# EVOLUTION
### INTERNATIONAL JOURNAL OF ORGANIC EVOLUTION

# Phylogenetic Uncertainty Revisited: Implications for Ecological Analyses

**(A)** True Tree (unknown)

**(B)** Consensus tree (missing species)

**(C)** Polytomic super-tree (no missing species)

**(D)** Replication of trees (missing species)

**(E)** Replication of trees (no missing species)

**(A)** True Tree

**(B)** Molecular Data

A: ATTCCGGGATTTCCGATC
B: ??????????????????????
C: TTTCCAAATCTCTCAAAG
D: CCAATTTAACCTGGGTAC
E: CCCAAATTAATCCCAAGG
F: CCAAACCCCATTATTAGG
G: ??????????????????????
H: ATTATTATGGAGGAGCCC

**(C)** Expert Opinion

"PUT B is sister of C or D, PUT G is sister of F or H"

"PUT B is sister of A, PUT G is sister of F"

**(D)** Backbone Trees (inferred)

**(E)** Operational Trees (randomized)

**(A)** MDCCs   **(B)** Polytomies   **(C)** Resolved
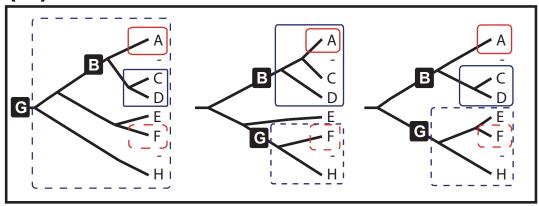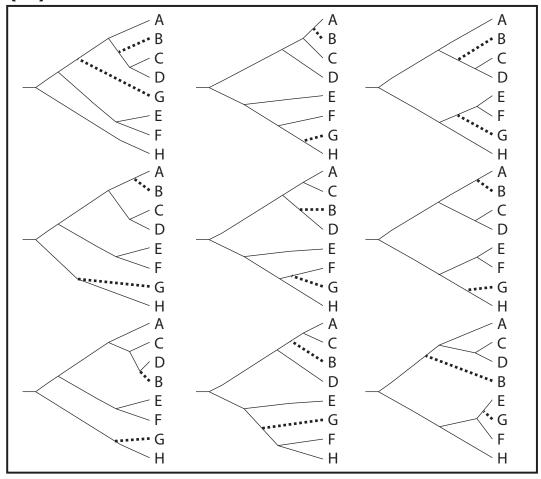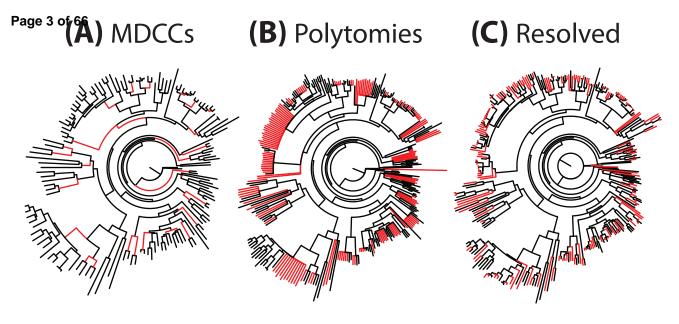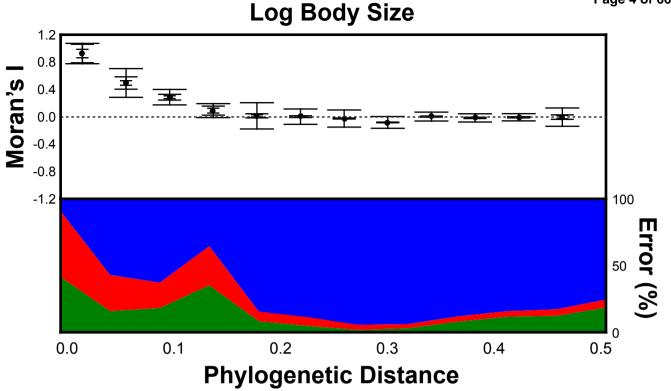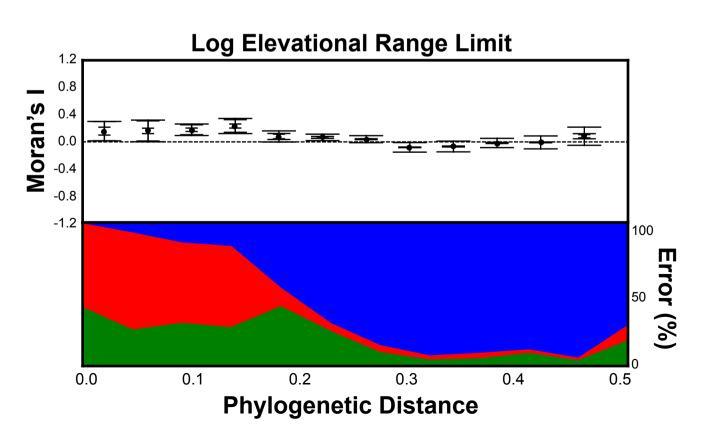
Log Body Size



Log Elevational Range Limit

For: Evolution

Correspondence to:
Thiago Fernando Rangel
Departamento de Ecologia
Universidade Federal de Goiás
Cx.P: 131, 74001-970
Goiânia, Goiás, Brasil
thiago.rangel@ufg.br

Article:

# Phylogenetic Uncertainty Revisited:

# Implications for Ecological Analyses

Thiago F. Rangel [1]
Robert K. Colwell [1,2,3]
Gary R. Graves [4,5]
Karolina Fučíková [2]
Carsten Rahbek [4]
José Alexandre F. Diniz-Filho [1]

1: Departmento de Ecologia, Universidade Federal de Goiás, CP 131, 74.001-970, Goiânia, GO, Brasil. TFR: thiago.rangel@ufg.br. JAFD-F: diniz@ufg.br.

2: Department of Ecology and Evolutionary Biology, University of Connecticut, 75 North Eagleville Rd., Unit 3043, 06269-3043, Storrs, CT, United States. RKC: robert.colwell@uconn.edu. KF: karolina.fucikova@uconn.edu.

