Whilst acknowledging that certain details of the results might
85	therefore be particular to these specific datasets, this study documents the degree to which
86	methodological decisions have the potential to influence tree space analysis.

## 87 Distances

88 This study considers distances that purport to quantify the similarity of relationships between
89 cladograms: the Robinson–Foulds (RF), matching split information (MS), phylogenetic

90	information (PI), clustering information (CI), path (Pt), Kendall–Colijn (KC) and quartet (Q)
91	metrics, and a new metric (SV) derived from vector representations of trees.
92	The Robinson–Foulds (symmetric partition) distance (Robinson and Foulds 1981) counts
93	the number of splits (loosely equivalent to edges or nodes) that occur in one tree but not the
94	other, making no allowance for the existence of splits that may be almost—but not quite—
95	identical. This distance is crude: it has a low resolution, is readily saturated, and is sensitive to
96	the relocation of a single group within a tree (Steel and Penny 1993).
97	Information-theoretic distances (Smith 2020a) generalize the Robinson–Foulds distance
98	to account for the differing information content of differently sized splits, and to acknowledge
99	similarities between pairs of splits that are not quite identical. These metrics construct a
100	'matching' that pairs splits between two trees so as to maximize the amount of information that
101	all paired splits hold in common; the amount of information not held in common gives the
102	distance. The clustering, phylogenetic or matching split concepts of information capture subtly
103	different aspects of similarity between relationships.
104	The quartet distance (Estabrook et al. 1985) counts whether the relationships between
105	each possible combination of four leaves are the same or different between two trees; it has a
106	similar objective to information theoretic distances, but is slower to calculate.
107	Euclidian vector-based tree distances are the square root of the sum of squared
108	differences between explicit vector representations of trees. The path distance (Steel and Penny
109	1993) constructs a vector such that for each pair of leaves $\{i, j\}$ , the entry of the vector $e_{ij}$ is the
110	number of edges between <i>i</i> and <i>j</i> . For the Kendall–Colijn (KC, Kendall and Colijn 2016)
111	distance with $\lambda = 0$ (which discards branch length information), $e_{ij}$ denotes the number of edges
112	separating the common ancestor of $i$ and $j$ from the root; taxa whose most recent common

113 ancestor is further from the root belong to a smaller taxonomic group. Setting  $e_{ii}$  to the number 114 of leaves in the smallest bipartition split containing both *i* and *j* provides an alternative measure 115 of the size of a taxonomic group that is defined for unrooted trees; the Euclidian distance 116 between such vectors defines a metric that I term the split size vector (SV) metric. 117 The KC metric is the only metric examined that assigns significance to the position of the 118 root of a tree. To establish the degree to which annotating the position of the root influences the 119 properties of tree space, all experiments with the clustering information distance are repeated 120 with and without the root node labelled. 121 I do not consider distances that incorporate branch length information (e.g. Billera et al. 122 2001; Speyer and Sturmfels 2004; Garba et al. 2018), whilst acknowledging that these can 123 produce 'natural' tree spaces with desirable properties (Gori et al. 2016; Monod et al. 2018; 124 Garba et al. 2021). Neither do I include 'edit'-based distances, which are difficult to calculate 125 exactly, and whose approximations exhibit undesirable properties (Smith 2020a). 126 Other distances, which capture other aspects of tree similarity, might also be used as the 127 basis for tree space construction: leaf-to-leaf distances (e.g. Leigh et al. 2011) emphasize branch 128 lengths over relationships; shape metrics (e.g. Mir et al. 2013; Colijn and Plazzotta 2018) 129 consider aspects of tree shape but not relationship information. As these distances do not denote 130 the similarity in the evolutionary relationships implied by cladograms in any straightforward 131 sense, I do not consider them further. 132 I have previously evaluated a number of tree distance metrics in their ability to assign 133 higher distances to cladograms that denote increasingly different evolutionary relationships 134 (Smith 2020a). In summary, these tests evaluate whether tree distances exhibit the following 135 desirable properties:

136	•	Moving a single subtree a greater distance results in a greater distance to the resulting
137		tree ('length of move');
138	•	Moving a small subtree represents a smaller change than moving a larger subtree the
139		same distance ('number of leaves moved');
140	•	Few pairs of trees exhibit the maximum possible distance ('saturation');
141	•	Few pairs of trees are allocated identical distance values ('sensitivity');
142	•	Tree shape is not correlated with tree distance ('shape independence');
143	•	Simulated clusters of trees can be recovered ('cluster recovery');
144	•	Trees inferred from progressively more degraded datasets are further from the reference
145		topology used to generate the pristine dataset, whether datasets are degraded by
146		subsampling characters ('bullseye subsampling') or by switching character states
147		between leaves ('bullseye miscoding');
148	•	Trees separated by more subtree pruning and regrafting rearrangements tend to exhibit
149		greater distances ('SPR rearrangement');
150	•	Random tree pairs exhibit a consistent score ('random distances interquartile range').
151	The pr	esent study evaluates the KC and SV metrics against these criteria (detailed in full in
152	Smith	(2020a)), and against a new benchmark designed to explore the sensitivity of metrics to
153	differe	nces in tree balance. This new 'balance independence' test uses 10 000 pairs of 25-leaf
154	trees d	rawn from a uniform distribution. I calculate the distance between each pair of trees using
155	each d	istance metric, and the degree of balance for each tree using the total cophenetic index

156 (TCI, Mir et al. 2013), using R function TreeTools::TotalCopheneticIndex() (Smith

157 2019a). Low TCI values denote a balanced tree, in which the left and right children of each node

158 exhibit an equal number of descendants. A lack of correlation  $(r^2)$  between a metric distance and

159 the difference in TCI values indicates that a metric is independent of tree balance.

#### 160 *Clustering*

- 161 I identify clusters of unique tree topologies using:
- the Hartigan–Wong K-means algorithm (Hartigan and Wong 1979, R function
- 163 kmeans()), with 3 random starts and up to 42 iterations;
- partitioning around medoids (cluster::pam(), Maechler et al. 2019), using 3 random
   starts and the algorithmic shortcuts of Schubert and Rousseeuw (2021);
- hierarchical clustering with minimax linkage (Murtagh 1983)
- 167 (protoclust::protoclust(), Bien and Tibshirani 2011) (chosen after outperforming
- 168 other linkage methods in initial informal analyses); and
- spectral clustering (using custom function TreeDist::SpectralEigens() alongside
   cluster::pam()).

171 I use silhouette coefficients to calculate the optimal clustering method and number of clusters for

172 each analysis (after Kaufman and Rousseeuw 1990). The silhouette value of a given tree

- 173 compares its cohesion its distance from each other tree within its cluster with its separation –
- 174 its distance from each tree that is not within its cluster. Values close to +1 denote a high
- 175 proximity to other trees within its cluster; values close to -1 indicate proximity to trees in other
- 176 clusters. The silhouette coefficient is the mean silhouette value of all trees. Following Kaufman

177	and Rousseeuw (1990), I interpret silhouette coefficients greater than 0.7 as representing 'strong'
178	structure; $> 0.5$ as 'reasonable' structure; $> 0.25$ as 'weak structure that may not be genuine'; and
179	< 0.25 as lacking clustering structure.
180	Clusterings (i.e. assignments of trees to clusters) are compared using their variation of
181	information (VI, Meilă 2007). Similar clusterings exhibit a low VI: the cluster to which a tree
182	belongs in one clustering strongly predicts which cluster it belongs to in the other. The VI of two
183	clusterings that each divide objects into two equally sized clusters will range from zero to two;
184	the maximum possible VI decreases if clusters are uneven in size, and increases where more
185	clusters are present in a clustering.
186	To evaluate whether clustering structure is preserved after mapping to two dimensions, I
187	consider all tree sets with 'reasonable' clustering structure (silhouette coefficient > 0.5). I
188	selected two mapping approaches – PCoA and t-SNE – for detailed (and computationally
189	expensive) investigation on the basis of preliminary analyses. After computing clusterings from
190	distances mapped into two dimensions, I record any change to the number of clusters, and
191	calculate the VI between clusterings computed from original and mapped distances.
192	Mapping
193	Distances are calculated using the R (R Core Team 2021) packages 'TreeDist' (Smith 2020b) and
194	'Quartet' (Sand et al. 2014; Smith 2019b) and mapped into 1-12 dimensions using a suite of
195	multidimensional scaling (MDS) approaches:
196	• principal coordinates analysis (PCoA, also termed classical MDS) (Gower 1966; R

