

Version dated: July 23, 2015

# Assessment of cladistic data availability for living mammals

THOMAS GUILLERME<sup>1,\*</sup> AND NATALIE COOPER<sup>1,2</sup>

<sup>1</sup>*School of Natural Sciences, Trinity College Dublin, Dublin 2, Ireland.*

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

**\*Corresponding author.** [guillert@tcd.ie](mailto:guillert@tcd.ie)

## 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 morphological data 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 neontological morphological data even in well-studied groups such as mammals.

We investigated the amount of morphological (cladistic) data available for living mammals and how this data was phylogenetically distributed across orders. 22 of 28 mammalian orders have <25% species with available morphological data; 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, morphological matrix, extinct, topology

## INTRODUCTION

21

22       There is an increasing consensus among biologists that studying both living and  
23 fossil taxa is essential for fully understanding macroevolutionary patterns and  
24 processes [1, 2]. To perform such analyses it is necessary to combine living and fossil  
25 taxa in phylogenetic trees. One increasingly popular method, the Total Evidence  
26 method [3, 4], combines molecular data from living taxa and morphological data from  
27 both living and fossil taxa in a supermatrix (e.g. [5, 4, 6, 1, 7]), producing a phylogeny  
28 with living and fossil taxa at the tips. A downside of this method is that it requires  
29 molecular data for living taxa and morphological data for both living and fossil taxa.  
30 Chunks of this data can be difficult, or impossible, to collect for every taxon in the  
31 analysis. For example, fossils rarely have molecular data and incomplete fossil  
32 preservation may restrict the amount of morphological data available. Additionally, it  
33 has become less common to collect morphological characters for living taxa when  
34 molecular data is available (e.g. in [8], only 13% of living taxa have coded  
35 morphological data). Unfortunately this missing data can lead to errors in phylogenetic  
36 inference. Simulations show that the ability of the Total Evidence method to recover the  
37 correct topology decreases when there is little overlap between morphological data in  
38 living and fossil taxa, and that the effect of missing data on topology is greatest when  
39 living taxa have few morphological data [9]. This is because (1) fossils cannot branch in  
40 the correct clade if it contains no morphological data for living taxa; and (2) fossils have  
41 a higher probability of branching within clades with more morphological data for

42 living taxa, regardless of whether this is the correct clade [9].

43           The issues above highlight that it is crucial to have sufficient morphological data  
44 for living taxa in a clade before using a Total Evidence approach. However, it is unclear  
45 how much morphological data for living taxa is actually available, i.e. already coded  
46 from museum specimens and deposited in phylogenetic matrices accessible online, and  
47 how this data is distributed across clades. Intuitively, most people assume this kind of  
48 data has already been collected, but empirical data suggest otherwise (e.g. in [4, 8, 7]).  
49 To investigate this further, we assess the amount of available morphological data for  
50 living mammals to determine whether sufficient data exists to build reliable Total  
51 Evidence phylogenies in this group. We also determine whether the available cladistic  
52 data is phylogenetically overdispersed or clustered across mammalian orders.

## 53   MATERIALS AND METHODS

### 54   *Data collection and standardisation*

55 We downloaded all cladistic matrices containing any living and/or fossil mammal taxa  
56 from three major public databases: MorphoBank (<http://www.morphobank.org/> [10]),  
57 Graeme Lloyd’s website ([graemelloyd.com/matrmmamm.html](http://graemelloyd.com/matrmmamm.html)) and Ross Mounce’s  
58 GitHub repository (<https://github.com/rossmounce/cladistic-data>). We also  
59 performed a systematic Google Scholar search for matrices that were not uploaded to  
60 these databases (see Supplementary Materials Section 1 for a detailed description of the

61 search procedure). In total, we downloaded 286 matrices containing 5228 unique  
62 operational taxonomic units (OTUs). We used OTUs rather than species since entries in  
63 the matrices ranged from species to families, and standardised the taxonomy as  
64 described in Supplementary Materials (section 1). We designated as “living” all OTUs  
65 that were either present in the phylogeny of [11] or the taxonomy of [12].

