

Assessment of available anatomical characters for phylogenetic analysis among living mammals

THOMAS GUILLERME^{1,*} AND NATALIE COOPER^{1,2}

¹*School of Natural Sciences, Trinity College Dublin, Dublin 2, Ireland.*

²*Department of Life Sciences, Natural History Museum, Cromwell Road, London, SW7 5BD, UK.*

***Corresponding author.** *t.guillerm@imperial.ac.uk*

Abstract

Analyses of living and fossil taxa are crucial for understanding changes in biodiversity through time. The Total Evidence method allows living and fossil taxa to be combined in phylogenies, by using molecular data for living taxa and morphological data for both living and fossil taxa. With this method, substantial overlap of anatomical characters among living and fossil taxa is crucial for accurately inferring topology. However, although molecular data for living species is widely available, scientists using and generating morphological data mainly focus on fossils. Therefore, there is a gap in our knowledge of anatomical characters available for living taxa even in well-studied groups such as mammals.

We investigated the amount of available anatomical characters for living mammals and how this data was phylogenetically distributed across orders. 22 of 28 mammalian orders have <25% species with available anatomical characters; this has implications for the accurate placement of fossil taxa, although the issue is less pronounced at higher taxonomic levels. In most orders, species with available data are randomly distributed across the phylogeny, which may reduce the impact of the problem. We suggest that increased morphological data collection efforts for living taxa are needed to produce accurate Total Evidence phylogenies.

Key words: Total Evidence method, phylogenetic clustering, cladistic matrix, extinct, topology

INTRODUCTION

There is an increasing consensus among biologists that studying both living and fossil taxa is essential for fully understanding macroevolutionary patterns and processes [1, 2]. To perform such analyses it is necessary to combine living and fossil taxa in phylogenetic trees. One increasingly popular method, the Total Evidence method [3], combines molecular data from living taxa and morphological data from both living and fossil taxa in a supermatrix (e.g. [4, 3, 5, 1, 6]), producing a phylogeny with living and fossil taxa at the tips. A downside of this method is that it requires molecular data for living taxa and morphological data for both living and fossil taxa. Sections of this data can be difficult, or impossible, to collect for every taxon in the analysis. For example, fossils rarely have molecular data and incomplete fossil preservation may restrict the amount of anatomical characters available. Additionally, it has become less common to collect anatomical characters for living taxa when molecular data is available (e.g. in [7], only 13% of living taxa have coded anatomical characters). Unfortunately this missing data can lead to errors in phylogenetic inference. Simulations show that the ability of the Total Evidence method to recover the correct topology decreases when there is little overlap between anatomical characters in living and fossil taxa, and that the effect of missing data on topology is greatest when living taxa have few anatomical characters available [8]. This is because (1) fossils will not be placed accurately within the correct clade if it contains no anatomical characters for living taxa; and (2) fossils have a higher probability of branching within clades with

more anatomical characters available for living taxa, regardless of whether this is the correct clade [8].

The issues above highlight that it is crucial to have sufficient anatomical characters available for living taxa in a clade before using a Total Evidence approach. However, it is unclear how much anatomical characters are actually available for living taxa, i.e. already coded from museum specimens and deposited in phylogenetic matrices accessible online, and how this data is distributed across clades. Intuitively, most people assume this kind of data has already been collected, but empirical data suggest otherwise (e.g. in [3, 7, 6]). To investigate this further, we assess the amount of available anatomical characters for living mammals to determine whether sufficient data exists to build reliable Total Evidence phylogenies in this group. We also determine whether the available anatomical characters are phylogenetically overdispersed or clustered across mammalian orders.

MATERIALS AND METHODS

Data collection and standardisation

We downloaded all cladistic matrices containing any living and/or fossil mammal taxa from three major public databases: MorphoBank (morphobank.org [9]), Graeme Lloyd's website (graemetlloyd.com/matrmamm.html) and Ross Mounce's GitHub repository (github.com/rossmounce/cladistic-data). We also performed a systematic Google

Scholar search for matrices that were not uploaded to these databases (see Electronic Supplementary Material (ESM) for a detailed description of the search procedure). In total, we downloaded 286 matrices containing 5228 unique operational taxonomic units (OTUs). We used OTUs rather than species because entries in the matrices ranged from species to families, and standardised the taxonomy as described in the ESM. We designated as “living” all OTUs that were either present in the phylogeny of [10] or the taxonomy of [11].