197 function stats::cmdscale(), R Core Team 2021);

198	• non-metric MDS with a Kruskal-1 stress function (Kruskal 1964) (MASS::isoMDS(),
199	Venables and Ripley 2002);
200	• Sammon's (1969) metric non-linear mapping (MASS::sammon(), Venables and Ripley
201	2002);
202	• curvilinear components analysis (CCA) (Demartines and Herault 1997; Sun et al. 2013)
203	(ProjectionBasedClustering::CCA(), Thrun et al. 2020), another metric MDS
204	method;
205	• diffusion mapping (Coifman and Lafon 2006) (diffusionMap::diffuse(), Richards
206	and Cannoodt 2019);
207	• Laplacian eigenmapping (Belkin and Niyogi 2003) (dimRed::embed(), Kraemer et al.
208	2018), a kernel eigenmap method; and
209	• t-distributed stochastic neighbour embedding (van der Maaten and Hinton 2008; van der
210	Maaten 2014) (Rtsne::Rtsne(), Krijthe 2015).
211	PCoA is a simple approach which essentially rotates a high-dimensional space such that as much
212	of the variance of the data as possible falls within the plotted dimensions (Thrun 2018). PCoA
213	requires Euclidean distances, and converting distances between phylogenetic trees into a
214	Euclidean space entails a loss of information (Nye 2011). To make the distances Euclidian, I
215	follow the standard practice of adding a constant to each distance (Cailliez 1983; Jombart et al.
216	2017), whilst noting that this might distort the relative magnitude of individual distances.
217	Kruskal-1 and Sammon MDS mappings minimize the normalized difference between
218	original and mapped distances, each using a separate stress function to quantify the normalized

219 difference. In the usual case where tree distances are metrics, Sammon MDS is expected to 220 closely resemble PCoA (Ekman and Blaalid 2011) – though it can emphasize accuracy in shorter 221 rather longer distances (van der Maaten et al. 2009), providing a clearer depiction of local 222 geometric features such as separation between clusters (Thrun 2018). CCA uses a stress function that implicitly assigns points in a high number of dimensions 223 224 to locations on a 'manifold' that can be readily represented in fewer dimensions - akin to 225 reconstructing original two-dimensional distances on a sheet of paper that has since been 226 crumpled into a three-dimensional ball. This is accomplished with a stress function that penalizes 227 distortion in distances that are short when mapped (*contra* short *original* distances, as in the 228 Sammon stress function), allowing longer distances to deform more readily. The length scale that 229 qualifies as 'short' decreases as the mapping is refined. 230 Diffusion mapping is a different manifold-learning approach. Rather than minimizing a 231 stress function, trees are represented as nodes on a graph, with each node connected to others by 232 edges whose lengths are a function of the distances between trees. A Markov chain constructed 233 over this graph generates a transition matrix; treating the eigenvectors of this Markov matrix as 234 coordinates results in a low-dimensional space that, when successful, captures the main structure 235 of the data, in particular preserving the spatial relationships of near neighbours (Coifman and 236 Lafon 2006). Laplacian eigenmapping is a special case of diffusion mapping that emphasizes the 237 influence of local density on the mapping, in part by connecting trees only to a number (here, 50) 238 of their nearest neighbours in the initial graph; it is considered particularly appropriate when data 239 contain meaningful clusters (Belkin and Niyogi 2003).

Finally, t-distributed stochastic neighbour embedding (t-SNE) constructs a probability
distribution whereby trees that lie close to a specified tree are more probable. A low-dimensional

mapping is selected in order that the equivalent treatment of mapped distances replicates thisprobability distribution as closely as possible.

#### 244 Distortion

245 To evaluate the susceptibility of a tree space to distortion on mapping, I calculate its correlation 246 dimension (Camastra and Vinciarelli 2002), a measure of its intrinsic dimensionality – that is, the 247 number of dimensions necessary in order to reproduce all the structure present in the tree space. 248 I evaluate the distortion of mappings using the product of the trustworthiness and 249 continuity metrics (Venna and Kaski 2001; Kaski et al. 2003), calculated using R package 250 'dreval' with k = 10 nearest neighbours. Trustworthiness measures the degree to which points 251 that are nearby in a mapping are truly close neighbours; continuity, the extent to which points 252 that are truly nearby retain their close spatial proximity when mapped. Their product gives a 253 composite score that encapsulates both aspects of quality. I also calculate the strength of correlation (Pearson's  $r^2$  and Kendall's  $\tau$ ) between original and mapped distances, which 254 255 corresponds to the goodness of fit of a Shepard (1962) plot. Pearson's  $r^2$  measures the degree to 256 which the original distance can be predicted from the magnitude of the mapped distance: it will 257 be zero if mapped distances are random with respect to original distance, and one where the ratio 258 between any two distances is identical before and after mapping. Kendall's  $\tau$  considers only the 259 ranking of distances; where  $\tau = 1$ , tree pairs will be ranked in the same sequence whether sorted 260 by original or mapped distances.

I graphically depict stress by plotting the minimum spanning tree (MST, Gower and Ross 1969) – the shortest path connecting all trees – for 350 trees uniformly selected from the list of all Bayesian and parsimony results. Tortuous paths indicate distortion in a mapping (Anderson 1971). To quantify the distortion thus shown, I calculate the 'MST extension factor,' which I

define as the ratio between the mapped length of the MST and the shortest length possible for
each mapping (i.e. the length of the MST calculated from mapped distances); in the absence of
distortion, this ratio will be unity.
The adequacy of PCoA mappings can be further evaluated by calculating the proportion

269 of variation retained, or through visual examination of scree plots (Jolliffe 2002); these

approaches were not systematically applied in this study.

## 271 Stratigraphic congruence

272 The distribution of fossil taxa in the stratigraphic record is independent of their morphology,

except insofar as both represent a single historical record of evolution (Sansom et al. 2018). The stratigraphic congruence of trees ought therefore to be reflected in the structure of any space that fully reflects the nature of the evolutionary histories implied by its constituent trees, even where the data used to assess stratigraphic fit are not used to construct the space.

277 Wright and Lloyd (2020) quantified stratigraphic congruence with the minimum implied 278 gap (MIG) statistic, calculated using fossil occurrence data from the Palaeobiology Database 279 after rooting each tree on a manually specified outgroup taxon. A 'gap' in the fossil record is a 280 period of time in which a taxon is inferred to exist, but is not represented by fossils. The MIG is 281 the sum of gaps across all edges, when each node is situated at the time that minimizes gaps. A 282 small MIG denotes a good fit with the stratigraphic record, and by implication an increased 283 likelihood that a tree faithfully represents evolutionary history. To establish the extent to which 284 mappings of tree space portray stratigraphic structure, I calculate the cumulative proportion of variance (adjusted  $r^2$ ) of stratigraphic consistency predicted by the first one to twelve 285 286 dimensions of each mapping.

287 Results

Six-dimensional mappings for each dataset, tree distance method and mapping method, with
evaluation of clusterings and depiction of stratigraphic fit, are provided in the online

290 supplementary material (Smith 2021). Results obtained under the clustering information distance

291 when trees were rooted do not materially differ from those when trees are treated as unrooted

292 (Smith 2021).

293 *Tree distance metric* 

The results of the tests devised to compare tree distances by Smith (2020a), plus the new

<sup>295</sup> 'balance independence' test, are presented in Table 1.

Of the metrics examined, only the quartet and information theoretic tree distances consistently reflect differences in the evolutionary relationships within trees (Table 1). Relative to these distances, Euclidian vector-based distances – the path, Kendall–Colijn and split size vector metrics – do a poor job of representing pre-defined structures in sets of trees. They are less effective at identifying known clusters of trees (Table 1, 'cluster recovery'), and more often fail to assign greater distances to trees that are increasingly far from a reference tree (Table 1, 'length of move,' 'bullseye' and 'SPR rearrangement' tests).