66 Matrices with few characters are problematic when comparing available data  
67 among matrices because (1) they have less chance of having characters that overlap  
68 with those of other matrices [13] and (2) they are more likely to contain a higher  
69 proportion of specific characters that are not-applicable across large clades (e.g. “antler  
70 ramifications” is a character that is only applicable to Cervidae not all mammals [14]).  
71 Therefore we selected only matrices containing  $>100$  characters for each OTU. This  
72 threshold was chosen to correspond with the number of characters used in [9] and [15].  
73 Results of analyses with no threshold are available in Supplementary Material. After  
74 removing matrices with  $<100$  characters, we retained 1074 unique living mammal  
75 OTUs from 126 matrices.

### 76 *Data availability and distribution*

77 To assess the availability of cladistic data for each mammalian order, we calculated the  
78 percentage of OTUs with cladistic data at three different taxonomic levels: family,  
79 genus and species. We consider orders with  $<25\%$  of living taxa with cladistic data as  
80 having low data coverage, and orders with  $>75\%$  of living taxa with cladistic data as  
81 having high data coverage.

82 We investigated whether the available cladistic data for each order was (i)  
83 randomly distributed, (ii) overdispersed or (iii) clustered, with respect to phylogeny,  
84 using two metrics from community phylogenetics: the Nearest Taxon Index (NTI; [16])  
85 and the Net Relatedness Index (NRI; [16]). NTI is most sensitive to clustering or  
86 overdispersion near the tips, whereas NRI is more sensitive to clustering or  
87 overdispersion across the whole phylogeny [17]. Both metrics were calculated using the  
88 *picante* package in R [18, 19].

89 NTI [16] is based on mean nearest neighbour distance (MNND) and is calculated  
90 as follows:

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

91 where  $\overline{MNND}_{obs}$  is the observed mean distance between each of  $n$  taxa with cladistic  
92 data and its nearest neighbour with cladistic data in the phylogeny,  $\overline{MNND}_n$  is the  
93 mean of 1000 mean MNND between  $n$  randomly drawn taxa, and  $\sigma(MNND_n)$  is the  
94 standard deviation of these 1000 random MNND values. NRI is calculated in the same  
95 way, but MNND is replaced by mean phylogenetic distance (MPD) as follows:

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

96 where  $\overline{MPD}_{obs}$  is the observed mean phylogenetic distance of the tree containing only  
97 the  $n$  taxa with cladistic data. Negative NTI and NRI values show that the focal taxa are  
98 more overdispersed across the phylogeny than expected by chance, and positive values  
99 reflect clustering.



100 We calculated NTI and NRI values for each mammalian order separately, at each

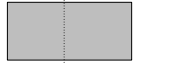





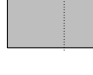
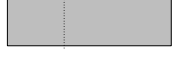
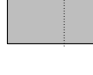



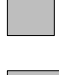

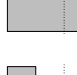


101 different taxonomic level. For each analysis our focal taxa were those with available  
102 cladistic data at that taxonomic level and the phylogeny was that of the order pruned  
103 from [11].