Matrices with few characters are problematic when comparing available data among matrices because (1) they have less chance of having overlapping characters with other matrices [12] and (2) they are more likely to contain specific characters that are not applicable across large clades (e.g. “antler ramifications” is a character that is only applicable to Cervidae [13]). Therefore we selected only matrices containing >100 characters for each OTU. This threshold was chosen to correspond with the number of characters used in [8] and [14]. Results of analyses with no threshold are available in the ESM. After removing matrices with <100 characters, we retained 1074 unique living OTUs from 126 matrices.

Data availability and distribution

To assess the availability of anatomical characters for each mammalian order, we calculated the percentage of OTUs with cladistic data at three different taxonomic levels: family, genus and species. We consider orders with <25% of living taxa with

available anatomical characters as having low data coverage, and orders with >75% of living taxa with available anatomical characters as having high data coverage.

For each order, we investigated whether the available anatomical characters were (i) randomly distributed, (ii) overdispersed or (iii) clustered, with respect to phylogeny, using two metrics from community phylogenetics: the Nearest Taxon Index (NTI; [15]) and the Net Relatedness Index (NRI; [15]). NTI is most sensitive to clustering or overdispersion near the tips, whereas NRI is more sensitive to it across the whole phylogeny [16]. Both metrics were calculated using the *picante* package in R [17, 18].

NTI [15] is based on mean nearest neighbour distance (*MNND*) and is calculated as follows:

$$NTI = - \left(\frac{\overline{MNND}_{obs} - \overline{MNND}_n}{\sigma(MNND_n)} \right) \quad (1)$$

where \overline{MNND}_{obs} is the observed mean branch length between each of n taxa with available anatomical characters and its nearest neighbour with available anatomical characters in the phylogeny, \overline{MNND}_n is the mean of 1000 *MNND* between n randomly drawn taxa, and $\sigma(MNND_n)$ is the standard deviation of these 1000 random *MNND* values. NRI is calculated in the same way, but using the mean phylogenetic distance (*MPD*):

$$NRI = - \left(\frac{\overline{MPD}_{obs} - \overline{MPD}_n}{\sigma(MPD_n)} \right) \quad (2)$$

where \overline{MPD}_{obs} is the observed mean phylogenetic branch length of the tree containing only the n taxa with available anatomical characters. Negative NTI and NRI values show that the focal taxa are more overdispersed across the phylogeny than expected by








chance, and positive values reflect clustering.




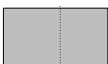
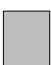




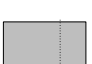






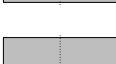
We calculated NTI and NRI values for each mammalian order separately, at each different taxonomic level. For each analysis our focal taxa were those with available anatomical characters at that taxonomic level and the phylogeny was that of the order pruned from [10].