The KC metric places a particular emphasis on differences in tree shape ( $r^2 = 0.35$  for eight-leaf trees; see Table 1, 'shape independence'), and thus downplays differences in the relationships between labelled leaves; 38% of the variation in the KC score between pairs of 25leaf trees can be attributed to differences in the degree of balance (Table 1, 'balance independence'), compared to < 3% for all other studied distances. The sensitivity of the KC metric to properties of trees that take no account of which leaf is which curtails its ability to

discriminate trees based on the evolutionary relationships they imply, reducing its relevance tophylogenetic questions.

The SV metric outperforms the KC metric against all but two of the examined benchmarks, but still performs poorly relative to the quartet distance and information theoretic distances (Table 1). As such, it is difficult to see a clear case for using Euclidian vector-based distances, whose values have no straightforward interpretation, to measure the phylogenetic similarity of trees.

316 Because different metrics capture different aspects of tree similarity, the tree spaces they 317 define can exhibit very different properties. The strong connection between tree shape and the 318 KC distance means that differences in the degree of tree balance are often the primary feature of 319 KC tree spaces, but do not characterize spaces constructed using other metrics (e.g. Fig. 1d-f). 320 Mappings of Robinson–Foulds spaces often stand out as particularly different to those of other 321 spaces; in many cases, the underlying RF space lacks structures, such as clusters and correlation 322 with stratigraphic fit, that are present in all other tree metric spaces (Fig. 1a–c, g–l; Smith 2021); though in other cases (e.g. Fig. 1d), RF mappings exhibit structure that is not evident in other 323 324 spaces.

#### 325 *Clusters*

If a dataset displays genuine clustering structure, then it is desirable for clusters to be clearly distinguished. Tree spaces constructed on the quartet, KC and SV metrics exhibit the most prominent clusters, whereas clustering is least defined in RF tree spaces (Fig. 2a). Better-defined clusters exhibit a higher silhouette coefficient, increasing the number of cases in which 'reasonable' clustering structure (silhouette coefficient > 0.5) can be identified (Fig. 2c). For tree spaces that exhibit 'reasonable' clustering, the clustering solution identified is very similar (VI  $\leq$ 

332 0.01 bits) under all distance metrics except the quartet, KC and SV metrics (VI with each other 333 metric ≥ 0.03, ≥ 0.02 and ≥ 0.01 bits respectively) (Fig. 2e).

The highest silhouette coefficients are typically obtained with hierarchical clustering (Fig. 2g, h). Where 'reasonable' structure is present, K-means and PAM tend to produce similar results to each other (VI = 0.0094 bits) and to hierarchical clustering (VI = 0.011 bits) (Fig. 2f). Spectral clustering tends to resolve clusterings that are somewhat different from those of other methods (VI  $\ge$  0.021 bits) (Fig. 2f), often with lower silhouette coefficients; these clusters often fall below the threshold for 'reasonable' structure, even in some instances where 'strong' structure (silhouette coefficient > 0.7) is recovered by other methods.

## 341 *Effects of mapping*

342 The degree of clustering is often exaggerated in two-dimensional mappings of tree space.

343 Silhouette coefficients on clusterings calculated from mapped distances are typically higher by 344 around 0.25 (Fig. 2a–d), with the effect that 'weak' structure in the original tree space often 345 appears 'reasonable' when mapped, and 'reasonable' structure often appears 'strong' (for an 346 extreme example see Fig. 1d, noting how the minimum spanning tree hints at a discrepancy

347 between mapped and original distances).

This said, the existence of 'reasonable' structure in the original tree space does not guarantee that clustering will be evident in a 2D mapping. Of the 116 tree spaces (11% of 128 datasets × eight distance metrics) with at least 'reasonable' clustering structure, 19 2D PCoA mappings and 66 2D t-SNE mappings display no more than 'weak' structure, meaning that genuine clusters cannot be distinguished.

Even where clustering structure exists in both the original tree space and its twodimensional mapping (97 PCoA mappings; 50 t-SNE mappings), dimensionality reduction often

355 changes the composition of clusters markedly (VI > 0.25 bits in 39% of PCoA and 94% of t-SNE 356 mappings) (Fig. 1g-i). Clusterings are identical in only 53% of PCoA and 6% of t-SNE 357 mappings (Fig. 3a–b). Changes in cluster composition are particularly pronounced in mappings 358 of the Euclidian vector-based and quartet distances, and in PCoA mappings of RF distances (Fig. 359 3a-b). 360 More broadly, tree spaces defined by different metrics have different propensities for 361 mapping. Mappings of RF tree spaces exhibit greater distortion than mappings of other spaces, 362 reflected by lower trustworthiness and continuity metrics, higher stress, more extended minimum 363 spanning trees, and less correlation with original distances (Fig. 3c-e). To obtain a trustworthy 364 and continuous mapping of RF distances, it is often necessary to plot at least one dimension 365 more than with other distance metrics (Fig. 3c). 366 Conversely, KC tree space, and to a lesser extent the quartet, path and SV spaces, can be 367 mapped in a more trustworthy and continuous fashion than information theoretic tree spaces, 368 often attaining the same degree of distortion with one or even two fewer dimensions (Fig. 3c) – 369 reflected by lower stress, less extension of the minimum spanning tree, and a higher correlation 370 between original and mapped distances (Fig. 3c-e). 371 Though mappings of KC, SV and quartet tree spaces are the most faithful to the original 372 distances, these mappings tend to exhibit a lower intrinsic dimensionality (Fig. 3g) and, for the 373 SV and quartet spaces, a correspondingly weaker correlation with stratigraphic congruence (Fig. 374 3f) – suggesting that the improved mapping may reflect a simpler original tree space that fails to 375 represent certain aspects of tree similarity.

376 *Mapping method* 

In most cases, PCoA, Kruskal-1 and Sammon mappings of tree space differ only in small details,
a recurrent theme being that Sammon maps often contain outliers plotted far from the majority of
trees (Smith 2021). These methods consistently attain the highest correlation with the original
distances and stratigraphic congruence, and high levels of trustworthiness and continuity (Fig. 4),
indicating that these methods map the original tree space with the least distortion.
The lower correlation between other methods and original distances reflects their
different motivations – for example, contraction of large distances may be seen as justified if it

allows a clearer mapping of spatial relationships on a more local scale. In the case of t-SNE, this

385 trade-off results in mappings with higher trustworthiness and continuity and with less-extended

386 MSTs. The opposite is true for Laplacian eigenmapping, diffusion mapping or CCA. t-SNE,

387 Laplacian eigenmapping and diffusion mapping each exhibit prominent and idiosyncratic

388 structure (which may or may not correspond to structure in the original tree space), whereas a

389 typical CCA map simply depicts a separate, approximately hyperspherical cloud corresponding

390 to each 'reasonable' cluster, with no clear evidence of any further structure (Smith 2021).

### 391 Number of dimensions

392 RF, path and information theoretic tree spaces have particularly high intrinsic dimensionalities 393 (median  $\approx 5$ ; Fig. 3g). Correspondingly, in the great majority of datasets considered herein, two-394 dimensional mappings exhibit low (<<0.95) trustworthiness and continuity values. Mapping 395 additional dimensions depicts distances more accurately and often reveals additional structure 396 (Figs 4, 5): it is not uncommon for a single high dimension of tree space to account for 50% of 397 the variance in stratigraphic fit (Fig. 5).

398 In contrast, the lower dimensionalities of quartet (3.8), KC (2.7) and SV (2.5) tree spaces 399 indicate that these spaces might often be mapped to three or even two dimensions with little 400 distortion. But even with these metrics, two dimensions are enough to produce mappings with 401 high values (> 0.95) of trustworthiness and continuity only where the number of distinct tree 402 topologies within the tree space is minimal (< 30). A third dimension is enough to attain these 403 values only in a minority of cases, and never in datasets containing trees with twenty or more 404 leaves; the majority of analyses require at least four to five dimensions for a trustworthy and 405 continuous representation (Figs 3c, 4a).