## 104 RESULTS

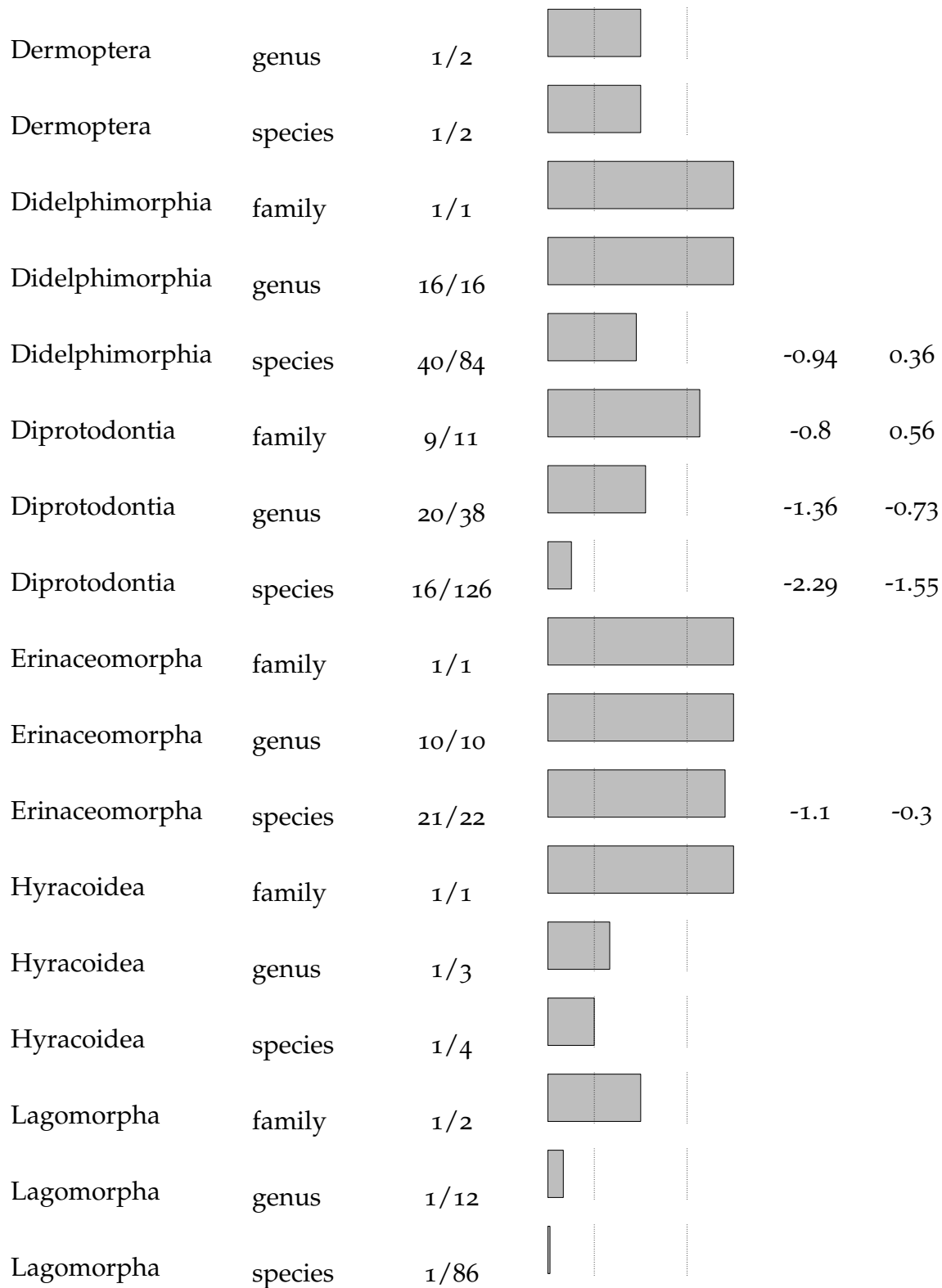
105 22 of 28 orders have low coverage (<25% species with cladistic data) and six have high  
106 coverage (>75% species with cladistic data) at the species-level. At the genus-level,  
107 three orders have low coverage and 12 have high coverage, and at the family-level, no  
108 orders have low coverage and 23 have high coverage (Table 1).