3. University of Colorado Museum of Natural History, Boulder, Colorado 80309, USA

4: Department of Vertebrate Zoology, National Museum of Natural History, Smithsonian Institution, Washington, DC 20013, United States. gravesg@si.edu

5: Center for Macroecology, Evolution and Climate, Department of Biology, University of Copenhagen, Universitetsparken 15, 2100, Copenhagen O, Denmark. crahbek@bio.ku.dk

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

1 **Abstract**

2 Ecologists and biogeographers usually rely on a single phylogenetic tree to study evolutionary

3 processes that affect macroecological patterns. This approach ignores the fact that each

4 phylogenetic tree is a hypothesis about the evolutionary history of a clade, and cannot be directly

5 observed in nature. Also, trees often leave out many extant species, or include missing species as

6 polytomies because of a lack of information on the relationship among taxa. Still, researchers

7 usually do not quantify the effects of phylogenetic uncertainty in ecological analyses. We propose

8 here a novel analytical strategy to maximizes the use of incomplete phylogenetic information, while

9 simultaneously accounting for several sources of phylogenetic uncertainty that may distort

10 statistical inferences about evolutionary processes. We illustrate the approach using a clade-wide

11 analysis of the hummingbirds, evaluating how different sources of uncertainty affect several

12 phylogenetic comparative analyses of trait evolution and biogeographic patterns. Although no

13 statistical approximation can fully substitute for a complete and robust phylogeny, the method we

14 describe and illustrate enables researchers to broaden the number of clades for which studies

15 informed by evolutionary relationships are possible, while allowing the estimation and control of

16 statistical error that arises from phylogenetic uncertainty. Software tools to carry out the necessary

17 computations are offered.

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

18   **Introduction**

19   Ecological and biogeographical studies often involve hundreds of species, representing both recent

20   and ancient lineages, and both locally endemic and cosmopolitan species. Although the geographic

21   distribution of most groups of terrestrial vertebrates is increasingly well known, species sampling

22   for phylogenetic analysis is rarely complete for larger clades. Even when a comprehensive

23   phylogeny is available for ecological analyses, key sources of phylogenetic uncertainty are seldom

24   taken into account. We begin this paper by outlining the sources of uncertainty in studies that rely

25   on phylogenetic hypotheses to infer evolutionary processes, and how these sources have been

26   considered or ignored in previous studies. With this motivation, we propose here a novel analytical

27   strategy to quantify and account for key sources of phylogenetic uncertainty in any study that uses

28   phylogenetic input data. As an example of the application of our methodology, we implement it to

29   investigate patterns of trait evolution and phylogenetic assemblages in hummingbirds. We contend

30   that any study that infers evolutionary hypotheses should aim to account for all sources of

31   phylogenetic uncertainty. Finally, we offer freely available software tools to help researchers

32   account for phylogenetic uncertainty in their own research.

33   Phylogenetic uncertainty may arise from two distinct sources (FitzJohn et al. 2009, Diniz-

34   Filho et al. 2013): (1) weak, missing, or conflicting empirical support for hypothesized relationships

35   among species in a given clade (e.g. tree topology, branch length estimation and absolute time

36   calibration), which can be expressed in the form of multiple alternative topologies (e.g., resulting

37   from different analysis methods or different data types), polytomic clades, or low branch support

38   values; and (2) incomplete and unrepresentative sampling of known species. Most ecological

39   studies implicitly assume absolute knowledge of phylogenetic history by simply ignoring the

40   uncertainty caused by the lack of empirical support for phylogenetic trees or issues related to branch

41   length estimation (Fig. 1). When contrasting phylogenetic hypotheses are available, ecologists

3

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

42 usually rely on a consensus tree to "average" phylogenetic information, thus failing to take account

43 of variation among trees (an expression of overall phylogenetic uncertainty). In addition, although

44 ecologists acknowledge that most polytomies are expressions of ambiguous or missing empirical

45 data, highly polytomic trees are nonetheless often used in ecological analyses without proper

46 quantification of phylogenetic uncertainty and without conducting the necessary sensitivity

47 analyses.

48     Ecologists have dealt with incomplete phylogenies in several ways. The first approach is to

49 focus only on clades for which relatively complete phylogenies are available, but this strategy

50 restricts ecological studies to a very small number of groups (Pagel 1999) and can undermine

51 assemblage-level or macroecological studies (Webb and Donoghue 2005). A favored method is to

52 assemble supertrees from smaller, overlapping trees and to fill gaps in phylogenies by placing

53 unsampled species in trees according to their taxonomic classification (Bininda-Emonds 2004,

54 Webb and Donoghue 2005, Hernandez and Vrba 2005, Davies et al. 2004, Ranwez et al. 2007).

55 Nevertheless, even for the best-studied taxa, such as mammals and birds, supertrees that span all

56 species are relatively recent achievements (e.g. Bininda-Emonds et al. 1999, Beck et al. 2006).

57 Unfortunately, supertrees are usually highly polytomic (not fully resolved). Recent approaches

58 based on building megatrees (Roquet et al. 2012) or complex addition of species on backbone trees

59 (Jetz et al. 2012) are still in their infancy. Phylomatic software (Webb and Donoghue 2005), for

60 instance, is a popular tool for assemblage-level ecological studies that constructs customized trees

61 of virtually any size. Although Phylomatic allows additional phylogenetic representations of the

62 clade when compared to the backbone phylogeny, terminal branches are usually based on

63 taxonomic relationships among inserted species and thus are commonly polytomic. Additional

64 drawbacks of supertrees are the absence of information on branch lengths (which require further

65 effort, and assume confidence in the fossil record for time calibration) and, again, the lack of

4

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

66    explicit conversion of phylogenetic uncertainty into explicit measures of error associated with

67    statistical inference and parameter estimates.

68         A second and more radical strategy consists of ignoring the species that are absent from the

69    available phylogeny, under the assumption that the species included in the analysis represent an

70    unbiased and representative sample of all species in the clade (FitzJohn et al. 2009). Of course, the

71    full evolutionary history of a clade of substantial age can be described only by a phylogeny with all

72    species (including extinct ones), although this ideal is achievable only in simulated scenarios

73    (Colwell and Rangel 2010). Estimating the degree of bias due to missing species can, in principle,

74    be achieved by replicating the analysis with random sub-samples of species that are present in the

75    phylogeny (rarefaction or "thinning", Davies et al. 2012), but unfortunately this approach is not

76    common in evolutionary or ecological studies.

77         Several other approaches to deal with phylogenetic uncertainty have been recently

78    developed, increasing awareness that species missing from phylogenetic trees may seriously affect

79    statistical inference of evolutionary processes, even when the missing species are inserted in the tree

80    in the form of polytomies (Davies et al. 2012, Diniz-Filho et al. 2013). For example, under the

81    Bayesian framework of molecular phylogenetic reconstruction, missing species can be inserted by

82    assigning empty sequences to missing species and by constraining their insertion using priors on the

83    tree topology (Huelsenbeck and Rannala 2003). Several recent studies have employed different

84    analytical strategies to account for phylogenetic uncertainty. For example, Isaac et al. (2007) dealt

85    with missing species in their analysis of evolutionary distinctiveness of mammals by allocating

86    missing species among their presumed closest relatives using a model of constant rate of speciation

87    and extinction. Day et al. (2008) studied the diversification rates of cichlid fish radiation of Lake

88    Tanganyika, accounting for the potential effects of missing species on the estimates of the timing

89    and rate of diversification. Kuhn et al. (2011) proposed a method and a computational tool that uses

90    a birth-death model of diversification to resolve polytomies and to define not only tree topology,

5

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

91    but also branch lengths in polytomic trees, and Jetz et al. (2012) built sets of phylogenetic trees for

92    the entire avian clade, combining genetic and taxonomic information and adding a simulation of

93    evolutionary diversification to better establish branch lengths and resolve polytomies. Davies et al.