The intrinsic dimensionality of a space also reflects properties of the datasets under 406 407 examination. Under all distance metrics, it is negatively correlated with the log ratio of the number of characters to the number of taxa ( $r^2 = 0.1 - 0.24$ ,  $p < 10^{-3}$ ; Supp. Fig. 1). 408 Dimensionality correlates positively with the total number of taxa in RF space ( $r^2 \approx 0.1$ , p < 0.1409  $10^{-3}$ ), and negatively in quartet space ( $r^2 \approx 0.1, p < 10^{-3}$ ), but displays no significant 410 411 correlation in other metric spaces. Tree space dimensionality is positively correlated with the number of unique trees under the RF ( $r^2 = 0.24$ ,  $p < 10^{-3}$ ) and information theoretic ( $r^2 =$ 412 0.11 - 0.14,  $p < 10^{-3}$ ) distances, and (more weakly) the path and KC distances ( $r^2 < 0.03$ , p =413 0.03 - 0.04); but no such correlation exists under the SV and Quartet distances. 414

415 DISCUSSION

416 When analysing the distribution of phylogenetic trees, three decisions prove highly

417 consequential: the distance metric used to construct a tree space; how clusters are identified; and

418 how tree space is visualized.

419 *Distance metric* 

Tree spaces are defined with reference to an underpinning distance metric. Fundamentally, a distance metric should afford smaller distances to trees that are more similar with respect to the properties under consideration – different metrics can impose profoundly different tree spaces (Fig. 1), so a tree space will only be illuminating if its underlying metric is relevant to its application.

425 *Robinson–Foulds spaces* 

The properties of the RF distance that produce poor performance in a range of practical settings [Table 1; Steel and Penny (1993); Smith (2020a)] are particularly relevant to the construction of tree spaces: its low resolution imposes an over-quantized and thus 'gappy' space; its ready saturation means that even quite similar trees can be assigned the maximum distance; and its sensitivity to the relocation of a single group or leaf means that a subset of otherwise similar trees will be allocated unrepresentatively large distances.

In part, the high intrinsic dimensionality of RF tree spaces (Fig. 3g) reflects the distortions necessary to accommodate these phenomena. Correspondingly, the RF mappings analysed in this study often contain artefacts, fail to depict structures that are apparent under other metrics, recover weaker clustering structure, and are highly distorted (Figs 1–3). As such, it is difficult to be confident that interpretations of RF tree spaces accurately represent any meaningful aspect of tree similarity.

438 *Euclidian vector-based spaces* 

439 At first blush, tree spaces defined on Euclidian vector-based metrics look like promising

440 alternatives – particularly with regard to the high fidelity of their low-dimensional mappings

441 (Figs 1–3). The particularly low intrinsic dimensionalities of the KC and SV metrics (Fig. 3g)

442 allow the majority of their tree space structure to be represented in three or even two dimensions (Fig. 3). These two metrics also stand out for the clear definition of their clustering structure 443 444 (Fig. 2a–d), even if this maps less faithfully into few dimensions (Fig. 3). 445 However, such clustering structure often fails to correspond to artificial structure known 446 to characterize the true distribution of trees (Table 1, 'cluster recovery'). Lower intrinsic 447 dimensionalities seem to be accomplished by downplaying phylogenetic differences between 448 trees, resulting in simplistic spaces whose structures emphasize the contribution of tree shape. A 449 substantial proportion of the variance in KC distances reflects the degree of tree balance (Table 450 1), meaning that KC tree spaces are often dominated by a single dimension that discriminates 451 balanced from unbalanced trees (as in Fig. 1e), independently from how leaves happen to be 452 labelled. The contribution of tree shape to the path and SV metrics, though more nuanced, results 453 in comparable behaviour. 454 Because the relative contributions of phylogenetic and shape-based factors are not

455 explicit in the definition of these vector-based metrics, it is difficult to disentangle their
456 contribution to the structure of tree space. Consequently, Euclidian vector-based tree distances,
457 and the KC metric in particular, are poorly suited to questions of evolutionary relationships.

458 *Quartet and information-theoretic metrics* 

459 Though each have subtly different emphases, quartet and information-theoretic distances

460 increase monotonically as tree topologies undergo increasing amounts of deformation (Table 1),

461 making them inherently relevant to questions concerning the similarity of evolutionary

462 relationships between cladograms (Smith 2020a).

463 The matching split information, phylogenetic information and clustering information
464 distances produce broadly similar tree spaces with similar clustering, dimensionality and

465	mapping characteristics (Figs 1-3), so are treated together here. Quartet tree spaces exhibit a
466	more pronounced clustering structure (Fig. 2a-d) and a lower intrinsic dimensionality (Fig. 3g)
467	than information-theoretic tree spaces, meaning that they can produce more information-rich
468	maps using fewer dimensions (Fig. 3).
469	In many cases, being able to obtain a tree space that discriminates clusters more readily
470	and which requires one fewer dimension to obtain a given level of trustworthiness and continuity
471	will more than offset the slightly poorer performance of the quartet metric against the
472	benchmarks of Smith (2020a) and Table 1, and justify its significantly greater running time –
473	measured in hours rather than minutes for many of the datasets examined here. On the other
474	hand, clusters obtained using information-theoretic distances are typically rendered more
475	faithfully in mappings. Confidence that interpreted structure genuinely characterizes the
476	underlying trees will be greatest if its presence can be demonstrated in both quartet and
477	information-theoretic tree spaces, which offer complementary views on the phylogenetic
478	similarity of trees.

### 479 *Clusters*

One motivation for tree space analysis is the identification of subsets of trees that are more similar with respect to the evolutionary histories they imply. This objective is most readily met when the distance from which clusters are calculated measures that property directly. Clusters identified through the visual inspection of 2D tree space mappings will group trees according to *mapped* distances, which are an opaque function of original tree-to-tree distances, distorted in a manner that is particular to each mapping technique and influenced by all other tree-to-tree distances under consideration. Such clusters thus have no straightforward interpretation in their

487 own right, except as approximations to the clustering structure imposed by the original,

488 undistorted distances.

489 My results show that clusters derived from mapped distances are poor approximations to 490 clusters based on measured distances. In the majority of RF, SV and quartet tree spaces in which 491 'reasonable' or better structure is present in both original and mapped spaces, clusters derived 492 from original versus mapped distances differ substantially in their constitution (median VI > 0.2493 bits; Fig. 3a). Mapping a tree space into two dimensions using PCoA consistently exaggerates 494 clustering structure, causing a mean increase in silhouette coefficient of 0.3 (Fig. 2a–b) – enough 495 that maps may depict 'reasonable' or even 'strong' structure where original, undistorted distances 496 exhibit only 'weak' structure that 'may not be genuine.' Correspondingly, many mappings depict 497 multiple clusters that lack 'reasonable' support in the underpinning tree space (Fig. 2c-d).

It is therefore inadvisable to assume that clusters interpreted from two-dimensional mappings represent genuine structure. Even if such clusters sometimes happen to group trees with certain characteristics in common, it is difficult to see how they would be preferable to a clustering derived from a direct and explicit measure of those specific characteristics.

502 Where a tree space does exhibit clustering structure, a secondary objective is to assign 503 trees to clusters in a fashion that minimizes overlap between clusters, thus maximizing the 504 silhouette coefficient. Hierarchical clustering usually performs best against this criterion (Fig. 505 2h), though partitioning around medoids and K-means clustering occasionally produce the best-506 defined clustering. Differences between the clusterings recovered by different methods tend to be 507 relatively small (Fig. 2f), and which method is most appropriate will depend on the specific 508 structure within a given dataset and the emphases of the particular clustering methods: for 509 example, K-means and PAM are very effective when clusters are consistent shapes or sizes, but

can produce unexpected results when this assumption is violated (MacKay 2003; Hastie et al.2009).

512 Such factors may explain contribute to the poor performance of spectral clustering in 513 these datasets (Fig. 2f–h), despite its accurate recovery of pre-defined clusters of trees in other 514 settings (Gori et al. 2016): the geometry of these artificial tree spaces may align better with the 515 strengths of spectral clustering. Though the use of a single clustering method is unlikely to 516 mislead, the use of multiple methods provides additional opportunities to maximize the 517 silhouette coefficient, and thus to better appreciate the clustering structure of a tree set.

518 *Visualizing tree spaces* 

519 Different mapping techniques have different motivations, and thus differ markedly in the
520 structure they depict. Mapping has an order of magnitude more impact on the clustering
521 structures perceived – the easiest aspect of structure to quantify – than the measurement of tree
522 distance or the method of cluster detection (Figs 2e–f; 3a–b).