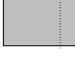


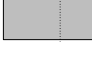






RESULTS

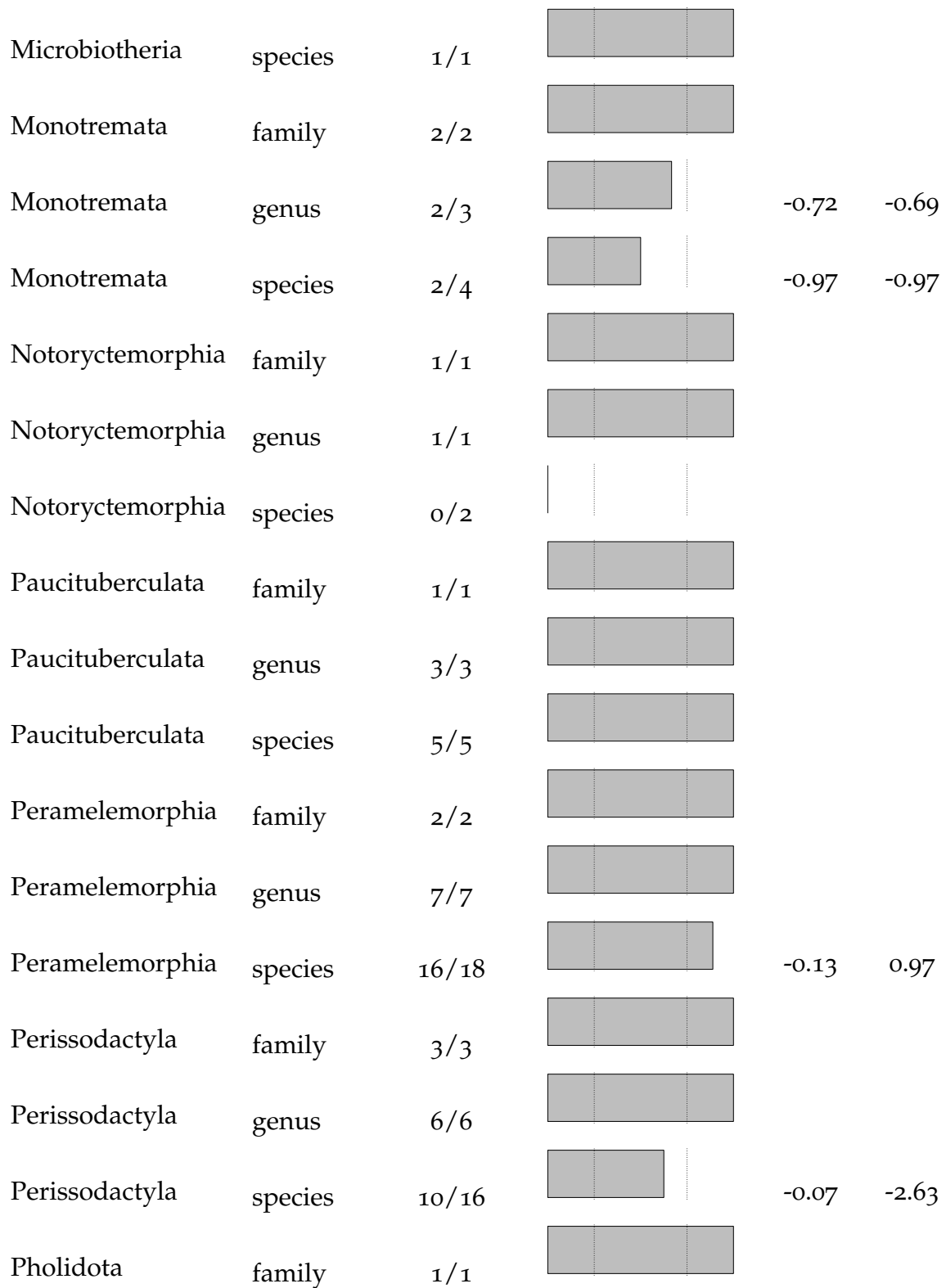
22 of 28 orders have low coverage ($<25\%$ species with available anatomical characters) and six have high coverage ($>75\%$ species with available anatomical characters) at the species-level. At the genus-level, three orders have low coverage and 12 have high coverage, and at the family-level, no orders have low coverage and 23 have high coverage (Table 1).











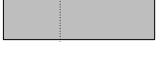
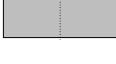



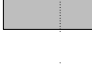

Table 1: Number of taxa with available cladistic data for mammalian orders at three taxonomic levels (without any character threshold; results from the 286 matrices). The coverage represents the proportion of taxa with available morphological data. The left vertical bar represents 25% of available data (“low” coverage if <25%); The right vertical bar represents 75% of available data (“high” coverage if >75%). When the Net Relatedness Index (NRI) and the Nearest Taxon Index (NTI) are negative, taxa are more phylogenetically dispersed than expected by chance; when NRI or NTI are positive, taxa are more phylogenetically clustered by expected by chance. Significant NRI or NTI are highlighted in bold. *p <0.05; **p <0.01; ***p <0.001.










Order	Taxo- nomic level	Propor- tion of taxa	Coverage	NRI	NTI
Afrosoricida	family	2/2			
Afrosoricida	genus	17/17			
Afrosoricida	species	23/42		1.75	1.08
Carnivora	family	14/15		0.63	0.6
Carnivora	genus	54/125		4.81**	1.78*
Carnivora	species	76/283		7.66**	0.85
Cetartiodactyla	family	21/21			

Cetartiodactyla	genus	100/128		0.85	0.94
Cetartiodactyla	species	171/310		1.92*	-0.46
Chiroptera	family	15/18		-0.28	0.56
Chiroptera	genus	93/202		13.47**	1.1
Chiroptera	species	215/1053		8.82**	1.22
Cingulata	family	1/1			
Cingulata	genus	8/9		1.51	-1.57
Cingulata	species	9/29		1.9*	0.11
Dasyuromorphia	family	2/2			
Dasyuromorphia	genus	8/22		-0.75	-1.07
Dasyuromorphia	species	9/64		-0.88	-0.34
Dermoptera	family	1/1			
Dermoptera	genus	1/2			
Dermoptera	species	1/2			
Didelphimorphia	family	1/1			
Didelphimorphia	genus	16/16			
Didelphimorphia	species	42/84		-1.65	0.2