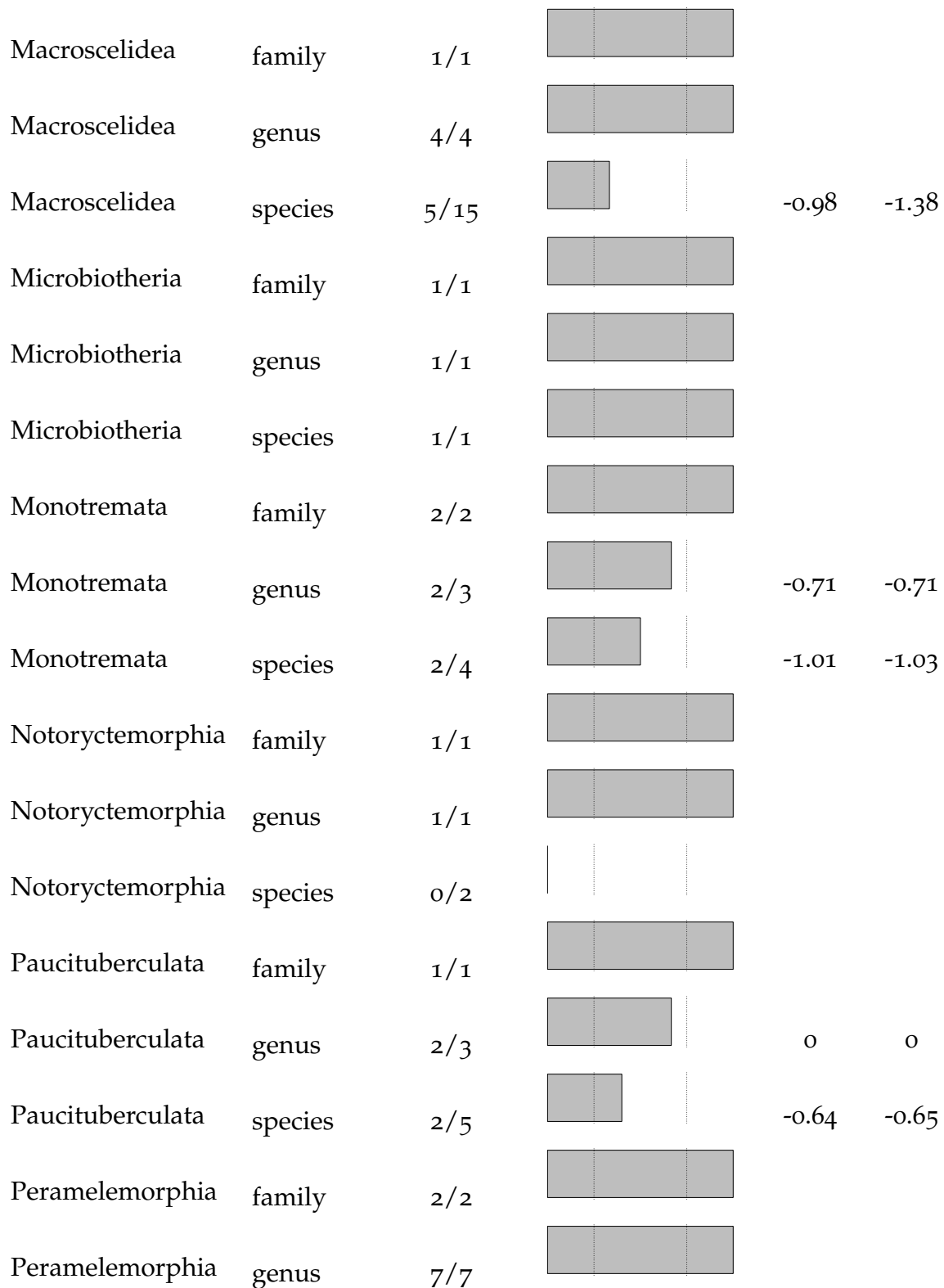
Table 1: Number of taxa with available cladistic data for mammalian orders at three taxonomic levels. The left vertical bar represents low coverage (<25%); the right vertical bar represents high coverage (>75%). Negative Net Relatedness Index (NRI) and Nearest Taxon Index (NTI) values indicate phylogenetic overdispersion; positive values indicate phylogenetic clustering. Significant NRI or NTI values are in bold. \*p < 0.05; \*\*p < 0.01; \*\*\*p < 0.001.