523 PCoA, Sammon and Kruskal-1 mappings have a similar philosophy: they seek to 524 minimize the stress induced by a mapping by minimizing a measure of distortion that penalizes 525 mismatches between original and mapped distances. Interpretation of such mappings is 526 straightforward: mapped distances are approximately proportional to the true distances between 527 trees. (This does not mean that mapped *areas* are proportional to original hypervolumes – see 528 Mammola (2019).) In line with this common principle, and despite the potential shortcomings of 529 PCoA (Lee and Verleysen 2007), these methods often result in very similar mappings -530 consistent with some other results from simulated and real datasets (van der Maaten et al. 2009; 531 Ekman and Blaalid 2011). As PCoA is significantly faster to calculate, its status as the most

532 widely used mapping method seems justified.

533	CCA mapping likewise seeks to minimize stress – but the cost function employed aims
534	not to faithfully reflect original distances, but rather to produce a "revealing representation" of
535	the data, with an emphasis on facilitating the visual recognition of clustering (Demartines and
536	Herault 1997). The clear depiction of clustering structure seems here to be obtained by largely
537	discarding other aspects of tree space structure. In contrast to the results of Wilgenbusch et al.
538	(2017), CCA-mapped distances exhibited lower correlation with original distances, and CCA
539	mappings exhibited lower trustworthiness and continuity than PCoA, Sammon and Kruskal-1
540	mappings (Fig. 4). This difference may reflect idiosyncrasies of the tree sets being examined: for
541	example, the Wilgenbusch et al. tree sets cluster according to the gene from which trees had been
542	inferred, so emphasizing the distinction between clusters captures relatively more of the variation
543	in tree-to-tree distances than in datasets with weak clustering structure.
544	Diffusion mappings and Laplacian eigenmappings attempt to identify a lower-
545	dimensional manifold that underlies the high dimensional space; the poor performance of
546	explicitly manifold-learning methods relative to other mapping techniques (Fig. 4) suggests that
547	the sets of optimal trees examined herein are not associated with any manifold, or sample the
548	manifold at too low a density for its inferred structure to improve mapping (Venna et al. 2010).
549	Finally, t-SNE obtains very high levels of trustworthiness and continuity, at the expense
550	of a weak correlation between mapped and original distances (Fig. 4). As such, the interpretation
551	of t-SNE mappings is not intuitive: distances between clusters, and the sizes of clusters, may not
552	be representative; and t-SNE mappings can display apparent structure in datasets known to be
553	homogeneous (Wattenberg et al. 2016). These properties characterize many of the t-SNE maps
554	generated herein (see for example studies 20, 36, 89 in Smith 2021), meaning that the capacity
555	of t-SNE mapping to represent local structure must be weighed against the danger of

misinterpretation. This risk can be reduced by exploring different values of the 'perplexity' and
'epsilon' parameters that govern the structure of t-SNE mappings, and comparing results to a
PCoA mapping.

559 Dimensionality

560 Tree spaces inherently exist in many dimensions. More trees tend to produce more complicated 561 structure with more intrinsic dimensions, as previously noted by Wilgenbusch et al. (2017), and 562 observed here under the RF and information-theoretic distances. Tree sets derived from datasets 563 with a low character: taxon ratio also tend to exhibit higher dimensionalities, perhaps because 564 matrices containing fewer characters constrain relationships less decisively: a broader posterior 565 distribution of tree topologies will encompass a larger region of tree space and thus encounter 566 more distortion when mapped (just as a map of the globe is more distorted than a cartographic 567 map of a smaller region).

The dimensionality of tree space is also influenced by the choice of distance metric (Fig. 3). Whereas tree spaces with more intrinsic dimensions have the capacity to contain more sophisticated and instructive structure, they are harder to faithfully depict in few dimensions. This trade-off does not have a natural optimum, as the utility of a tree space is not a simple function of its dimensionality.

573 KC and SV tree spaces obtain a low dimensionality by downplaying phylogenetic aspects 574 of tree similarity. Although mappings produced after discarding relevant features of tree space 575 may be less distorted, this is unlikely to compensate for the concomitant loss of information: at 576 the extreme, a dimensionality of zero can be obtained by a metric that assigns all pairs of trees a 577 distance of zero.

578 On the other hand, the mere fact that more dimensions are present need not make a tree 579 space more instructive: a metric that assigns all pairs of trees a unit distance can produce a 580 meaningless tree space with many dimensions. Analogously, the low sensitivity and rapid 581 saturation of the RF distance (Table 1) mean that it allocates many tree pairs identical scores, 582 potentially increasing the number of dimensions necessary to map RF spaces without distortion 583 (Fig. 3c–e) without a corresponding gain in utility.

In contrast, the lower dimensionality of quartet tree space relative to tree spaces defined by information-theoretic distances seems not to reflect a substantial difference in how well the metrics measure the phylogenetic similarity of cladograms (Table 1). Consequently, the lower amount of distortion introduced when quartet spaces are mapped (Fig. 3c–f) is a reason to prefer this distance for visualization, so long as examination of higher dimensions of both quartet and information-theoretic spaces confirms that a low-dimensional mapping adequately summarizes structure.

Even under the quartet distance, however, the majority of datasets require more than three (median: five) dimensions to attain levels of trustworthiness and continuity greater than 0.95 (Fig. 3c). Humans are less able to perceive metric distances in three-dimensional visualizations than in two dimensions (Kjellin et al. 2010), and 3D displays are ineffective for estimating relative positions (Tory et al. 2006). Mappings that require multiple dimensions may thus fail in their objective of making distances easier to visualize.

Wilgenbusch et al. (2017) take the more optimistic position that that two-dimensional mappings capture the most important aspects of tree space structure, including clustering. In practice, I suspect that individual datasets each occupy their own position between these extremes. Although 2D maps tend to exaggerate the degree of clustering – leading to the

601 misidentification of clusters (Figs 1d, 2a–d, 3a–b), and in some cases the failure to depict aspects 602 of tree space structure that are relevant to interpretation (Fig. 5) – whatever structure is portrayed 603 by a 2D plot can at least by deciphered at a glance, in contrast to more cognitively taxing 604 portrayals of higher-dimensional space, which must still ultimately be perceived through the two 605 dimensions of the retina. The potential for misinterpretation can be reduced by plotting the 606 minimum spanning tree (e.g. Fig. 1d), by marking clusters that are statistically supported by the 607 original distances, by evaluating how well a low-dimensional mapping conveys tree-to-tree 608 distances, and by carefully examining higher dimensions for evidence of additional structure.

#### 609 Recommendations

610 In summary, commonly used practices are generally inadequate for the interpretation of the 611 phylogenetic tree spaces explored herein. The Kendall-Colijn and Robinson-Foulds metrics do 612 not directly measure trees' phylogenetic similarity; their associated tree spaces are poorly suited 613 to phylogenetic questions. Clusters identified by visual inspection of mappings are likely to 614 misrepresent the true structure of a tree set. Two dimensions are seldom sufficient to convey the 615 full structure of tree space, and 2D mappings should be viewed with suspicion unless shown to 616 exhibit high values of continuity and trustworthiness; a low correlation between original and 617 mapped distances indicates that the interpretation of a mapping may require additional care. 618 The 128 tree sets studied herein include multiple examples where standard practice would 619 lead to invalid conclusions. For instance, Wright and Lloyd (2020) correctly interpret the RF 620 space of trees from Yates (2003) (Fig. 1a; cf. fig. 1C in Wright and Lloyd (2020)) as exhibiting 621 no relationship with stratigraphic congruence – yet a strong relationship is present in all other 622 metric spaces (see Fig. 1b and Smith (2021)). Similarity, 2D mappings of Russell and Dong 623 (1993) or Xu et al. (2018) tree spaces (Fig. 5) contain no hint of the significant correlation with

624 stratigraphic congruence that exists in higher dimensions. Strong clustering in the mapped RF 625 space of trees from Fischer et al. (2016) is entirely an artefact of mapping: no corresponding 626 structure exists in the original distance matrix. These are not isolated instances, but examples that 627 illustrate recurrent patterns evident across all these studies; and there is no obvious reason that 628 the tree sets analysed here should be particularly intractable to tree space analysis. With the 629 caveat that 'landscapes' of trees selected using different optimality criteria or from different 630 sources of data may exhibit different properties, these results raise serious concerns over the 631 validity of previous presentations of tree spaces.

