

MACROEVOLUTION WITH FOSSIL AND LIVING TAXA

by

THOMAS GUILLERME

B.Sc., Université Montpellier 2, 2010

M.Sc., Université Montpellier 2, 2012

A thesis submitted in partial fulfillment of
the requirements for the degree of

DOCTOR OF PHILOSOPHY

School of Natural Sciences
(Zoology)

Trinity College Dublin

OCTOBER 2015

© Thomas Guillaume, 2015

DECLARATION

I declare that this thesis has not been submitted as an exercise for a degree at this or any other University and it is, unless otherwise referenced, entirely my own work. I agree to deposit this thesis in the University's open access institutional repository or allow the library to do so on my behalf, subject to Irish Copyright Legislation and Trinity College Library conditions of use and acknowledgement.

Thomas Guillerme

SUMMARY

Hablbalbalbalba. Make it short!

ACKNOWLEDGEMENTS

Thanks folks!

TABLE OF CONTENTS

DECLARATION.	2
SUMMARY.	3
ACKNOWLEDGEMENTS	4
TABLE OF CONTENTS	5
LIST OF TABLES	6
LIST OF FIGURES	7
1 INTRODUCTION	1
2 MISSING DATA IN LIVING MAMMALS	2
2.1 Summary	2
2.2 Introduction	3
2.3 Materials and Methods	4
2.3.1 Data collection and standardisation	4
2.3.2 Data availability and distribution	6
2.4 Results	7
2.5 Discussion	11
3 DISCUSSION	15
BIBLIOGRAPHY	16

LIST OF TABLES

TABLE 2.1	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%). A negative Net Relatedness Index (NRI) and Nearest Taxon Index (NTI) shows more phylogenetically dispersed taxa than expected by chance; a positive value shows more phylogenetically clustered taxa than expected by chance. Significant NRI or NTI values are highlighted in bold. One star (*) signifies a p-value between 0.05 and 0.005; two stars between 0.005 and 0.0005 and three stars <0.0005.	7
-----------	--	---

LIST OF FIGURES

FIGURE 2.1 Phylogenetic distribution of species with available cladistic data across two mammalian orders (A: Primates; B: Carnivora). Edges are colored in grey when no cladistic data is available for a species and in blue when data is available.	11
--	----

CHAPTER 1

INTRODUCTION

Let's start.

Assessment of cladistic data availability for living mammals

Key words: Total Evidence method, data structure, phylogenetic, fossil, topology

A shorter version (2500 words) as been submitted to Biology Letters as an invited submission for a special issue blablabla

2.1 SUMMARY

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.

There is an increasing consensus among evolutionary biologists that studying both living and fossil taxa is essential for fully understanding macroevolutionary patterns and processes (Slater and Harmon, 2013; Fritz et al., 2013; Wood et al., 2013). For example, including both living and fossil taxa in evolutionary studies can improve the accuracy of timing diversification events (e.g. Ronquist et al., 2012), our understanding of relationships among lineages (e.g. Beck and Lee, 2014), and our ability to infer biogeographical patterns through time (e.g. Meseguer et al., 2015). To perform such analyses it is necessary to combine living and fossil taxa in phylogenetic trees. One increasingly popular method, the Total Evidence method (Eernisse and Kluge, 1993; Ronquist et al., 2012), combines molecular data from living taxa and morphological data from both living and fossil taxa in a supermatrix (e.g. Pyron, 2011; Ronquist et al., 2012; Schrago et al., 2013; Slater and Harmon, 2013; Beck and Lee, 2014; Meseguer et al., 2015), producing a phylogeny with living and fossil taxa at the tips. These phylogenies can be dated using methods such as tip-dating (Ronquist et al., 2012; Wood et al., 2013) and incorporated into macroevolutionary studies (e.g. Ronquist et al., 2012; Wood et al., 2013; Slater, 2013).