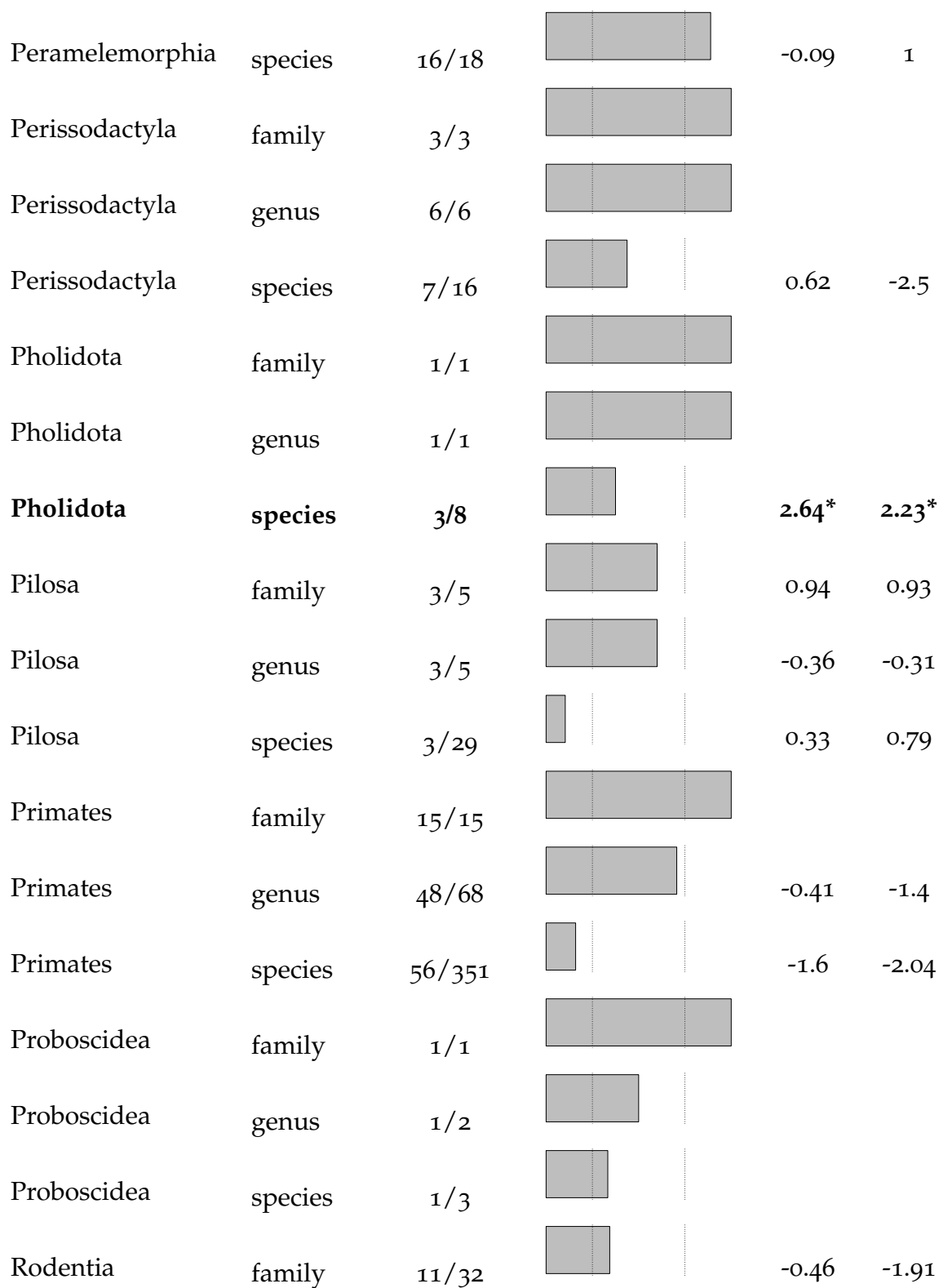
Order	Taxonomic level	Proportion of taxa	Coverage	NRI	NTI
Afrosoricida	family	2/2			
Afrosoricida	genus	17/17			









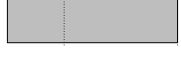
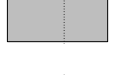




<b>Afrosoricida</b>	<b>species</b>	<b>23/42</b>		<b>1.89*</b>	1.19
Carnivora	family	11/15		0.43	1.68
<b>Carnivora</b>	<b>genus</b>	<b>30/125</b>		<b>4.14**</b>	<b>1.81*</b>
<b>Carnivora</b>	<b>species</b>	<b>42/283</b>		<b>18.64**</b>	<b>3.02**</b>
Cetartiodactyla	family	21/21			
<b>Cetartiodactyla</b>	<b>genus</b>	<b>77/128</b>		0.87	<b>1.77*</b>
<b>Cetartiodactyla</b>	<b>species</b>	<b>129/310</b>		<b>2.72*</b>	0.04
Chiroptera	family	13/18		0.55	0.63
<b>Chiroptera</b>	<b>genus</b>	<b>85/202</b>		<b>16.91**</b>	<b>2.85**</b>
<b>Chiroptera</b>	<b>species</b>	<b>165/1053</b>		<b>14.55**</b>	<b>3.44**</b>
Cingulata	family	1/1			
Cingulata	genus	8/9		1.49	-1.63
Cingulata	species	6/29		1.43	0.36
Dasyuromorphia	family	2/2			
Dasyuromorphia	genus	7/22		-1	-1.45
Dasyuromorphia	species	8/64		-1.15	-0.62
Dermoptera	family	1/1			