94    (2012) showed that phylogenetic uncertainty can dramatically inflate estimates of phylogenetic

95    signal and proposed a rarefaction-based method to guide inference. After accounting for

96    phylogenetic uncertainty, Batista et al. (2013) showed that potential loss of evolutionary history

97    caused by extinction of threatened Western Hemisphere anurans would not be significantly higher

98    than that of non-threated anurans. These few examples clearly show the importance of dealing with

99    phylogenetic uncertainty, and demonstrate that systematists and comparative biologists have

100   already begun to develop and employ methods to incorporate estimates of uncertainty in the

101   inference of diversification rates (Day et al. 2008, FitzJohn et al. 2009), character evolution (Losos

102   1994, Huelsenbeck et al. 2000, Housworth and Martins 2001, Huelsenbeck and Rannala 2003, Ives

103   et al. 2007), reconstruction of ancestral states (Ronquist 2004), and tree topology (Felsenstein 1985,

104   Holder and Lewis 2003).

105          Here we propose a unified analytical approach to estimate and account for multiple sources

106   of phylogenetic uncertainty. Our method consists of partitioning variance among estimated

107   parameters of evolutionary processes, using random sampling from the universe of probable

108   phylogenies. We illustrate the approach using a clade-wide analysis of the hummingbirds,

109   evaluating how different sources of uncertainty affect several phylogenetic comparative analyses of

110   trait evolution and biogeographic patterns.

111

## Accounting for multiple sources of phylogenetic uncertainty

113   We propose here an empirical approach, based on simulations, that maximizes the use of

114   incomplete phylogenetic information in ecological studies, while simultaneously accounting for

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

115    several sources of phylogenetic uncertainty that may distort statistical inferences about evolutionary

116    and ecological processes.

117

118

119    *Missing species*

120    Here we define a *phylogenetically uncertain taxon* (henceforth a "PUT" for a single taxon, and

121    "PUTs" for multiple taxa) as a taxonomic unit (e.g. population, species, genus) that is recognized as

122    valid and is accepted as belonging to a particular clade, but is missing from the available

123    phylogenetic tree(s) for that clade. The absence from the tree could be due to multiple causes, such

124    as unavailability of molecular and morphological data (Fig. 2b). However, although the data

125    required to formally reconstruct the phylogenetic relationships of a PUT may be missing,

126    evolutionary information about the PUT is rarely completely non-existent, as no authority would

127    dispute that some other taxon in the clade must be the closest relative of the PUT at some level in

128    the phylogeny. In fact, additional available information, such as taxonomic, morphological or

129    behavioral data, may be useful to define a list of the most likely sister species or lineages of any

130    given PUT (Fig. 2c). Therefore, it makes sense to use all acceptable information available to

131    conservatively define what we will call the *most derived consensus clade* (MDCC), i.e., the node

132    that unequivocally contains each PUT.

133        Thus, an MDCC can be defined as the most recent common ancestor of all unquestioned

134    candidates for being the closest relative of the PUT (Fig. 2d). Of course, the validity of information

135    (e.g. taxonomic, biogeographical, behavioral, morphological, etc) used to define a MDCC can be

136    questioned. However, if two taxonomists, for example, disagree on the proper MDCC for a given

137    PUT, instead of eliminating the PUT from the statistical analysis, thereby ignoring phylogenetic

138    uncertainty, the MDCC should be conservatively redefined as the most recent common ancestor of

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

139    the two candidate MDCCs for the PUT, each championed by a different taxonomist (Fig. 2). Thus,

140    it is important to stress that the goal of defining a MDCC to a PUT is not to replace the formal

141    methods of phylogenetic reconstruction, but to conservatively use all available information to allow

142    biologists to make inference of ecological and evolutionary processes using all available

143    information, even with incomplete phylogenies.

144        We propose here a simple process of building a possible phylogeny that includes all extant

145    species. This process relies on the definition of an MDCC for each PUT and the insertion of the

146    PUT in a random position within its MDCC. We begin by randomizing the order in which PUTs are

147    to be added to the tree. Next, within each MDCC, we assign the PUT to a point along one a branch

148    of the clade. The choice of insertion point along a branch can be uniform random if no additional

149    information is available, or it may be guided by a biologically realistic model of diversification

150    (Kuhn et al. 2011, Davies et al. 2012, Jetz et al. 2012). If the baseline tree is ultrametric, the branch

151    length for the inserted PUT is simply the distance from the attachment point to the end of all other

152    tips. If, however, the baseline tree is not ultrametric, the branch length of the PUT can be sampled

153    from a distribution of possible branch length values. Once a PUT has been inserted, its own branch

154    may serve as a potential insertion point for subsequent species assigned from the PUT queue. Our

155    algorithm iterates until each PUT has been sequentially added to the appropriate MDCC, producing

156    a complete phylogeny of known species (Martins et al. 2013, Batista et al. 2013).

157

158    *Polytomies*

159    If all polytomies in a tree are the consequence of phylogenetic uncertainty (i.e. "soft" polytomies),

160    then it is necessary to ensure that such uncertainty is quantified and accounted for in statistical

161    analyses that use such a tree (Lewis et al. 2005). To explore the space of all possible dichotomous

162    trees one must resolve the polytomies, assuming no additional information about the evolutionary

8

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

163      history of the clade (Batista et al. 2013), or employing a model of diversification (Martins 1996,

164      Housworth and Martins 2001, Kuhn et al. 2011).

165          Our approach consists in producing fully dichotomous trees by resolving polytomies

166      stochastically. For each node in the tree with three or more branches, we choose two branches at

167      random and reassign them to the same node. Each remaining branch from the original polytomy is

168      then inserted sequentially, in random order, within the clade constructed from the former polytomy,

169      at a randomly chosen position along the length of the existing branches. Dichotomous nodes in the

170      original phylogeny are not changed, thereby preserving the phylogenetic information in the original

171      baseline tree. This process guarantees that the resulting tree is fully resolved and that the

172      phylogenetic uncertainty arising from polytomies is also taken into account when multiple

173      randomized trees are compared (Batista et al. 2013). As with sampling multiple phylogenetic trees

174      (below) or randomizing the position of the PUTs in the phylogeny, not every possible tree topology

175      is examined. However, given a large enough number of replicates, our stochastic procedure ensures

176      that the parameter space will be explored, thereby accounting for phylogenetic uncertainty.