A downside of the Total Evidence method is that it requires a lot of data. One must collect molecular data for living taxa and morphological data for both living and fossil taxa; two types of data that require fairly different technical skills (e.g. molecular sequencing *vs.* anatomical description). Additionally, large chunks of this data can be difficult, or even impossible, to collect for every taxon present in the analysis. For example, fossils very rarely have molecular data and incomplete fossil preservation (e.g. soft *vs.* hard tissues) may restrict the amount of morphological data available (Sansom and Wills, 2013). Additionally, since the molecular phylogenetics revolution, it has become less common to collect morphological characters for living taxa when molecular data is available (e.g. in (Slater, 2013), only 13% of the 169 living taxa have coded morphological data). Unfortunately this missing data can lead to errors in phylogenetic inference; in fact, simulations show that the ability of the Total Evidence method to recover the correct phylogenetic topology decreases when there is a low overlap between morphological data in the living and fossil taxa (Guillerme and Cooper, 2015), regardless the overall amount of morphological data available for the fossils (or the amount of molecular data available for the living species). The effect of missing data on topology is greatest when living taxa have few

morphological data. This is because (1) a fossil cannot branch in the correct clade if there is no overlapping morphological data in the clade; and (2) a fossil has a higher probability of branching within a clade with more morphological data available for living taxa, regardless of whether this is the correct clade or not (Guillerme and Cooper, 2015).

The issues above highlight that it is crucial to have sufficient morphological data for living taxa in a clade before using a Total Evidence approach. However, it is unclear how much morphological data for living taxa is actually available (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 (Ronquist et al., 2012; Slater, 2013; Beck and Lee, 2014). To investigate this further, we assess the amount of available morphological data for living mammals to determine whether sufficient data exists to build reliable Total Evidence phylogenies in this group. We collected cladistic data (i.e. discrete morphological characters used in phylogenetics) from 286 phylogenetic matrices available online and measured the proportion of cladistic data available for each mammalian order. Additionally, because missing morphological data in living species can influence tree topology as described above, we determined whether the available cladistic data was phylogenetically overdispersed or clustered in the mammalian orders where data was missing. We find that available morphological data for living mammals is scarce but generally randomly distributed across phylogenies. We recommend that efforts be made to collect and share more cladistic data for living species to improve the accuracy of Total Evidence phylogenies.

2.3 MATERIALS AND METHODS

2.3.1 *Data collection and standardisation*

We downloaded all cladistic matrices containing any living and/or fossil mammal taxa from three major public databases (accessed 10th of June 2015): Morphobank (<http://www.morphobank.org/>) (O’Leary and Kaufman, 2011), Graeme Lloyd’s website (graemetlloyd.com/matrmamm.html) and Ross Mounce’s GitHub repository (<https://github.com/rossmounce/cladistic-data>). We also performed a systematic Google Scholar search (accessed 11th of June 2015) for matrices that were not uploaded to these databases (see Supplementary Materials for a detailed description of the

CHAPTER 4
search procedure). In total, we downloaded 286 matrices containing a total of 11010 operational taxonomic units (OTUs) of which 5228 were unique. In this study, we refer to OTUs rather than species since the entries in the downloaded matrices were not standardised and ranged from specific individual specimen names (i.e. the name of a collection item) to the family-level. Where possible, we considered OTUs at their lowest valid taxonomic level (i.e. species) but some OTUs were only valid at a higher taxonomic level (e.g. genus or family). Therefore for some orders, we sampled more genera than species (Table ??).

To select the lowest valid taxonomic level for each OTU, we standardised their taxonomy by correcting species names so they matched standard taxonomic nomenclature (e.g., *H. sapiens* was transformed to *Homo sapiens*). We designated as “living” all OTUs that were either present in the phylogeny of (Bininda-Emonds et al., 2007) or the taxonomy of (Wilson and Reeder, 2005), and designated as “fossil” all OTUs that were present in the Paleobiology database (<https://paleobiodb.org/>). For OTUs that did not appear in these three sources, we first decomposed the name (i.e. *Homo sapiens* became *Homo* and *sapiens* and tried to match the first element with a higher taxonomic level (family, genus etc.). Any OTUs that still had no matches in the sources above were designated as non-applicable (NA; see Supplementary Material for more details).

The number of characters in each matrix ranged from 6 to 4541. Matrices with few characters are problematic when comparing available data among matrices because (1) they have less chance of having characters that overlap with those of other matrices (Wagner, 2000) and (2) they are more likely to contain a higher proportion of specific characters that are not-applicable across large clades (e.g. “antler ramifications” is a character that is only applicable to Cervidae not all mammals (Brazeau, 2011)). Therefore we selected only matrices containing >100 characters for each OTU. This threshold was chosen to correspond with the number of characters used in (Guillerme and Cooper, 2015) and (Harrison and Larsson, 2015). Note that results of analyses with no character threshold are available in Supplementary Material. After removing all matrices with <100 characters, we retained 1074 unique living mammal OTUs from 126 matrices for our analyses.