Rodentia	genus	21/450		-2.11	0.3
Rodentia	species	15/2094		-1.65	-2.55
Scandentia	family	2/2			
Scandentia	genus	2/5		-0.77	-0.76
Scandentia	species	2/20		-1.79	-1.99
Sirenia	family	2/2			
Sirenia	genus	2/2			
Sirenia	species	4/4			
Soricomorpha	family	3/4		-0.93	-0.92
<b>Soricomorpha</b>	<b>genus</b>	<b>19/43</b>		<b>6.98**</b>	<b>2.49*</b>
<b>Soricomorpha</b>	<b>species</b>	<b>19/392</b>		<b>13.19**</b>	<b>3.89**</b>
Tubulidentata	family	1/1			
Tubulidentata	genus	1/1			
Tubulidentata	species	1/1			

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

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

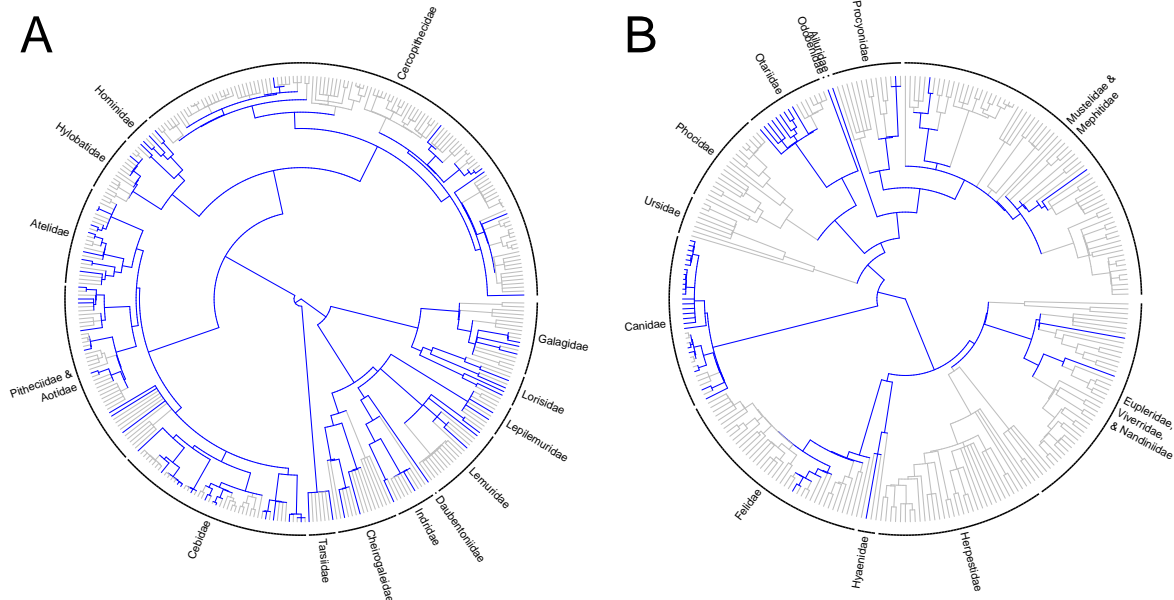


Figure 1: Phylogenetic distribution of species with available cladistic data across two orders (A: Primates; B: Carnivora). Blue branches indicate available cladistic data for the species.

115

## DISCUSSION

116 Our results show that although phylogenetic relationships among living mammals are  
117 well-resolved (e.g. [11, 20]) , most of the data used to build these phylogenies is  
118 molecular, and very little cladistic data is available for living mammals compared to  
119 fossil mammals (e.g. [21, 22]). This has implications for building Total Evidence

120 phylogenies containing both living and fossil mammals, as without sufficient cladistic  
121 data for living species, fossil placements in these trees are very uncertain [9].