177

178      *Multiple phylogenetic trees*

179      Modern methods of phylogenetic reconstruction and inference are based on searching the universe

180      of possible phylogenetic trees. The search is guided by empirical data among possible trees, which

181      are judged by their ability to predict the observed data, given a model of molecular evolution

182      (Holder and Lewis 2003). However, because of scarce, ambiguous or missing empirical data,

183      searches are often unable to unambiguously rank all trees within an entire group of equally likely

184      trees, therefore yielding a large set of possible trees. Moreover, especially analyses of large,

185      genome-scale data sets may yield multiple yet well-supported topologies depending on the

186      analytical approach or on selected subsets of data (e.g., Xi et al. 2014, examples in Cooper 2014).

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

187    To account for the phylogenetic uncertainty that arises from lack of definitive empirical support for

188    a single phylogenetic hypothesis, one must not limit the statistical analysis to a single sampled tree

189    or to a majority-rule consensus tree that "averages" a larger set of possible trees (Fig. 1). Sampling

190    or averaging among possible trees can potentially improve a point-estimate of a parameter, but does

191    not account for statistical error in evolutionary inference, thereby masking phylogenetic uncertainty

192    associated with tree reconstruction. In practice, variation in the results of ecological analyses that

193    use different phylogenetic trees is an expression of uncertainty in phylogenetic reconstruction.

194    Thus, it is necessary to replicate these analyses over a large number of possible phylogenetic trees,

195    and subsequently study the variation in parameter estimates that arises as a consequence of

196    differences among trees (Fig. 2).

197

198    *Uncertainty and sensitivity analysis*

199    Computer simulations have long been used to estimate the magnitude of statistical error and to

200    quantify the probability of a given outcome when a system is not fully known (Doubilet et al. 1998,

201    Saltelli et al. 2008). Usually run in tandem with uncertainty quantification, a sensitivity analysis

202    allows one to determine the degree to which the sources of uncertainty in model inputs are

203    responsible for the uncertainty in the model output (Pannell 1997). Coupling uncertainty and

204    sensitivity analysis offers us the opportunity to quantify the robustness of models in the presence of

205    phylogenetic uncertainty, account for uncertainty in evolutionary inference, and enhance

206    communication of the magnitude of statistical error in the study.

207        Martins (1996) proposed a way to carry out phylogenetic comparative studies when the

208    phylogenetic relationships among species are unknown. In her method, a large sample of trees is

209    generated by randomly resolving a fully polytomic tree ("star phylogeny"), using models of

210    phenotypic evolution and diversification rates (so that uncertainty in branch lengths is also

10

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

211    implicitly incorporated). Species traits are then analyzed on each of the possible trees. Although the

212    mean results of such analyses will converge to a non-phylogenetic analysis (Abouheif 1998), the

213    approach of Martins (1996) represents a landmark in phylogenetic inference under uncertainty

214    present in an unresolved phylogeny (see also Housworth and Martins 2001).  The mean of the

215    squared standard error of the calculated evolutionary statistic (e.g. the correlation between two

216    species' traits), known as $V_s$, estimates the true variance of the statistic, which is due to sample

217    variance. The variance of the evolutionary statistic calculated among the randomly generated trees,

218    known as $V_p$, estimates the variance due to phylogenetic uncertainty. Thus, inferences based on

219    parameter estimation must account for both sources of error, which can be done by computing

220    confidence intervals or $P$-values based on the sum of $V_p$ and $V_s$.

221        We expand Martins' (1996) strategy for incorporating phylogenetic uncertainty in parameter

222    estimation further to accommodate multiple sources of phylogenetic uncertainty. We treat each

223    source of uncertainty as a factor in an experimental design, while parameter estimates are treated as

224    the response variable under study. Thus, we can ask how different phylogenetic trees, alternative

225    resolutions of polytomies, and/or probable configurations of PUTs would affect parameter

226    estimates. We partition the amount of variance in a parameter estimate arising from each source of

227    phylogenetic uncertainty with an Analysis of Variance (ANOVA), which not only isolates variance

228    among sources but also calculate the magnitude of standard error in the parameter estimate (a

229    separate Martins' (1996) $V_p$ for each source of uncertainty).

230        We view the sources of phylogenetic uncertainty we have outlined here as a hierarchy to be

231    employed in the design of an experiment used to assess uncertainty and sensitivity. If multiple

232    empirical phylogenies are available, then this source of uncertainty can be regarded as the highest

233    level of uncertainty in the analysis. Because each empirical phylogeny is unique, its polytomies

234    must be treated individually. Thus, polytomies can be regarded as the second level of phylogenetic

235    uncertainty, immediately below the multiple empirical phylogenies. Of course, for each available

11

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

236  empirical phylogeny there is a virtually unlimited number of alternative ways to resolve polytomies

237  randomly or by following some diversification model. Finally, for each combination of original

238  phylogeny and polytomy resolution, PUTs can be inserted according to the procedure described

239  above. Thus, insertion of PUTs is the lowest level in the hierarchy of sources of phylogenetic

240  uncertainty. Because sources of uncertainty are nested within a higher level of classification, and

241  groups representing subordinate levels are randomly chosen, a nested (hierarchic) analysis of

242  variance is ideal for the analysis of sensitivity and uncertainty.

243

244  **Application example: trait evolution and phylogenetic assemblage patterns in**

245  **hummingbirds**

246  We illustrate our method by applying it to ecological hypotheses for the hummingbirds

247  (Trochilidae), a large, monophyletic clade (~330 species) with a rich natural history literature and

248  ongoing phylogenetic, morphological, behavioral, and ecological research. McGuire et al. (2007)

249  published a multilocus molecular phylogeny for the hummingbirds that includes only 146 species,

250  but encompasses all the higher-level trochilid diversity (73 of the approximately 104 recognized

251  genera, Schuchmann 1999). Although a more complete hummingbird phylogeny is now available

252  (McGuire et al. 2014), we use the earlier phylogeny(McGuire et al. 2007) to demonstrate our

253  methods because it is typical of current level of phylogenetic knowledge for many comparably

254  diverse taxa.

255       Replicating the phylogenetic reconstruction methodology of McGuire et al. (2007), we

256  sampled 25,000 trees from the posterior distribution (hereafter called "backbone" trees), after

257  discarding a burn-in of 5,000 steps. We relied upon the taxonomic classification of Schuchmann

258  (1999) to designate the most likely MDCC for each PUT (Fig. 3).

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

259   From the total of about 330 described species of hummingbirds, we compiled geographical

260 distributions and ecologically important morphological data for 304 species. Species endemic to

261 islands (the West Indies and the Juan Fernández Archipelago) were not included in the analysis. We

262 compiled estimates of average body mass (intersexual mean) for each of the 304 species from

263 Dunning (2007) and Schuchmann (1999). Although large intersexual and geographical variation in

264 hummingbird body size is well documented (Colwell 2000), the averages used here are useful for

265 broad taxonomic analyses. Based on the literature, we also recorded the maximum known

266 elevational range limits for each species (see Appendix A for references). Body mass (as a measure

267 of body size) and maximum elevational range limit were log-transformed to conform to a normal

268 distribution.