To minimize artefacts when analysing the distribution of cladograms, I recommend that tree space analysis employs the clustering information or quartet distances – ideally, both. These distances are sensitive to differences in the evolutionary relationships implied by cladograms, but not to factors such as tree shape that are irrelevant to most phylogenetic questions.

Information-theoretic distances (particularly the clustering information distance) measure the similarity between cladograms more effectively than the quartet distance, and have a higher intrinsic dimensionality. Insofar as this higher dimensionality denotes a more information-rich tree space, the clustering information distance is well suited to the identification of clusters of trees; moreover, these clusters tend to retain their identity when mapped. On the other hand, the lower dimensionality of quartet tree space means that it suffers less distortion when mapped. Structure evident in both metric spaces might warrant additional confidence.

Except where clustering is conducted for a separate purpose (Gori et al. 2016) – for instance, when using clusterings to generate summary trees (Stockham et al. 2002) – clusters should be identified objectively from original tree distances. The clustering with the highest silhouette coefficient can be considered the best representation of the underlying structure,

647 provided that this coefficient is high enough to indicate that the structure is meaningful (> 0.5). 648 Hierarchical clustering often finds the best clustering, but as the optimal method depends on the 649 nature of clustering structure, I encourage the use of multiple clustering methods. Depicting the 650 best clustering on mappings (as in Fig. 1g–i) reduces the potential for misinterpretation where 651 mappings do not reflect the structure of the original tree space.

The optimal mapping method will depend on the purpose of the visualization. PCoA maps – which tend to closely resemble Sammon and Kruskal mappings, but are much faster to compute – tend to reproduce original tree-to-tree distances most faithfully, making them easy to interpret, whilst also depicting structure consistently (high trustworthiness and continuity); as such, they are an obvious choice for instructive mappings. Alternatively, t-SNE maps emphasize local structural relationships, though their interpretation can be counter-intuitive; whereas CCA maps depict cluster membership whilst downplaying other structural features.

659 Whichever mapping method is employed, it is important to evaluate the quality of the 660 mapping: high (> 0.95) values of the trustworthiness and continuity measures are desirable, as is 661 a good correlation with original distance metrics. Even if minimum spanning trees can help to 662 visually assess the degree of distortion, it is not possible to be confident that any apparent 663 structure is genuine unless the quality of a mapping is explicitly documented. Of course, these 664 metrics will be invalid, and distances misrepresented, unless plotting software is configured to 665 plot *x* and *y* axes to the same scale.

666 These recommendations are drawn from a limited sample of morphological datasets; it is 667 likely that tree sets obtained from different datasets using different methods will occupy tree 668 spaces with different properties. Nevertheless, the degree to which methodological decisions can 669 influence the interpretation of tree space represents a strong argument for conducting and

- 670 documenting basic checks to establish that presented results truly represent the underlying
- 671 structure of the high-dimensional tree space.
- To facilitate best practice in the construction, evaluation and interpretation of tree space, I
- 673 have produced a 'point-and-click' graphical interface within R installed using
- 674 install.packages('TreeDist') and launched by executing the command
- 675 TreeDist::MapTrees(). This software allows users to upload trees, select tree distance,
- 676 mapping and clustering methods, and generate high-dimensional mappings, with real-time
- 677 evaluations of mapping and clustering quality to ensure that interpretations truly reflect the
- 678 underlying distribution of phylogenetic trees.

### 679 ACKNOWLEDGEMENTS

- 680 Reviews from Michelle Kendall and four anonymous referees, comments from editor Sebastian
- 681 Höhna, and discussions with Andrew Millard, Matthias Sinnesael and April Wright, stimulated
- the research and improved the manuscript. Nura Kawa provided scripts for spectral clustering.

#### 683 SUPPLEMENTARY MATERIAL

Data available from the Dryad Digital Repository: https://dx.doi.org/10.5061/dryad.kh1893240

## 685 DISCLOSURE STATEMENT

686 I disclose no conflicts of interest.

- 687 **References**
- 688 Amenta N., Klingner J. 2002. Case study: Visualizing sets of evolutionary trees. IEEE

689 Symposium on Information Visualization, 2002. INFOVIS 2002.:71–74.

- 690 Anderson A.J.B. 1971. Ordination methods in ecology. Journal of Ecology. 59:713–726.
- Bastert O., Rockmore D., Stadler P.F., Tinhofer G. 2002. Landscapes on spaces of trees. Applied
- 692 Mathematics and Computation. 131:439–459.
- 693 Belkin M., Niyogi P. 2003. Laplacian eigenmaps for dimensionality reduction and data

representation. Neural Computation. 15:1373–1396.

- Bien J., Tibshirani R. 2011. Hierarchical clustering with prototypes via minimax linkage. Journal
  of the American Statistical Association. 106:1075–1084.
- Billera L.J., Holmes S.P., Vogtmann K. 2001. Geometry of the space of phylogenetic trees.
- Advances in Applied Mathematics. 27:733–767.
- Cailliez F. 1983. The analytical solution of the additive constant problem. Psychometrika.
  48:305–308.
- 701 Camastra F., Vinciarelli A. 2002. Estimating the intrinsic dimension of data with a fractal-based
- method. IEEE Transactions on pattern analysis and machine intelligence. 24:1404–1407.
- 703 Carpenter K. 2001. Phylogenetic analysis of the Ankylosauria. The Armored Dinosaurs.
- Bloomington: Indiana University Press. p. 455–483.
- 705 Chakerian J., Holmes S. 2012. Computational tools for evaluating phylogenetic and hierarchical
- clustering trees. Journal of Computational and Graphical Statistics. 21:581–599.

707	Coifman R.R., Lafon S. 2006. Diffusion maps. Applied and Computational Harmonic Analysis.
708	21:5–30.
709	Colijn C., Plazzotta G. 2018. A metric on phylogenetic tree shapes. Systematic Biology. 67:14.
710	Demartines P., Herault J. 1997. Curvilinear component analysis: A self-organizing neural
711	network for nonlinear mapping of data sets. IEEE Transactions on Neural Networks.
712	8:148–154.
713	Ekman S., Blaalid R. 2011. The devil in the details: Interactions between the branch-length prior
714	and likelihood model affect node support and branch lengths in the phylogeny of the
715	Psoraceae. Systematic Biology. 60:541–561.
716	Estabrook G.F., McMorris F.R., Meacham C.A. 1985. Comparison of undirected phylogenetic
717	trees based on subtrees of four evolutionary units. Systematic Zoology. 34:193-200.
718	Fischer V., Bardet N., Benson R.B.J., Arkhangelsky M.S., Friedman M. 2016. Extinction of fish-
719	shaped marine reptiles associated with reduced evolutionary rates and global environmental
720	volatility. Nature Communications. 7:10825.
721	Garba M.K., Nye T.M.W., Boys R.J. 2018. Probabilistic distances between trees. Systematic
722	Biology. 67:320–327.
723	Garba M.K., Nye T.M.W., Lueg J., Huckemann S.F. 2021. Information geometry for
724	phylogenetic trees. Journal of Mathematical Biology. 82:19.
725	Goloboff P.A., Catalano S.A. 2016. TNT version 1.5, including a full implementation of
726	phylogenetic morphometrics. Cladistics. 32:221–238.

727	Gori K., Suchan T., Alvarez N., Goldman N., Dessimoz C. 2016. Clustering genes of common
728	evolutionary history. Molecular Biology and Evolution. 33:1590–1605.

- 729 Gower J.C. 1966. Some distance properties of latent root and vector methods used in
- 730 multivariate analysis. Biometrika. 53:325–338.
- 731 Gower J.C., Ross G.J.S. 1969. Minimum spanning trees and single linkage cluster analysis.

Journal of the Royal Statistical Society. Series C (Applied Statistics). 18:54–64.

Hartigan J.A., Wong M.A. 1979. Algorithm AS 136: A K-means clustering algorithm. Journal of

the Royal Statistical Society. Series C (Applied Statistics). 28:100–108.