2.3.2^{CHAPTER 2} Data availability and distribution

To assess the availability of cladistic data 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 cladistic data as having poor data coverage (“low” coverage), and orders with >75% of living taxa with cladistic data as having good data coverage (hereafter “high” coverage).

For orders with <100% cladistic data coverage at any taxonomic level, we investigated whether the available cladistic data was (i) randomly distributed, (ii) overdispersed or (iii) clustered, with respect to phylogeny, using two metrics from community phylogenetics: the Nearest Taxon Index (NTI; (Webb et al., 2002) and the Net Relatedness Index (NRI; (Webb et al., 2002). NTI is most sensitive to clustering or overdispersion near the tips, whereas NRI is more sensitive to clustering or overdispersion across the whole phylogeny (Cooper et al., 2008). Both metrics were calculated using the *picante* package in R (Kembel et al., 2010; R Core Team, 2015).

NTI (Webb et al., 2002) is based on mean nearest neighbour distance (MNND) and is calculated as follows:

$$NTI = - \left(\frac{\overline{MNND}_{obs} - \overline{MNND}_n}{\sigma(MNND_n)} \right) \quad (2.1)$$

where \overline{MNND}_{obs} is the observed mean distance between each of n taxa with cladistic data and its nearest neighbour with cladistic data in the phylogeny, \overline{MNND}_n is the mean of 1000 mean MNND between n randomly drawn taxa, and $\sigma(MNND_n)$ is the standard deviation of these 1000 random MNND values. NRI is similar but is based on mean phylogenetic distance (MPD) as follows:

$$NRI = - \left(\frac{\overline{MPD}_{obs} - \overline{MPD}_n}{\sigma(MPD_n)} \right) \quad (2.2)$$








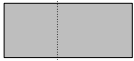

where \overline{MPD}_{obs} is the observed mean phylogenetic distance of the tree containing only the n taxa with cladistic data, \overline{MPD}_n is the expected random MPD for n taxa estimated by calculating the MPD from n taxa randomly drawn from the phylogeny and repeated 1000 times, and $\sigma(MPD_n)$ is the standard deviation of the 1000 random MPD values. Negative NTI and NRI values show that the focal taxa are more overdispersed across the phylogeny than expected by chance, and positive values reflect significant clustering.














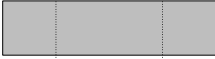




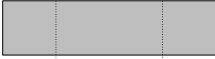





We calculated NTI and NRI values for each mammalian order separately, at each different taxonomic level. For each analysis our focal taxa were those with available cladistic data at that taxonomic level and the phylogeny was the phylogeny of the order pruned from (Bininda-Emonds et al., 2007).












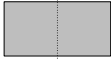

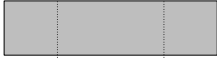


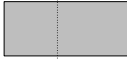

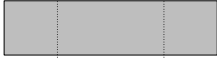





2.4 RESULTS

















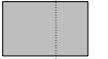

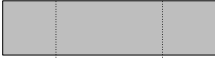





Across the 126 cladistic matrices we extracted, 22 out of 28 mammalian orders have low coverage (<25% of species with cladistic data) and six have high coverage (>75% of species with cladistic data) at the species-level. At the genus-level, three orders have low coverage and 12 have high coverage, and at the family-level, no orders have low coverage and 23 have high coverage (Table ??).

TABLE 2.1: 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%). A negative Net Relatedness Index (NRI) and Nearest Taxon Index (NTI) shows more phylogenetically dispersed taxa than expected by chance; a positive value shows more phylogenetically clustered taxa than expected by chance. Significant NRI or NTI values are highlighted in bold. One star (*) signifies a p-value between 0.05 and 0.005; two stars between 0.005 and 0.0005 and three stars <0.0005.