269   We extracted distributional data for 304 species from an updated version (16 July 2010) of

270 the comprehensive database for all land and fresh-water birds known to have breeding populations

271 in the Western Hemisphere, complied by Rahbek and Graves (2000, 2001), and mapped the

272 geographical range of each species on a gridded map at a resolution of 1º x 1º (latitude-longitude).

273 These maps represent a conservative extent-of-occurrence estimate of the breeding range based on

274 museum specimens, published sight records, and spatial distribution of habitats based on

275 documented records for South America.

276

277 *Uncertainty quantification*

278 Variability in input data is a fundamental source of uncertainty. To estimate the degree of

279 phylogenetic uncertainty in the hummingbird data due to phylogeny reconstruction and missing

280 species we randomly sampled 100 trees from the posterior (here after called "backbone" trees), and

281 for each tree we replicated the insertion of PUTs 100 times, thereby generating a total of 10,000

282 fully resolved phylogenetic trees. Next, we calculated the tree-to-tree pairwise distance matrix using

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

283     the weighted Robinson-Foulds (wRF) metric, which measures the minimum number of internal

284     branches that must be collapsed or expanded to make two trees identical. The weighted version of

285     Robinson-Foulds distance accounts for differences in branch lengths between trees (Steel and

286     Penny 1993). Distance matrices of phylogenetic trees have been used in phylogenetic reconstruction

287     to calculate consensus trees (Swofford 1991), produce supertrees (Bansal et al. 2010), and visualize

288     multi-dimensional space of possible trees (Hillis et al. 2005).

289          The average wRF pair-wise distance between the 10,000 trees is 2.7423, with a standard

290     deviation of 0.136. We used a multivariate analysis of variance (PERMANOVA, Anderson 2001)

291     to partition variance among phylogenetic trees through the pair-wise wRF distance matrix. The

292     estimated components of variance indicated that 17.02% of total variance among trees is

293     attributable to phylogenetic reconstruction (i.e. differences among backbone trees), whereas 82.97%

294     of total variance is due to phylogenetic uncertainty of missing species (PUTs).

295

296     *Phylogenetic autocorrelation in species traits*

297     We used Moran's I correlograms to estimate the magnitude of phylogenetic autocorrelation in body

298     size and elevational range limits among hummingbird species (see Gittleman and Kot 1990; Diniz-

299     Filho 2001; Pavoine and Ricota 2012). Moran's I is larger when species within a given phylogenetic

300     distance interval have similar trait values and smaller when species within the same distance

301     interval have very different trait values. To account for uncertainty in phylogenetic reconstruction

302     and missing species, the calculation of correlograms was repeated 1,000,000 times (i.e., replicating

303     the insertion of PUTs 1,000 times in each of the 1,000 randomly sampled backbone phylogenies).

304     We used a nested ANOVA to partition sources of uncertainty for each of the 12 distance classes in

305     the correlogram. Applying our method, the average Moran's I within a distance class is the

14

306   parameter estimate, whereas additive confidence intervals are calculated around the estimate for

307   each source of uncertainty.

308       Results of these analyses showed that phylogenetic autocorrelation in body size and

309   maximum elevational range limit both decline with phylogenetic distance, although at different

310   rates (Fig. 4, top panel in each plot). Closely related species tend to have similar body sizes but less

311   similar elevational range limits. Accounting for phylogenetic uncertainty widens confidence

312   intervals by at least a factor of two, regardless of trait (Fig. 4, confidence intervals in upper panel;

313   see the caption for details). For all distance classes, standard confidence intervals due to sampling

314   error are always narrower than confidence intervals due to phylogenetic uncertainty. In addition, the

315   magnitude of error attributed to each source of uncertainty is not constant over phylogenetic

316   distance (Fig. 4, bottom panels). The relative proportion of error caused by missing species tends to

317   be higher in short distance classes, as taxonomic information has been used to assign PUTs to

318   relatively derived positions in the phylogeny (MDCCs). In contrast, large phylogenetic confidence

319   intervals at intermediate and long distance classes arise from uncertainty in the empirical phylogeny

320   located at the base of the coquettes and brilliants clades, which together form the Andean clade

321   (McGuire et al. 2007).

322

323   *Phylogenetic signal and evolutionary models in species traits*

324   We used Blomberg et al.'s (2003) *K* statistic (hereafter *K*) to estimate the magnitude of deviation

325   from Brownian motion in body size and elevational range limits. *K* measures the degree of

326   similarity in species' traits in relation to the similarity expected under a Brownian motion model of

327   phenotypic evolution, given a phylogenetic hypothesis. *K* is based on the ratio between two

328   measures of distance: $MSE_0$, the squared distance between trait values and the phylogenetically

329   corrected mean trait value, and MSE, the squared distance between trait values estimated from a

15

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

330    variance-covariance matrix derived from a phylogenetic hypothesis. Thus, large values of

331    $MSE_0/MSE$ indicate a strong phylogenetic signal. To allow comparisons between different traits

332    and trees, observed $MSE_0/MSE$ is standardized by expected $MSE_0/MSE$, assuming that the trait

333    evolves under a Brownian motion model of phenotypic evolution. *K* values less than one indicate

334    that species are less similar for a given trait than expected under Brownian motion evolution,

335    whereas *K* values greater than one indicate that species are more similar for a given trait than

336    expected under Brownian motion evolution.

337        To account for uncertainty in phylogenetic reconstruction and missing species, we again

338    calculated 10,000 possible *K* values, replicating the insertion of PUTs 100 times in each of the 100

339    randomly sampled phylogenies. To estimate the accuracy (standard error) of *K* within each tree we

340    used jackknife permutation for each combination of PUT insertion and sampled phylogeny (Efron

341    and Tibshirani 1994). Finally, we used a nested ANOVA to partition error in estimated *K* among

342    different sources of uncertainty.

343        Our results revealed a surprisingly low phylogenetic signal in body size among

344    hummingbirds ($\overline{K} = 0.0597$); thus body size evolution can hardly be explained by a simple model

345    of Brownian evolution. However, the standard error of *K* due to sample size ($s_{KSamp}=0.0011$) is

346    relatively small (2.23%) compared to the error arising from phylogenetic uncertainty due to missing

347    species ($s_{KSamp+PUT}=0.0494$), as this source of uncertainty represents 97.1% of the total error in *K*.

348    Finally, uncertainty due to phylogenetic reconstruction represents only 0.67% of total error of the

349    evolutionary inference ($s_{KSamp+PUT+Phy}=0.0497$).