Diprotodontia	family	11/11			
Diprotodontia	genus	25/38		-1.13	-1.31
Diprotodontia	species	31/126		0.48	-1.77
Erinaceomorpha	family	1/1			
Erinaceomorpha	genus	10/10			
Erinaceomorpha	species	21/22		-1.07	-0.2
Hyracoidea	family	1/1			
Hyracoidea	genus	1/3			
Hyracoidea	species	1/4			
Lagomorpha	family	2/2			
Lagomorpha	genus	5/12		-1.06	-0.95
Lagomorpha	species	12/86		-0.62	-1.88
Macroscelidea	family	1/1			
Macroscelidea	genus	4/4			
Macroscelidea	species	12/15		-1.3	-1.06
Microbiotheria	family	1/1			
Microbiotheria	genus	1/1			



Pholidota	genus	1/1			
Pholidota	species	4/8		1.18	0.94
Pilosa	family	4/5		1.87	2
Pilosa	genus	4/5		-0.96	0.36
Pilosa	species	5/29		1.28	2.38*
Primates	family	15/15			
Primates	genus	48/68		-0.35	-1.33
Primates	species	64/351		-0.67	-1.27
Proboscidea	family	1/1			
Proboscidea	genus	2/2			
Proboscidea	species	2/3		-0.69	-0.69
Rodentia	family	18/32		0.66	-0.98
Rodentia	genus	82/450		-1.66	1.55
Rodentia	species	90/2094		2.76*	2.34*
Scandentia	family	2/2			
Scandentia	genus	2/5		-0.74	-0.74
Scandentia	species	3/20		-1.88	-0.84

Sirenia	family	2/2			
Sirenia	genus	2/2			
Sirenia	species	4/4			
Soricomorpha	family	3/4		-0.98	-0.99
Soricomorpha	genus	19/43		7.11**	2.59**
Soricomorpha	species	21/392		10.65**	3.56**
Tubulidentata	family	1/1			
Tubulidentata	genus	1/1			
Tubulidentata	species	1/1			

Only six orders had significantly clustered data (Afrosoricida and Pholidota at the species-level, and Carnivora, Cetartiodactyla, Chiroptera and Soricomorpha at both species- and genus-level) and none had significantly overdispersed data (Table 1).

Figure 1 shows randomly distributed OTUs with cladistic data in Primates (Figure 1A) and phylogenetically clustered OTUs with cladistic data in Carnivora (mainly Canidae; Figure 1B).

DISCUSSION

Our results show that although phylogenetic relationships among living mammals are