Order	Taxo- nomic level	Propor- tion of taxa	Coverage	NRI	NTI
Afrosoricida	family	2/2			
Afrosoricida	genus	17/17			
Afrosoricida	species	23/42		1.89*	1.19
Carnivora	family	11/15		0.43	1.68
Carnivora	genus	30/125		4.14**	1.81*
Carnivora	species	42/283		18.64**	3.02**
Cetartiodactyla	family	21/21			
Cetartiodactyla	genus	77/128		0.87	1.77*
Cetartiodactyla	species	129/310		2.72*	0.04

CHAPTER 2					
Chiroptera	family	13/18		0.55	0.63
Chiroptera	genus	85/202		16.91**	2.85**
Chiroptera	species	165/1053		14.55**	3.44**
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			
Dermoptera	genus	1/2			
Dermoptera	species	1/2			
Didelphimorphia	family	1/1			
Didelphimorphia	genus	16/16			
Didelphimorphia	species	40/84		-0.94	0.36
Diprotodontia	family	9/11		-0.8	0.56
Diprotodontia	genus	20/38		-1.36	-0.73
Diprotodontia	species	16/126		-2.29	-1.55
Erinaceomorpha	family	1/1			
Erinaceomorpha	genus	10/10			
Erinaceomorpha	species	21/22		-1.1	-0.3
Hyracoidea	family	1/1			
Hyracoidea	genus	1/3			
Hyracoidea	species	1/4			

			CHAPTER 2		
Lagomorpha	family	1/2			
Lagomorpha	genus	1/12			
Lagomorpha	species	1/86			
Macroscelidea	family	1/1			
Macroscelidea	genus	4/4			
Macroscelidea	species	5/15		-0.98	-1.38
Microbiotheria	family	1/1			
Microbiotheria	genus	1/1			
Microbiotheria	species	1/1			
Monotremata	family	2/2			
Monotremata	genus	2/3		-0.71	-0.71
Monotremata	species	2/4		-1.01	-1.03
Notoryctemorphia	family	1/1			
Notoryctemorphia	genus	1/1			
Notoryctemorphia	species	0/2			
Paucituberculata	family	1/1			
Paucituberculata	genus	2/3		0	0
Paucituberculata	species	2/5		-0.64	-0.65
Peramelemorphia	family	2/2			
Peramelemorphia	genus	7/7			
Peramelemorphia	species	16/18		-0.09	1
Perissodactyla	family	3/3			
Perissodactyla	genus	6/6			
Perissodactyla	species	7/16		0.62	-2.5

Pholidota	family	1/1			
Pholidota	genus	1/1			
Pholidota	species	3/8		2.64*	2.23*
Pilosa	family	3/5		0.94	0.93
Pilosa	genus	3/5		-0.36	-0.31
Pilosa	species	3/29		0.33	0.79
Primates	family	15/15			
Primates	genus	48/68		-0.41	-1.4
Primates	species	56/351		-1.6	-2.04
Proboscidea	family	1/1			
Proboscidea	genus	1/2			
Proboscidea	species	1/3			
Rodentia	family	11/32		-0.46	-1.91
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
Soricomorpha	genus	19/43		6.98**	2.49*
Soricomorpha	species	19/392		13.19**	3.89**

Tubulidentata	family	1/1	<div><div></div><div></div><div></div></div>
Tubulidentata	genus	1/1	<div><div></div><div></div><div></div></div>
Tubulidentata	species	1/1	<div><div></div><div></div><div></div></div>

Among the mammalian orders containing OTUs with no available cladistic data, only six orders had significantly clustered data (Carnivora, Cetartiodactyla, Chiroptera and Soricomorpha at both species- and genus-level and Afrosoricida and Pholidota at the species-level only) and no order had significantly overdispersed data at any taxonomic level (Table ??).

Two contrasting results are shown in Figure 2.1 with randomly distributed OTUs with cladistic data in Primates (Figure 2.1A) and phylogenetically clustered OTUs with cladistic data in Carnivora (mainly Canidae; Figure 2.1B).