122         The number of living mammalian taxa with no available cladistic data was  
123 surprisingly high at the species-level: only six out of 28 orders have a high coverage of  
124 taxa with available cladistic data. This high coverage threshold of 75% of taxa with  
125 available cladistic data represents the minimum amount of data required before  
126 missing data has a significant effect on the topology of Total Evidence trees [9]. Beyond  
127 this threshold, there is considerable displacement of wildcard taxa (*sensu* [23]) and  
128 decreased clade conservation [9]. Therefore we expect difficulties in placement of fossil  
129 taxa at the species-level in most mammalian orders, but fewer issues at higher  
130 taxonomic levels. This point is important from a practical point of view because of the  
131 slight discrepancy between neontological and palaeontological species concepts. While  
132 neontological species are described using morphology, genes, distribution etc.;  
133 palaeontological species can be based only on morphological, spatial and temporal data  
134 (e.g. [22]). Therefore, most palaeontological studies use genus as their smallest OTU  
135 (e.g. [22, 21]), so data availability at the genus-level in living mammals should be our  
136 primary concern when building phylogenies of living and fossil taxa.

137         When few species have available cladistic data, the ideal scenario is for them to  
138 be phylogenetically overdispersed to maximize the possibilities of a fossil branching  
139 from the right clade. The second best scenario is that species with cladistic data are  
140 randomly distributed across the phylogeny. Here we expect no special bias in the

141 placement of fossils [9], it is therefore encouraging that for most orders, species with  
142 cladistic data were randomly distributed across the phylogeny. The worst case scenario  
143 for fossil placement is that species with cladistic data are phylogenetically clustered.  
144 Then we expect two major biases to occur: first, fossils will not be able to branch within  
145 a clade containing no data, and second, fossils will have higher probability of branching  
146 within the most sampled clade by chance. Our results suggest that this may be  
147 problematic at the genus-level in Carnivora, Cetartiodactyla, Chiroptera and  
148 Soricomorpha. For example, a Carnivora fossil will be unable to branch in the  
149 Herpestidae, and will have more chance to randomly branch within Canidae (Figure  
150 1B).

151 Despite the absence of good cladistic data coverage for living mammals, the  
152 Total Evidence method still seems to be the most promising way of combining living  
153 and fossil data for macroevolutionary analyses. Following the recommendations in [9],  
154 we need to code cladistic characters for as many living species possible. Fortunately,  
155 data for living mammals is usually readily available in natural history collections,  
156 therefore, we propose that an increased effort be put into coding morphological  
157 characters from living species, possibly by engaging in collaborative data collection  
158 projects. Such an effort would be valuable not only to phylogeneticists, but also to any  
159 researcher focusing understanding macroevolutionary patterns and processes.

## 160 ETHICS STATEMENT

161 N/A

162

## DATA ACCESSIBILITY STATEMENT

163 All data and analysis code is available on GitHub

164 ([https://github.com/TGuillerme/Missing\\_living\\_mammals](https://github.com/TGuillerme/Missing_living_mammals)).

165

## AUTHORS' CONTRIBUTIONS

166 T.G. and N.C conceived and designed the experiments. T.G. performed the experiments

167 and analysed the data. T.G. and N.C. contributed to the writing of the manuscript. All

168 authors approved the final version of the manuscript.