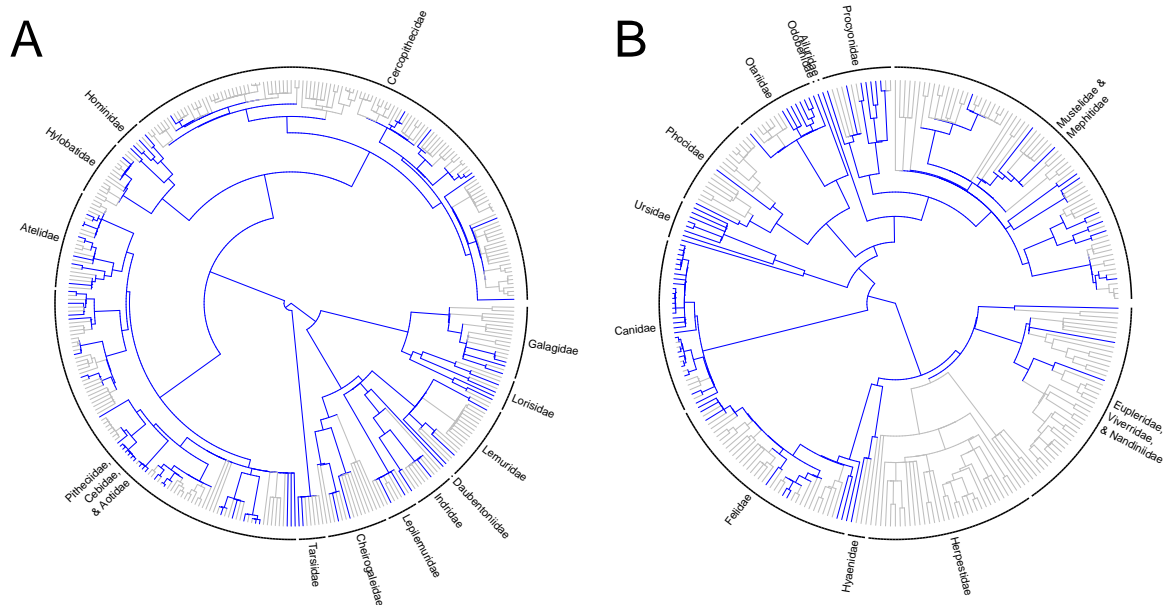


Figure 1: Phylogenetic distribution of species with available anatomical characters across two orders (A: Primates; B: Carnivora). Blue branches indicate available anatomical characters for the species.

well-resolved (e.g. [10, 19]), most of the data used to build these phylogenies is molecular, and few anatomical characters are available for living mammals compared to fossils (e.g. [20, 21]). This has implications for building Total Evidence phylogenies containing both living and fossil mammals, as without sufficient available anatomical characters for living species, fossil placements in these trees are very uncertain [8].

The number of living mammalian OTUs with no available anatomical characters was surprisingly high at the species-level: only six out of 28 orders have a high coverage of taxa with available anatomical characters. This high coverage threshold of 75% of taxa with available anatomical characters represents the minimum amount of data required before missing data has a significant effect on the topology of Total Evidence trees [8]. Beyond this threshold, there is considerable displacement of wildcard taxa and decreased clade conservation [8]. Therefore we expect difficulties in placement of fossils at the species-level in most mammalian orders, but fewer issues at higher taxonomic levels. This point is important from a practical point of view because of the slight discrepancy between neontological and palaeontological species concepts. While neontological species are described using morphology, genes, distribution etc.; palaeontological species can be based only on morphological, spatial and temporal data (e.g. [21]). Therefore, most palaeontological studies use the genus as their smallest OTU (e.g. [21, 20]), so data availability at the genus-level in living mammals should be our primary concern when building phylogenies of living and fossil taxa.

When few species have available anatomical characters, the ideal scenario is for

140 them to be evenly distributed (as measured by phylogenetic overdispersion) to
141 maximize the possibilities of a fossil branching in the right clade. The second best
142 scenario is that species with available anatomical characters are randomly distributed
143 across the phylogenies. Here we expect no special bias in the placement of fossils [8], it
144 is therefore encouraging that for most orders, species with available anatomical
145 characters were randomly distributed across the phylogeny. The worst case scenario for
146 fossil placement is that species with available anatomical characters are
147 phylogenetically clustered. Then we expect two major biases to occur: first, fossils will
148 not be able to branch within a clade containing no data, and second, fossils will have
149 higher probability of branching within the most sampled clade by chance. Our results
150 suggest that this may be problematic at the genus-level in Carnivora, Cetartiodactyla,
151 Chiroptera and Soricomorpha. For example, a Carnivora fossil will be unable to be
152 placed in the Herpestidae clade, and will have more chance to randomly branch within
153 Canidae (Figure 1B).

154 Despite the absence of good cladistic data coverage for living mammals, the
155 Total Evidence method still seems to be the most promising way of combining living
156 and fossil data for macroevolutionary analyses. Following the recommendations in [8],
157 we need to code anatomical characters for as many living species possible. Fortunately,
158 data for living mammals is usually readily available in natural history collections,
159 therefore, we propose that an increased effort be put into coding anatomical characters
160 from living species, possibly by engaging in collaborative data collection projects. Such

an effort would be valuable not only to phylogeneticists, but also to any researcher focusing understanding macroevolutionary patterns and processes.

ETHICS STATEMENT

N/A

DATA ACCESSIBILITY STATEMENT

All data and code are available on GitHub (https://github.com/TGuillerme/Missing_living_mammals).

AUTHORS' CONTRIBUTIONS

T.G. and N.C designed the study. T.G. analysed the data. T.G. and N.C. wrote the manuscript.

COMPETING INTERESTS

We have no competing interests.

ACKNOWLEDGMENTS

We thank David Bapst, Graeme Lloyd, Nick Matzke and April Wright. Thanks to one anonymous reviewer, Graham Slater, Peter Wagner and Graeme Lloyd for their useful comments.

FUNDING STATEMENT

This work was funded by a European Commission CORDIS Seventh Framework Programme (FP7) Marie Curie CIG grant (proposal number: 321696).