350        The relative importance of sources of uncertainty shift discordantly when maximum

351    elevational range limit is considered, although this trait also has very little phylogenetic signal ($\overline{K} =$

352    0.0221). Standard error of *K* due to sampling size contributes 59.7% of total error ($s_{KSamp} = 0.0179$),

16

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

353    whereas error due to missing species represents 40.1% of total error ($s_{KSamp+PUT} = 0.0300$), and error

354    due to phylogenetic reconstruction represents only 0.2% of total error ($s_{KSamp+PUT+Phy}=0.03001$).

355            To test the hypothesis that estimated $K$ values are not significantly different from random

356    expectation, we calculated 100 $Ks$ for each combination of PUT insertion and sampled phylogeny,

357    randomizing trait values among species (Blomberg et al. 2003, Revell et al. 2008). We employed

358    the same nested ANOVA design to estimate standard error of the null distribution due to each

359    source of uncertainty. Finally, we used Welch's $t$-test to evaluate the hypothesis that estimated $Ks$

360    do not differ significantly from the null expectation:

$$t = \frac{\overline{K}_{trait} - \overline{K}_{H0}}{\sqrt{\dfrac{s^2_{K_{trait}}}{n_{K_{trait}}} + \dfrac{s^2_{K_{H0}}}{n_{K_{H0}}}}}$$

361

362    Estimated $K$ for hummingbird body mass is significantly larger than expected under the null

363    expectation ($\overline{K}_{H0} = 0.0237$), and accounting for phylogenetic uncertainty does not affect the

364    hypothesis test ($t_{Samp} = 2589.92$, $t_{Samp+PUT} = 62.5$, $t_{Samp+PUT+Phy} = 76.72$, all $P$-values $< 0.001$).

365    Conversely, phylogenetic signal in maximum elevational range limit is significantly smaller than

366    expected under the null hypothesis ($\overline{K}_{H0} = 0.0333$). Accounting for phylogenetic uncertainty also

367    does not change the inference of statistical significance in hummingbird maximum elevational

368    range limit ($t_{Samp} = -275.36$, $t_{Samp+PUT} = -34.74$, $t_{Samp+PUT+Phy} = -44.47$, all $P$-values $< 0.001$).

369            These results indicate that, although body size evolution deviates significantly from

370    Brownian motion, it carries a small phylogenetic signal. Thus, evolutionary constraints and niche

371    conservatism can be invoked for body size (for instance, evolution under an O-U process – see

372    Hansen et al. 2008), whereas elevational range has less phylogenetic signal than expected by

373    chance.

374

375    *Phylogenetic community structure at the macroecological scale*

376    Phylogenetic species variability (PSV, Helmus et al. 2007) of an assemblage is maximized (PSV =

377    1) when the assemblage is composed of the least related species in a clade, and minimized (PSV =

378    0) when the most related species coexist in an assemblage. We calculated PSV for the 1979

379    assemblages in the Western Hemisphere based on mapped species ranges with at least two

380    hummingbird species, and tested statistical significance of each PSV value using 300 permutations

381    of species identities. However, because phylogenetic distance between hummingbird species is not

382    known with certainty, we generated 300 possible phylogenies with different PUT insertions, and

383    replicated the calculation of PSV for each of the 1979 assemblages using the each of the 300

384    randomly selected phylogenetic trees.

385        Some hummingbird assemblages are composed of species with non-random phylogenetic

386    relationships (Fig. 5). However, identifying any significant departure from randomness requires

387    accounting for all sources of uncertainty. Standard null models to test phylogenetic structure of

388    communities (e.g. Graham et al. 2009), which randomize species identities while preserving species

389    richness, are designed to account only for non-random "sampling" from the phylogeny. For

390    hummingbird assemblages across the Western Hemisphere, an analysis with such a null model

391    would lead to the inference of significant phylogenetic dispersion along the middle and upper

392    Andes (contrary to Graham et al. 2009), and significant phylogenetic clustering across the Pacific

393    coast of Central and North America (Fig 5A). However, uncertainty caused by missing species and

394    phylogenetic reconstruction may seriously affect pattern detectability. Because of substantial

395    uncertainty in the phylogenetic relationship of the two Andean clades (Coquettes and Brilliants),

396    when phylogenetic uncertainty is taken into account no significant phylogenetic dispersion in

397    Andean assemblages is detected. Notice that this result is concordant with the lack of phylogenetic

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

398    signal and lack of phylogenetic autocorrelation for elevational range at the species level, as we

399    previously discussed. Moreover, once phylogenetic uncertainly is taken into account, many

400    assemblages in Central America can no longer be considered phylogenetically clustered (Fig. 5B).

401    Because species are neither randomly distributed in the phylogeny nor in geographic space,

402    the magnitude of error in the analysis of phylogenetic structure of assemblages tends to be strongly

403    spatially autocorrelated. In addition, sampling effect is usually higher in species-poor assemblages.

404    On the other hand, assemblages with a higher proportion of species that are members of clades

405    characterized by phylogenetic uncertainty in deep (basal) nodes are subject to a higher proportion of

406    error due to phylogenetic reconstruction. Figure 6 depicts the relative contribution of each source of

407    error in analysis the analysis of phylogenetic structure of hummingbird assemblages. Because of

408    phylogenetic uncertainty in the reconstruction of the relationship between the two, basal, Andean

409    clades (Coquetes and Brilliants), statistical error is substantially higher than the other two sources of

410    error for Andean assemblages. Conversely, statistical error in North American assemblages is

411    mostly due to sampling, as species richness is very low. Uncertainty due to missing species (PUTs)

412    is not particularly pronounced in any assemblage, as no assemblage is composed by more than 17%

413    PUTs.

414

415    **Concluding Remarks**

416    Phylogenetic uncertainty is not evenly distributed across time and space. As a consequence, the

417    effects of phylogenetic uncertainty in the statistical analysis of phylogenetic data cannot be

418    estimated without a thorough sensitivity analysis on a case-by-case basis. To illustrate the new

419    methods we propose to account for phylogenetic uncertainty, in this paper we used phylogenetic

420    hypotheses derived from molecular data to analyze the hummingbird clade, a well-studied

421    taxonomic group widely accepted as monophyletic. Partition of variance among sources of

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

422     uncertainty reveals that variance among simulated phylogenies is caused primarily by missing

423     species, in this case, even using the best information available to assign phylogenetically uncertain

424     taxa (PUTs). Conversely, variance among backbone phylogenies, which are estimated through

425     molecular data, is substantially smaller, indicating that efforts to improve the knowledge of

426     evolutionary history of hummingbirds (and other clades that share similar patterns of uncertainty)

427     should be concentrated on gathering molecular data for additional species (e.g. McGuire et al.

428     2014).