- Hastie T., Tibshirani R., Friedman J. 2009. The Elements of Statistical Learning: Data Mining,
  Inference, and Prediction, Second Edition. New York: Springer-Verlag.
- Hillis D.M., Heath T.A., St. John K. 2005. Analysis and visualization of tree space. Systematic
  Biology. 54:471–482.
- Holmes S. 2006. Visualising data. Statistical problems in particle physics, astrophysics and
   cosmology, Proceedings of PHYSTAT05. London: Imperial College Press. p. 197–208.
- 741 Höhna S., Landis M.J., Heath T.A., Boussau B., Lartillot N., Moore B.R., Huelsenbeck J.P.,
- 742 Ronquist F. 2016. RevBayes: Bayesian phylogenetic inference using graphical models and
- an interactive model-specification language. Systematic Biology. 65:726–736.
- Huang H., Li Y. 2013. MASTtreedist: Visualization of tree space based on maximum agreement
  subtree. Journal of Computational Biology. 20:42–49.
- 746 Jolliffe I.T. 2002. Principal Component Analysis. New York: Springer-Verlag.

747	Jombart T., Kendall M., Almagro-Garcia J., Colijn C. 2017. Treespace: Statistical exploration of
748	landscapes of phylogenetic trees. Molecular Ecology Resources. 17:1385–1392.
749	Kaski S., Nikkilä J., Oja M., Venna J., Törönen P., Castrén E. 2003. Trustworthiness and metrics
750	in visualizing similarity of gene expression. BMC Bioinformatics. 4:48.
751	Kaufman L., Rousseeuw P.J. 1990. Partitioning around medoids (Program PAM). Finding
752	groups in data: An introduction to cluster analysis. John Wiley & Sons, Ltd. p. 68–125.
753	Kendall M., Colijn C. 2016. Mapping phylogenetic trees to reveal distinct patterns of evolution.
754	Molecular Biology and Evolution. 33:2735–2743.
755	Kjellin A., Pettersson L.W., Seipel S., Lind M. 2010. Evaluating 2D and 3D visualizations of
756	spatiotemporal information. ACM Transactions on Applied Perception. 7:1–23.
757	Kraemer G., Reichstein M., Mahecha M.D. 2018. dimRed and coRanking—unifying
758	dimensionality reduction in R. The R Journal. 10:342–358.
759	Krijthe J.H. 2015. Rtsne: t-distributed stochastic neighbor embedding using a Barnes-Hut
760	implementation.
761	Kruskal J.B. 1964. Multidimensional scaling by optimizing goodness of fit to a nonmetric
762	hypothesis. Psychometrika. 29:1–27.
763	Lee J.A., Verleysen M. 2007. Nonlinear dimensionality reduction. Springer Science & Business
764	Media.

- Leigh J.W., Schliep K., Lopez P., Bapteste E. 2011. Let them fall where they may: Congruence
  analysis in massive phylogenetically messy data sets. Molecular Biology and Evolution.
  28:2773–2785.
- Lewis P.O. 2001. A likelihood approach to estimating phylogeny from discrete morphological
  character data. Systematic Biology. 50:913–925.
- Lloyd G.T., Wright A.M. 2020. Data from: Bayesian analyses in phylogenetic palaeontology:
  Interpreting the posterior sample. Dryad, doi:10.5061/dryad.zs7h44j4w.
- van der Maaten L.J.P. 2014. Accelerating t-SNE using tree-based algorithms. Journal of Machine
  Learning Research. 15:3221–3245.
- van der Maaten L.J.P., Hinton G.E. 2008. Visualizing high-dimensional data using t-SNE.

Journal of Machine Learning Research. 9:2579–2605.

- van der Maaten L.J.P., Postma E., van den Herik J. 2009. Dimensionality reduction: A
- comparative review. Journal of Machine Learning Research. 10:66–71.
- MacKay D.J.C. 2003. Information Theory, Inference, and Learning Algorithms. Cambridge:
  Cambridge University Press.
- Maddison D.R. 1991. The discovery and importance of multiple islands of most-parsimonious
  trees. Systematic Biology. 40:315–328.
- Maechler M., Rousseeuw P., Struyf A., Hubert M., Hornik K. 2019. Cluster: Cluster analysis
  basics and extensions.

- Mammola S. 2019. Assessing similarity of *n*-dimensional hypervolumes: Which metric to use?
  Journal of Biogeography. 46:2012–2023.
- Meilă M. 2007. Comparing clusterings—an information based distance. Journal of Multivariate
  Analysis. 98:873–895.
- 788 Mir A., Rosselló F., Rotger L.A. 2013. A new balance index for phylogenetic trees.
- 789 Mathematical Biosciences. 241:125–136.
- 790 Monod A., Lin B., Yoshida R., Kang Q. 2018. Tropical geometry of phylogenetic tree space: A

statistical perspective. arXiv.:1805.12400.

- Murtagh F. 1983. A survey of recent advances in hierarchical clustering algorithms. The
  Computer Journal. 26:354–359.
- Nixon K.C., Carpenter J.M. 1996. On consensus, collapsibility, and clade concordance.
- 795 Cladistics. 12:305–321.
- Nye T.M.W. 2011. Principal components analysis in the space of phylogenetic trees. The Annals
  of Statistics. 39:2716–2739.
- R Core Team. 2021. R: A language and environment for statistical computing.
- 799 Richards J., Cannoodt R. 2019. diffusionMap: Diffusion map.
- Robinson D.F., Foulds L.R. 1981. Comparison of phylogenetic trees. Mathematical Biosciences.
  53:131–147.

802	Russell D.A., Dong ZM. 1993. The affinities of a new theropod from the Alxa Desert, Inner			
803	Mongolia, People's Republic of China. Canadian Journal of Earth Sciences. 30:2107–2127.			
804	Sammon J.W. 1969. A nonlinear mapping for data structure analysis. IEEE Transactions on			
805	Computers. C-18:401–409.			
806	Sand A., Holt M.K., Johansen J., Brodal G.S., Mailund T., Pedersen C.N.S. 2014. tqDist: A			
807	library for computing the quartet and triplet distances between binary or general trees.			
808	Bioinformatics. 30:2079–2080.			
809	Sansom R.S., Choate P.G., Keating J.N., Randle E. 2018. Parsimony, not Bayesian analysis,			
810	recovers more stratigraphically congruent phylogenetic trees. Biology Letters.			
811	14:20180263.			
812	Schoch R.R., Milner A.R. 2008. The intrarelationships and evolutionary history of the			
813	temnospondyl family Branchiosauridae. Journal of Systematic Palaeontology. 6:409-431.			
814	Schubert E., Rousseeuw P.J. 2021. Fast and eager k-medoids clustering: O(k) runtime			
815	improvement of the PAM, CLARA, and CLARANS algorithms. Information Systems.			
816	101:101804.			
817	Shepard R.N. 1962. The analysis of proximities: Multidimensional scaling with an unknown			
818	distance function. II. Psychometrika. 27:219–246.			
819	Smith M.R. 2019a. TreeTools: Create, modify and analyse phylogenetic trees. Comprehensive R			
820	Archive Network. doi: 10.5281/zenodo.3522725.			

821	Smith M.R. 2019b. Quartet: Comparison of phylogenetic trees using quartet and split measures.
822	Comprehensive R Archive Network. doi:10.5281/zenodo.2536318.
823	Smith M.R. 2020a. Information theoretic Generalized Robinson–Foulds metrics for comparing
824	phylogenetic trees. Bioinformatics. 36:5007–5013.
825	Smith M.R. 2020b. TreeDist: Distances Between Phylogenetic Trees. Comprehensive R Archive
826	Network. doi:10.5281/zenodo.3528123.
827	Smith M.R. 2021. Six-dimensional tree space projections of Wright and Lloyd (2020) datasets.
828	Dryad, doi:10.5061/dryad.kh1893240.
829	Speyer D., Sturmfels B. 2004. The tropical Grassmannian. Advances in Geometry. 4:389–411.
830	Steel M.A., Penny D. 1993. Distributions of tree comparison metrics—some new results.
831	Systematic Biology. 42:126–141.
832	Stockham C., Wang LS., Warnow T. 2002. Statistically based postprocessing of phylogenetic
833	analysis by clustering. Bioinformatics. 18:S285–S293.
834	Sun J., Crowe M., Fyfe C. 2013. Incorporating visualisation quality measures to curvilinear
835	component analysis. Information Sciences. 223:75-101.
836	Thrun M.C. 2018. Methods of projection. Projection-Based Clustering through Self-Organization
837	and Swarm Intelligence. Wiesbaden: Springer Fachmedien Wiesbaden. p. 33-42.
838	Thrun M., Lerch F., Pape F., Schreier T., Winckelmann L. 2020. ProjectionBasedClustering:
839	Projection based clustering.