REFERENCES

*

References

- [1] Slater GJ, Harmon LJ. Unifying fossils and phylogenies for comparative analyses of diversification and trait evolution. *Methods Ecol Evol.* 2013;4(8):699–702.
- [2] Fritz SA, Schnitzler J, Eronen JT, Hof C, Böhning-Gaese K, Graham CH. Diversity in time and space: wanted dead and alive. *Trends Ecol Evol.* 2013;28(9):509 – 516.
- [3] Ronquist F, Klopstein S, Vilhelmsen L, Schulmeister S, Murray D, Rasnitsyn A. A total-evidence approach to dating with fossils, applied to the early radiation of the Hymenoptera. *Syst Biol.* 2012;61(6):973–999.
- [4] Pyron R. Divergence time estimation using fossils as terminal taxa and the origins of Lissamphibia. *Syst Biol.* 2011;60(4):466–481.
- [5] Schrago C, Mello B, Soares A. Combining fossil and molecular data to date the diversification of New World Primates. *J Evolution Biol.* 2013;26(11):2438–2446.

- [6] Beck RM, Lee MS. Ancient dates or accelerated rates? Morphological clocks and the antiquity of placental mammals. *P Roy Soc B-Biol Sci.* 2014;281(20141278):1–10.
- [7] Slater GJ. Phylogenetic evidence for a shift in the mode of mammalian body size evolution at the Cretaceous-Palaeogene boundary. *Methods Ecol Evol.* 2013;4(8):734–744.
- [8] Guillerme T, Cooper N. Effects of missing data on topological inference using a Total Evidence approach. *Mol Phylogenet Evol.* 2016;94, Part A:146 – 158.
Available from:
<http://www.sciencedirect.com/science/article/pii/S1055790315002547>.
- [9] O’Leary MA, Kaufman S. MorphoBank: phylophenomics in the cloud. *Cladistics.* 2011;27(5):529–537.
- [10] Bininda-Emonds ORP, Cardillo M, Jones KE, MacPhee RDE, Beck RMD, Grenyer R, et al. The delayed rise of present-day mammals. *Nature.* 2007 03;446(7135):507–512.
- [11] Wilson DE, Reeder DM. Mammal species of the world: a taxonomic and geographic reference. vol. 1. JHU Press; 2005.
- [12] Wagner PJ. Exhaustion of morphologic character states among fossil taxa. *Evolution.* 2000;54(2):365–386.
- [13] Brazeau MD. Problematic character coding methods in morphology and their effects. *Biol J Linn Soc.* 2011;104(3):489–498.

- 213 [14] Harrison LB, Larsson HCE. Among-Character Rate Variation Distributions in
214 Phylogenetic Analysis of Discrete Morphological Characters. *Syst Biol*.
215 2015;64(2):307–324.
- 216 [15] Webb CO, Ackerly DD, McPeck MA, Donoghue MJ. Phylogenies and community
217 ecology. *Ann Rev Ecol Syst*. 2002;p. 475–505.
- 218 [16] Cooper N, Rodríguez J, Purvis A. A common tendency for phylogenetic
219 overdispersion in mammalian assemblages. *P Roy Soc B-Biol Sci*.
220 2008;275(1646):2031–2037.
- 221 [17] Kembel SW, Cowan PD, Helmus MR, Cornwell WK, Morlon H, Ackerly DD, et al.
222 Picante: R tools for integrating phylogenies and ecology. *Bioinformatics*.
223 2010;26:1463–1464.
- 224 [18] R Core Team. R: a language and environment for statistical computing. Vienna,
225 Austria; 2015. Available from: <http://www.R-project.org>.
- 226 [19] Meredith R, Janečka J, Gatesy J, Ryder O, Fisher C, Teeling E, et al. Impacts of the
227 Cretaceous terrestrial revolution and KPg extinction on mammal diversification.
228 *Science*. 2011;334(6055):521–524.
- 229 [20] O’Leary MA, Bloch JL, Flynn JJ, Gaudin TJ, Giallombardo A, Giannini NP, et al.
230 The placental mammal ancestor and the postK-Pg radiation of placentals. *Science*.
231 2013;339(6120):662–667.

232 [21] Ni X, Gebo DL, Dagosto M, Meng J, Tafforeau P, Flynn JJ, et al. The oldest known
233 primate skeleton and early haplorhine evolution. *Nature*. 2013;498(7452):60–64.