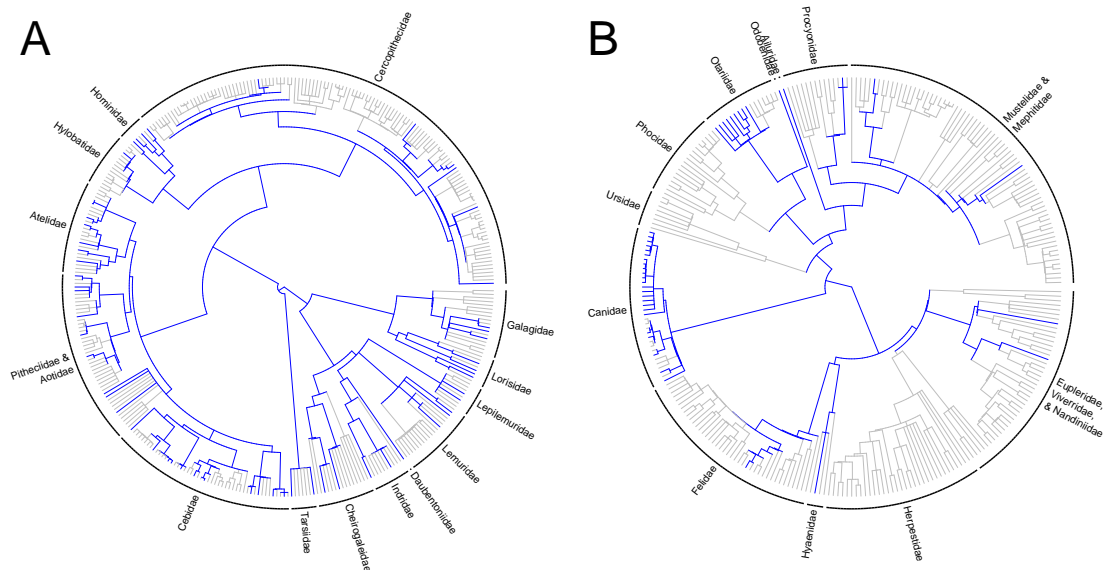


FIGURE 2.1: Phylogenetic distribution of species with available cladistic data across two mammalian orders (A: Primates; B: Carnivora). Edges are colored in grey when no cladistic data is available for a species and in blue when data is available.

2.5 DISCUSSION

Our results show that although phylogenetic relationships among living mammals are well-resolved (e.g. Bininda-Emonds et al., 2007; Meredith et al., 2011) , most of the data used to build these phylogenies is molecular, and very little cladistic data is available for living mammals compared to fossil mammals (e.g. O’Leary et al.,

2013; Ni et al., 2013). This has implications for building Total Evidence phylogenies containing both living and fossil mammals, as without sufficient cladistic data for living species, fossil placements in these trees are very uncertain (Guillerme and Cooper, 2015). Cladistic data coverage in living mammals varies across taxonomic levels and in its phylogenetic distribution. Higher taxonomic levels are always better sampled than lower ones and within these taxonomic levels, the available data is mostly randomly distributed across the phylogeny, apart from in six orders).

The number of living mammalian taxa with no available cladistic data was surprisingly high at the species-level: only six out of 28 orders have a high coverage of taxa with available cladistic data (and two of the 28 orders are monospecific!). This high coverage threshold of 75% of taxa with available cladistic data represents the minimum amount of data required before missing data has a significant effect on the topology of Total Evidence trees (Guillerme and Cooper, 2015). Beyond this threshold, there is considerable displacement of wildcard taxa (*sensu* (Kearney, 2002) and decreases in clade conservation (Guillerme and Cooper, 2015). Therefore we expect a high probability of topological artefacts for the placement of fossil taxa at the species-level in most mammalian orders. However, data coverage seems to be less of an issue at higher taxonomic levels (i.e. genus- and family-level). This point is important from a practical point of view because of the slight discrepancy between the neontological and palaeontological concept of species. While neontological species are described using morphology, genetic distance, spatial distribution and even behaviour, palaeontological species can be based only on morphological, spatial and temporal data (e.g. Ni et al., 2013). Because of this, most palaeontological studies are using the genus as their smallest OTU (e.g. Ni et al., 2013; O'Leary et al., 2013). Thus data availability at the genus-level in living mammals should be our primary concern when aiming to build phylogenies of living and fossil taxa.

When only a few species with cladistic data are available, the ideal scenario is for them to be phylogenetically overdispersed (i.e. that there is data for at least every sub-clade) to maximize the possibilities of a fossil branching from the right clade. The second best scenario is that species with cladistic data are randomly distributed across the phylogeny. In this scenario we expect no special bias in the placement of the fossil (Guillerme and Cooper, 2015), it is therefore encouraging that for most orders, species with cladistic data were randomly distributed across the phylogeny of each order. The worst case scenario for fossil placement is that species with cladistic data are phylogenetically clustered. In this situation we expect two major