429          Because phylogenetic uncertainty is expected to be relatively more concentrated within

430     some clades of a phylogeny than within others, statistical analyses of species assemblages or

431     temporal segments of the phylogeny will be affected differently. Thus, statistical analyses of

432     different traits for the same group of species, using the same phylogenetic information, may be

433     differently affected by phylogenetic uncertainty, both in intensity and direction of bias. Of course,

434     comparative analyses among multiple taxonomic groups, based on independently built phylogenetic

435     hypotheses, requires additional caution, as the heterogeneity of variance among groups may

436     seriously distort results in unpredictable directions.

437          Because the effects of phylogenetic uncertainty in ecological and evolutionary analyses are

438     not subject to generalization, quantifying and accounting for phylogenetic uncertainty through

439     sensitivity analysis is required in all ecological studies. Modern techniques of sensitivity analyses

440     typically involve the application of Monte Carlo methods. Although general simulation methods,

441     such as the one used in this study, could be modified to account for phylogenetic uncertainty in

442     most evolutionary and ecological analyses, the framework for uncertainty quantification and

443     sensitivity analysis should be tailored to the purpose of the study (e.g. estimate of diversification

444     rates, community phylogenetics, comparative analysis of species traits).

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

445       Finally, the approach proposed here can be used to quantify the full spectrum of components

446    of phylogenetic uncertainty, guiding sampling strategies for future studies and allowing more

447    reliable interpretations of the relative magnitude of historical and phylogenetic components of

448    biodiversity patterns.

449

## 450    Software Tools

451       We provide two software toolkits to enable the application of the analytical strategy

452    proposed here. The first software is SUNPLIN (Martins et al. 2013;

453    https://sourceforge.net/projects/sunplin), which is capable of generating randomized phylogenies

454    after the insertion PUTs into backbone trees, with MDCCs assigned. The generated trees can then

455    be applied in any analysis that requires phylogenies as input data. SUNPLIN can be used as an

456    online web service (http://wsmartins.net/sunplin/), as a library that connects through APIs to any

457    compiled software, or directly integrated into R (http://www.ecoevol.ufg.br/pam).

458       The second software toolkit, PAM (*Phylogenetic Analysis in Macroecology*,

459    http://www.ecoevol.ufg.br/pam), is a compiled computational platform for inference of ecological

460    and evolutionary processes in a spatially explicit context. In PAM, users can not only generate

461    replicates of phylogenetic trees to be used in other software applications, but can also run several

462    statistical analyses commonly used in biodiversity analysis, while estimating and accounting for

463    multiple sources of uncertainty using the analytical framework proposed here. PAM is a work in

464    progress and will be continuously expanded in the future.

465

## 466    Acknowledgements

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

478

479    **References**

480    Abouheif, E. 1998. Random trees and the comparative method: a cautionary tale. Evolution

481        52:1197-1204.

482    Anderson, M.J. 2001. A new method for non-parametric multivariate analysis of variance. Austral

483        Ecology 26:32-46.

484    Bansal, M. S., J. G. Burleigh, O. Eulenstein, and D. Fernández-Baca. 2010. Robinson-Foulds

485        Supertrees. Algorithms for Molecular Biology 5:18.

486    Batista, M. C. G, S. F. Gouveia, D. L. Silvano and T. F. Rangel. 2013. Spatially explicit analyses

487        highlight idiosyncrasies: species extinctions and the loss of evolutionary history. Diversity

488        and Distributions (*in press*).

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

489    Beck, R. M. D., O. R. P. Bininda-Emonds, M. Cardillo, F.-G., R. Liu and A. Purvis. 2006. A

490            higher-level MRP supertree of placental mammals. Evolutionary Biology 6:93-107.

491    Bininda-Emonds, O. R., J. L. Gittleman and A. Purvis. 1999. Building large trees by combining

492            phylogenetic information: a complete phylogeny of the extant Carnivora (Mammalia).

493            Biological Reviews 74: 143-175.

494    Bininda-Emonds, O. R. P. 2004. The evolution of supertrees. Trends in Ecology and Evolution

495            19:315-322.

496    Blomberg, S. P., T. Garland, and A. R. Ives. 2003. Testing for phylogenetic signal in comparative

497            data: behavioral traits are move labile. Evolution 57:717-745.

498    Colwell, R. K. 2000. Rensch's Rule crosses the line: Convergent allometry of sexual size

499            dimorphism in hummingbirds and flower mites. American Naturalist 156:495-510.

500    Colwell, R. K. and T. F. Rangel. 2010. Modelling quaternary range shifts and richness on tropical

501            elevational gradients. Phylosofical Transactions of the Royal Society, B 365: 3695-3707.

502    Cooper E.D. 2014. Overly simplistic substitution models obscure green plant phylogeny. Trends in

503            Plant Science 19: 576-582.

504    Davies, T. J., N. J. B. Kraft, N. Salamin and E. M. Wolkovich. 2012. Incompletely resolved

505            phylogenetic trees inflate estimates of phylogenetic conservatism. Ecology 93:242-247

506    Davies, T.J., T.G. Barraclough, M.W. Chase, P.S. Soltis, D.E. Soltis and V. Savolainen. 2004

507            Darwin's abominable mystery: insights from a supertree of the angiosperms. Proceedings of

508            the National Academy of Sciences 107, 1904–1909.

509    Day, J. D., J. A. Cotton and T. G. Barraclough. 2008. Tempo and Mode of Diversification of Lake

510            Tanganyika Cichlid Fishes. PLoS One 3:e1730

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

511    Diniz-Filho, J. A. F. 2001. Phylogenetic autocorrelation under distinct evolutionary processes.

512         Evolution 55:1104-1109.

513    Diniz-Filho, J. A. F., R. D. Loyola, P. Raia, A. O. Mooers and L. M. Bini. 2014. Darwinian

514         shortfalls in biodiversity conservation. Trends in Ecology and Evolution, *in press*.

515    Doubilet P., C.B. Begg, M. C. Weinstein, P. Braun, B. J. McNeil. 1985. Probabilistic sensitivity

516         analysis using Monte Carlo simulation. A practical approach. Medical Decision Making

517         5:157-177.

518    Dunning, J. B. 2007. CRC handbook of avian body masses. 2nd ed. CRC press.

519    Efron, B. and R. J. Tibshirani. 1994. An introduction to the bootstrap. Chapman & Hall.

520    Felsensetein, J. 1985. Phylogenies and the comparative method. American Naturalist 125:1-15.

521    FitzJohn, R., W. Maddison, and S. Otto. 2009. Estimating trait-dependent speciation and extinction

522         rates from incompletely resolved phylogenies. Systematic Biology 58:595–611.

523    Gittleman, J. L. and M. Kot. 1990. Adaptation: Statistics and a null model for estimating

524         phylogenetic effects. Systematic Zoology 39:227-241.