169

## COMPETING INTERESTS

170 We have no competing interests.

171

## ACKNOWLEDGMENTS

172 We thank David Bapst, Graeme Lloyd, Nick Matzke and April Wright.

173

## FUNDING STATEMENT

174 This work was funded by a European Commission CORDIS Seventh Framework

175 Programme (FP7) Marie Curie CIG grant (proposal number: 321696).



\*

176

177 References

- 178 [1] Slater GJ, Harmon LJ. Unifying fossils and phylogenies for comparative analyses  
179 of diversification and trait evolution. *Methods Ecol Evol.* 2013;4(8):699–702.
- 180 [2] Fritz SA, Schnitzler J, Eronen JT, Hof C, Böhning-Gaese K, Graham CH. Diversity  
181 in time and space: wanted dead and alive. *Trends Ecol Evol.* 2013;28(9):509 – 516.
- 182 [3] Eernisse D, Kluge A. Taxonomic congruence versus total evidence, and amniote  
183 phylogeny inferred from fossils, molecules, and morphology. *Mol Biol Evol.*  
184 1993;10(6):1170–1195.
- 185 [4] Ronquist F, Klopfstein S, Vilhelmsen L, Schulmeister S, Murray D, Rasnitsyn A. A  
186 total-evidence approach to dating with fossils, applied to the early radiation of the  
187 Hymenoptera. *Syst Biol.* 2012;61(6):973–999.
- 188 [5] Pyron R. Divergence time estimation using fossils as terminal taxa and the origins  
189 of Lissamphibia. *Syst Biol.* 2011;60(4):466–481.
- 190 [6] Schrago C, Mello B, Soares A. Combining fossil and molecular data to date the  
191 diversification of New World Primates. *J Evolution Biol.* 2013;26(11):2438–2446.
- 192 [7] Beck RM, Lee MS. Ancient dates or accelerated rates? Morphological clocks and  
193 the antiquity of placental mammals. *P Roy Soc B-Biol Sci.* 2014;281(20141278):1–10.

- 194 [8] Slater GJ. Phylogenetic evidence for a shift in the mode of mammalian body size  
195 evolution at the Cretaceous-Palaeogene boundary. *Methods Ecol Evol.*  
196 2013;4(8):734–744.
- 197 [9] Guillerme T, Cooper N. Effects of missing data on topological inference using a  
198 Total Evidence approach. *Mol Phylogenet Evol.* In review;X(X):X.
- 199 [10] O’Leary MA, Kaufman S. MorphoBank: phylophenomics in the cloud. *Cladistics.*  
200 2011;27(5):529–537.
- 201 [11] Bininda-Emonds ORP, Cardillo M, Jones KE, MacPhee RDE, Beck RMD, Grenyer R,  
202 et al. The delayed rise of present-day mammals. *Nature.* 2007 03;446(7135):507–512.
- 203 [12] Wilson DE, Reeder DM. *Mammal species of the world: a taxonomic and*  
204 *geographic reference.* vol. 1. JHU Press; 2005.
- 205 [13] Wagner PJ. Exhaustion of morphologic character states among fossil taxa.  
206 *Evolution.* 2000;54(2):365–386.
- 207 [14] Brazeau MD. Problematic character coding methods in morphology and their  
208 effects. *Biol J Linn Soc.* 2011;104(3):489–498.
- 209 [15] Harrison LB, Larsson HCE. Among-Character Rate Variation Distributions in  
210 *Phylogenetic Analysis of Discrete Morphological Characters.* *Syst Biol.*  
211 2015;64(2):307–324.

- 212 [16] Webb CO, Ackerly DD, McPeck MA, Donoghue MJ. Phylogenies and community  
213 ecology. *Ann Rev Ecol Syst.* 2002;p. 475–505.
- 214 [17] Cooper N, Rodríguez J, Purvis A. A common tendency for phylogenetic  
215 overdispersion in mammalian assemblages. *P Roy Soc B-Biol Sci.*  
216 2008;275(1646):2031–2037.
- 217 [18] Kembel SW, Cowan PD, Helmus MR, Cornwell WK, Morlon H, Ackerly DD, et al.  
218 *Picante: R tools for integrating phylogenies and ecology.* *Bioinformatics.*  
219 2010;26:1463–1464.
- 220 [19] R Core Team. *R: a language and environment for statistical computing.* Vienna,  
221 Austria; 2015. Available from: <http://www.R-project.org>.
- 222 [20] Meredith R, Janečka J, Gatesy J, Ryder O, Fisher C, Teeling E, et al. Impacts of the  
223 Cretaceous terrestrial revolution and KPg extinction on mammal diversification.  
224 *Science.* 2011;334(6055):521–524.
- 225 [21] O’Leary MA, Bloch JL, Flynn JJ, Gaudin TJ, Giallombardo A, Giannini NP, et al.  
226 The placental mammal ancestor and the postK-Pg radiation of placentals. *Science.*  
227 2013;339(6120):662–667.
- 228 [22] Ni X, Gebo DL, Dagosto M, Meng J, Tafforeau P, Flynn JJ, et al. The oldest known  
229 primate skeleton and early haplorhine evolution. *Nature.* 2013;498(7452):60–64.
- 230 [23] Kearney M. Fragmentary taxa, missing data, and ambiguity: mistaken  
231 assumptions and conclusions. *Syst Biol.* 2002;51(2):369–381.