biases to occur: first, the fossil will not be able to branch within a clade containing no data, and second, the fossil will have a higher probability, at random, of branching within the clade containing most of the available data. This means that fossils with uncertain phylogenetic affinities (*incertae sedis*) will have a higher probability of branching within the most sampled clade just by chance. Our results suggest that this may be an issue, at the genus-level, in Carnivora, Cetartiodactyla, Chiroptera and Soricomorpha. For example, a Carnivora fossil will be unable to branch in the Herpestidae that has no species with cladistic data, and will also have more chance to branch, randomly, within the Canidae clade than any other clade in Carnivora (Figure 2.1B). Thus, in Total Evidence trees, placements of some carnivoran fossils should be considered with caution. In this study, we treated all cladistic matrices as equal in a similar way to molecular matrices. For example, if matrix A contained 100 characters for taxa X and Y, and matrix B contained 50 different characters for taxa X and Z, we assumed that both matrices can be combined in a supermatrix containing 150 independent characters for taxon X, 100 for taxon Y and 50 for taxon Z. Unfortunately, cladistic data cannot always be treated in this way because some characters may overlap. For example, if matrix A has a character coding for the shape of a particular morphological feature and matrix B has a character coding for the presence of this same morphological feature and a second character coding for its shape, then these three characters are non-independent compound characters (Brazeau, 2011). However, in reasonably sized matrices (>100 characters; (Guillerme and Cooper, 2015; Harrison and Larsson, 2015)) it is more likely that a number of characters are consistently conserved among the different matrices and thus easily combinable. For example, within the Primate cladistic literature, the character *p7* - the size of the 4th lower premolar paraconid - has been used consistently for >15 years (e.g. Ross et al., 1998; ?; Ni et al., 2013) and can therefore be combined among the matrices. A conservative approach to avoid compound characters would be to select only the most recent matrix for each group, but this would result in the loss of a lot of data.

Despite the absence of good cladistic data coverage for living mammals, the Total Evidence methods still seems to be the most promising way of combining living and fossil data for macroevolutionary analyses. Following the recommendations in (Guillerme and Cooper, 2015), we need to code cladistic characters for as many living species possible. Fortunately, data for living mammals is usually readily available in natural history collections, therefore, we propose that an increased effort be

put into coding morphological characters from living species, possibly by engaging in collaborative data collection projects through web portals such as *Morphobank* (O'Leary and Kaufman, 2011). Such an effort would be valuable not only to phylogeneticists, but also to any researcher focusing understanding macroevolutionary patterns and processes.

CHAPTER 3

DISCUSSION

Let's talk...

BIBLIOGRAPHY

- Beck, R. M. and M. S. Lee. 2014. Ancient dates or accelerated rates? Morphological clocks and the antiquity of placental mammals. *Proceedings of the Royal Society B: Biological Sciences* 281:1–10.
- Bininda-Emonds, O. R. P., M. Cardillo, K. E. Jones, R. D. E. MacPhee, R. M. D. Beck, R. Grenyer, S. A. Price, R. A. Vos, J. L. Gittleman, and A. Purvis. 2007. The delayed rise of present-day mammals. *Nature* 446:507–512.
- Brazeau, M. D. 2011. Problematic character coding methods in morphology and their effects. *Biological Journal of the Linnean Society* 104:489–498.
- Cooper, N., J. Rodríguez, and A. Purvis. 2008. A common tendency for phylogenetic overdispersion in mammalian assemblages. *Proceedings of the Royal Society of London B: Biological Sciences* 275:2031–2037.
- Eernisse, D. and A. Kluge. 1993. Taxonomic congruence versus total evidence, and amniote phylogeny inferred from fossils, molecules, and morphology. *Molecular Biology and Evolution* 10:1170–1195.
- Fritz, S. A., J. Schnitzler, J. T. Eronen, C. Hof, K. Böhning-Gaese, and C. H. Graham. 2013. Diversity in time and space: wanted dead and alive. *Trends in Ecology and Evolution* 28:509 – 516.
- Guillerme, T. and N. Cooper. 2015. Effects of missing data on topological inference using a total evidence approach. *Molecular Phylogenetics and Evolution* X:X.
- Harrison, L. B. and H. C. E. Larsson. 2015. Among-character rate variation distributions in phylogenetic analysis of discrete morphological characters. *Systematic Biology* 64:307–324.
- Kearney, M. 2002. Fragmentary taxa, missing data, and ambiguity: mistaken assumptions and conclusions. *Systematic Biology* 51:369–381.
- Kembel, S., P. Cowan, M. Helmus, W. Cornwell, H. Morlon, D. Ackerly, S. Blomberg, and C. Webb. 2010. Picante: R tools for integrating phylogenies and ecology. *Bioinformatics* 26:1463–1464.
- Meredith, R., J. Janečka, J. Gatesy, O. Ryder, C. Fisher, E. Teeling, A. Goodbla, E. Eizirik, T. L. Simão, T. Stadler, D. Rabosky, R. Honeycutt, J. Flynn, C. Ingram, C. Steiner, T. Williams, T. Robinson, B. Angela, M. Westerman, N. Ayoub, M. Springer, and W. Murphy. 2011. Impacts of the Cretaceous terrestrial revolution and KPg extinction on mammal diversification. *Science* 334:521–524.
- Meseguer, A. S., J. M. Lobo, R. Ree, D. J. Beerling, and I. Sanmartín. 2015. Integrating fossils, phylogenies, and niche models into biogeography to reveal ancient evolutionary history: The case of hypericum (hypericaceae). *Systematic Biology* 64:215–232.
- Ni, X., D. L. Gebo, M. Dagosto, J. Meng, P. Tafforeau, J. J. Flynn, and K. C. Beard. 2013. The oldest known primate skeleton and early haplorhine evolution. *Nature* 498:60–64.