- 840 Tory M., Kirkpatrick A.E., Atkins M.S., Moller T. 2006. Visualization task performance with
- 2D, 3D, and combination displays. IEEE Transactions on Visualization and Computer
  Graphics. 12:2–13.
- 843 Venables W.N., Ripley B.D. 2002. Modern applied statistics with S. New York: Springer.
- 844 Venna J., Kaski S. 2001. Neighborhood preservation in nonlinear projection methods: An

experimental study. Artificial Neural Networks ICANN 2001.:485–491.

- 846 Venna J., Peltonen J., Nybo K., Aidos H., Kaski S. 2010. Information retrieval perspective to
- 847 nonlinear dimensionality reduction for data visualization. Journal of Machine Learning848 Research. 11:40.
- 849 Wattenberg M., Viégas F., Johnson I. 2016. How to use t-SNE effectively. Distill.
- 850 Whidden C., Matsen F.A. 2015. Quantifying MCMC exploration of phylogenetic tree space.
- 851 Systematic Biology. 64:472–491.
- Wilgenbusch J.C., Huang W., Gallivan K.A. 2017. Visualizing phylogenetic tree landscapes.
  BMC Bioinformatics. 18:85.
- Wilkinson M. 1994. Common cladistic information and its consensus representation: Reduced
- Adams and reduced cladistic consensus trees and profiles. Systematic Biology. 43:343–368.
- Willis A., Bell R. 2018. Uncertainty in phylogenetic tree estimates. Journal of Computational
  and Graphical Statistics. 27:542–552.
- 858 Wright A.M., Lloyd G.T. 2020. Bayesian analyses in phylogenetic palaeontology: Interpreting
- the posterior sample. Palaeontology. 63:997–1006.

860	Xu X., Tan Q., Gao Y., Bao Z.,	Yin Z., Guo B.,	Wang J., Tan L., Zhang	Y., Xing H. 2018. A

861 large-sized basal ankylopollexian from East Asia, shedding light on early biogeographic

history of Iguanodontia. Science Bulletin. 63:556–563.

- 863 Yates A.M. 2003. The species taxonomy of the sauropodomorph dinosaurs from the Löwenstein
- 864 Formation (Norian, Late Triassic) of Germany. Palaeontology. 46:317–337.

866 FIGURE CAPTIONS

### 867 Figure 1: Different distances can impose tree spaces with different characteristics.

868 First two dimensions of PCoA mappings of tree spaces, with minimum spanning tree of 350

- 869 points (solid lines). Higher dimensions depicted in supplementary information (Smith 2021). (a-
- b), 2 500 Bayesian (dots) and 100 parsimony (rings) trees from analysis of Yates (2003),

871 coloured by stratigraphic congruence (MIG, millions of years); (a) RF tree space does not exhibit

872 clear structure; MST indicates that the two apparent clusters do not correspond to clusters in the

873 original tree space, and that the mapping is highly distorted (MST extension factor = 22.4); (b)

path distance tree space (MST extension factor = 12.1), showing stratigraphic structure and clear

875 separation of parsimony and Bayesian trees; (c), cumulative correlation of stratigraphic fit (MIG)

876 with first *n* tree space axes; (d–f), 2 500 Bayesian and 54 parsimony trees from analysis of

877 Carpenter (2001); points coloured by tree balance (total cophenetic index; dark = balanced): (d),

strong clustering in RF mapping (silhouette coefficient = 0.70) has no underlying basis

(silhouette coefficient =  $0.040 \ll 0.2$ ), as suggested by tortuous minimum spanning tree

880 (extension factor = 33.3); (e) vertical axis in KC mapping (MST extension factor = 9.61) shows

clear correspondence with tree balance; (f) quartet mapping (MST extension factor = 13.4)

faithfully represents the absence of clustering and tree balance correlation present in the original

space; (g-i) trees from analysis of Fischer et al. (2016), showing (lack of) correspondence

884 between original clusters (point colour, corresponding to Bayesian vs. parsimony trees) and

885 clusters identified from mappings (using hierarchical clustering; dashed lines = convex hulls); (j-

1) trees from analysis of Schoch and Milner (2008), coloured by stratigraphic fit; different

887 metrics result in spaces with different (non-clustering) structures, whose validity is supported by

888 inspection of MST and of higher dimensions. Abbreviations: CID, clustering information

distance tree space; KC, Kendall–Colijn tree space; MIG, Minimum implied gap; Q'tet, Quartet
tree space; RF, Robinson–Foulds tree space; SV, split size vector tree space; TCI, total
cophenetic index.

- 892
- 893 Figure 2: Methods for optimal clusterings.

894 (a-b), strength of clustering (silhouette coefficient) across all 128 tree sets under each tree 895 distance metric in (a) original tree space; (b) 2D PCoA mapping. Box plots denote median and 896 interquartile range; strong evidence that medians differ exists where notches do not overlap. (c-897 d), Number of clusters in optimal clustering under each tree distance method, calculated from (c) 898 original distances; (d) 2D PCoA mapping. Tree sets lacking 'reasonable' structure (i.e. silhouette 899 coefficients < 0.5) are taken to exhibit a single cluster. (e–f), Mean difference (variation of 900 information) between optimal clusterings obtained under (e), each tree distance metric, (f), each 901 clustering method, from datasets exhibiting at least 'reasonable' clustering structure. Brighter 902 colours represent greater differences. (g), Definition (silhouette coefficient) of optimal clustering obtained under each clustering method, summarized for all distances and tree sets. Bars denote 903 904 medians and interquartile ranges. (h), Method obtaining clustering with highest silhouette 905 coefficient, across all tree spaces with at least 'reasonable' clustering structure (silhouette 906 coefficient > 0.5).

907

### 908 Figure 3: Quality of tree space mappings based on underlying tree distance method.

909 (a–b), difference between original and mapped clusterings, in tree spaces that contain at least

910 'reasonable' clustering structure (silhouette coefficient > 0.5); (c), number of dimensions where

911 trustworthiness and continuity are each > 0.95; (d), length of minimum spanning tree relative to

912 shortest possible in 3D PCoA mappings; increasing values indicate more distorted mappings; (e-

913 f), cumulative correlation coefficient (r<sup>2</sup>) between Sammon mapping axes and (e), original tree

914 distance; (f), stratigraphic congruence (MIG); (g), correlation dimension of tree spaces.

915 Box and whisker plots depict medians and interquartile ranges; where notches do not 916 overlap, strong evidence exists that medians differ.

917

# 918 Figure 4: Effectiveness of mapping methods.

(a), trustworthiness × continuity; (b), minimum spanning tree extension factor; (c), correlation
with original distances (adjusted r<sup>2</sup>); (d), correlation with stratigraphic congruence (MIG,
adjusted r<sup>2</sup>). Lines depict median and interquartile range. Kruskal-1 mappings (omitted for
clarity) behave equivalently to PCoA. t-SNE mapping results only available for first three
dimensions.

924

#### 925 Figure 5: Structure 'hidden' in higher dimensions.

First six dimensions of phylogenetic information distance PCoA tree space, showing 2500
Bayesian trees (dots) and single most parsimonious tree (circles). Results from: bottom left,
Russell and Dong (1993); top right, Xu et al. (2018). Structural features evident in higher
dimensions are not apparent within the first two dimensions (top left corner), particularly with
regard to stratigraphic congruence, which is strongly correlated with higher dimensions of tree
space (bottom, right).

- 932 TABLE CAPTIONS
- 933 Table 1. Performance of selected tree distance metrics against tests of tree distance behaviour.
- 934 Parentheses denote range of possible scores for each measure (best to worst). Note that random
- 935 tree pairs obtain the maximum possible RF distance, resulting in a zero interquartile range (\*).
- 936 Full details and results in Smith (2020a) and Smith (2021).
- 937 *Acronyms:* CID: clustering information distance; KC: Kendall–Colijn distance; QD:
- 938 quartet distance; RF: Robinson–Foulds distance; SV: split size vector distance.