- O'Leary, M. A., J. I. Bloch, J. J. Flynn, T. J. Gaudin, A. Giallombardo, N. P. Giannini, S. L. Goldberg, B. P. Kraatz, Z.-X. Luo, J. Meng, X. Ni, M. J. Novacek, F. A. Perini, Z. S. Randall, G. W. Rougier, E. J. Sargis, M. T. Silcox, N. B. Simmons, M. Spaulding, P. M. Velazco, M. Weksler, J. R. Wible, and A. L. Cirranello. 2013. The placental mammal ancestor and the post-K-Pg radiation of placentals. *Science* 339:662–667.
- O'Leary, M. A. and S. Kaufman. 2011. Morphobank: phylophenomics in the cloud. *Cladistics* 27:529–537.
- Pyrón, R. 2011. Divergence time estimation using fossils as terminal taxa and the origins of Lissamphibia. *Systematic Biology* 60:466–481.
- R Core Team. 2015. R: a language and environment for statistical computing. R Foundation for Statistical Computing Vienna, Austria.
- Ronquist, F., S. Klopstein, L. Villhelmsen, S. Schulmeister, D. Murray, and A. Rasnitsyn. 2012. A total-evidence approach to dating with fossils, applied to the early radiation of the Hymenoptera. *Systematic Biology* 61:973–999.
- Ross, C., B. Williams, and R. F. Kay. 1998. Phylogenetic analysis of anthropoid relationships. *Journal of Human Evolution* 35:221–306.
- Sansom, R. S. and M. A. Wills. 2013. Fossilization causes organisms to appear erroneously primitive by distorting evolutionary trees. *Scientific Reports* 3:1–5.
- Schrägo, C., B. Mello, and A. Soares. 2013. Combining fossil and molecular data to date the diversification of New World Primates. *Journal of Evolutionary Biology* 26:2438–2446.
- Slater, G. J. 2013. Phylogenetic evidence for a shift in the mode of mammalian body size evolution at the cretaceous-palaeogene boundary. *Methods in Ecology and Evolution* 4:734–744.
- Slater, G. J. and L. J. Harmon. 2013. Unifying fossils and phylogenies for comparative analyses of diversification and trait evolution. *Methods in Ecology and Evolution* 4:699–702.
- Wagner, P. J. 2000. Exhaustion of morphologic character states among fossil taxa. *Evolution* 54:365–386.
- Webb, C. O., D. D. Ackerly, M. A. McPeck, and M. J. Donoghue. 2002. Phylogenies and community ecology. *Annual review of ecology and systematics* Pages 475–505.
- Wilson, D. E. and D. M. Reeder. 2005. Mammal species of the world: a taxonomic and geographic reference vol. 1. JHU Press.
- Wood, H. M., N. J. Matzke, R. G. Gillespie, and C. E. Griswold. 2013. Treating fossils as terminal taxa in divergence time estimation reveals ancient vicariance patterns in the palpimanoid spiders. *Systematic Biology* 62:264–284.