525    Graham, C. H., J. L. Parra, C. Rahbek, and J. A. McGuire. 2009. Phylogenetic structure in tropical

526         hummingbird communities. Proceedings of the National Academy of Sciences of the United

527         States of America 106:19673-19678.

528    Helmus, M. R., T. J. Bland, C. K. Williams, and A. R. Ives. 2007. Phylogenetic measures of

529         biodiversity. American Naturalist 169:E68-E83.

530    Hernandez, F.M. and E. S. Vrba. 2005. A complete estimate of the phylogenetic relationships in

531         Ruminantia: a dated species-level supertree of the extant ruminants. Biological Reviews 80:

532         269–302.

533 Hillis, D. M., T. Heath and K. St. John. 2005. Analysis and Visualization of Tree Space. Systematic

534         Biology 54:1-12.

535 Housworth, E. A. and E. P. Martins. 2001. Random sampling of constrained phylogenies:

536         conducting phylogenetic analyses when the phylogeny is partially known. Systematic

537         Biology 50:628-639.

538 Holder, M. and P. O. Lewis. 2003. Phylogeny Estimation: Traditional and Bayesian Approaches.

539         Nature Reviews Genetics 4:275-284.

540 Huelsenbeck, J. P. and B. Rannala. 2003. Detecting correlation between characters in a comparative

541         analysis with uncertain phylogeny. Evolution 57:1237-1247.

542 Huelsenbeck, J. P., B. Rannala, and J. P. Masly. 2000. Accommodating Phylogenetic Uncertainty in

543         Evolutionary Studies. Science 288:2349-2350.

544 Ives, A. R., P. E. Midford and T. Garland Jr. 2007. Within-species variation and measuremente

545         Error in Phylogenetic Comparative Methods. Systematic Biology 56:252-270.

546 Isaac, N. J. B., S. T. Turvey, B. Collen, C. Waterman, and J. E. M. Baillie. 2007. Mammals on the

547         EDGE: Conservation Priorities Based on Threat and Phylogeny. PLoS ONE 2:e296.

548 Kuhn, T. S., A. O. Mooers, and G. H. Thomas. 2011. A simple polytomy resolver for dated

549         phylogenis. Methods in Ecology and Evolution 2:427-436.

550 Lewis, P. O., M. T. Holder and K. E. Holsinger. 2005. Polytomies and Bayesian Phylogenetic

551         Inference. Systematic Biology 54:241-253.

552 Losos, J. B. 1994. An approach to the Analysis of Comparative Data When a Phylogeny is

553         Unavailable or Incomplete. Systematic Biology 43:117-123.

554 Martins, E. P. 1996. Conducting phylogenetic comparative studies when the phylogeny is not

555         known. Evolution 50:12-22.

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

556 Martins, W.S., W. C. Carmo, H. J. Longo, T. C. Rosa and T. F. Rangel. 2013. SUNPLIN:

557     Simulation with Uncertainty for Phylogenetic Investigations. BMC BioInformatics 14:324-

558     335.

559 McGuire, J. A., C. C. Witt, D. L. Altshuler, and J. V. Remsen Jr. 2007. Phylogenetic systematics

560     and biogeography of hummingbirds: Bayesian and maximum likelihood analyses of

561     partitioned data and selection of an appropriate partitioning strategy. Systematic Biology

562     56:837-856.

563 McGuire, J. A., C. C. Witt, J. Remsen Jr, A. Corl, D. L. Rabosky, D. L. Altshuler, and R. Dudley.

564     2014. Molecular phylogenetics and the diversification of hummingbirds. Current Biology

565     24:910-916.

566 Jetz, W., G. H. Thomas, J. B. Joy, K. Hartmann and A. O. Mooers. 2013. The global diversity of

567     birds in space and time. Nature 491:444-448.

568 Pagel, M. D. 1999. Inferring the historical patterns of biological evolution. Nature 401:877-884.

569 Pannell, D.J. 1997. Sensitivity analysis of normative economic models: Theoretical framework and

570     practical strategies. Agricultural Economics 16: 139-152.

571 Pavoine, S. and C. Ricota. 2012. Testing for phylogenetic signal in biological traits: the ubiquity of

572     cross-product statistics. Evolution, 67:828-840.

573 Rahbek, C. and G. R. Graves. 2000. Detection of macro-ecological patterns in South American

574     hummingbirds is affected by spatial scale. Proceedings of the Royal Society of London, B

575     267:2259-2265.

576 Rahbek, C. and G. R. Graves. 2001. Multiscale assessment of patterns of avian species richness.

577     Proceedings of the National Academy of Sciences of the USA 98:4534-4539.

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

578    Ranwez, V., V. Berry, A. Criscuolo, P. H. Fabre, S. Guillemot, C. Scornavacca and E. J. P.

579        Douzery. 2007. PhySIC: a veto supertree method with desirable properties. Systematic

580        Biology 56:798–817.

581    Revell, L. J., L. J. Harmon, and D. C. Collar. 2008. Phylogenetic signal, evolutionary process, and

582        rate. Systematic Biology 57:591-601.

583    Ronquist, F. 2004. Bayesian Inference of Character Evolution. Trends in Ecology and Evolution

584        19:475-481.

585    Roquet, C., W. Thuiller and S. Lavergne. 2012. Building megaphylogenies for macroecology:

586        taking up the challenge. Ecography 36:13-26.

587    Saltelli, A., M. Ratto, T. Andres, F. Campolongo, J. Cariboni, D. Gatelli, M. Saisana, and S.

588        Tarantola. 2008. Global Sensitivity Analysis: The Primer. John Wiley & Sons.

589    Schuchmann, K.L. 1999. Family Trochilidae. Pages 468-680 in J. del Hoyo, A. Elliott, and J.

590        Sargatal, editors. Handbook of the Birds of the World. Vol. 5. Barn-owls to Hummingbirds.

591        Lynx Edicions, Barcelona.

592    Steel, M. A. and D. Penny. 1993. Distributions of tree comparison metrics - some new results.

593        Systematic Biology 42:126-141.

594    Swofford, D. L. 1991. When are phylogeny estimates from molecular and morphological data

595        incongruent? Pages 295-333 in M. M. Miyamoto and J. Cracraft, editors. Phylogenetic

596        analysis of DNA sequences. Oxford Univ. Press, New York.

597    Xi Z., L. Liu, J.S. Rest and C.S. Davis. 2014. Coalescent versus concatenation methods and the

598        placement of *Amborella* as sister to water lilies. Syst. Biol. (*in press*).

599    Webb, C. O. and M. J. Donoghue. 2005. Phylomatic: tree retrieval for applied phylogenetics.

600        Molecular Ecology Notes 5:181-183.

601                                                        27

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

602  **Figure Legends:**

603

604   Figure 1: A conceptual example of the effect of different sources of uncertainty in phylogenetic

605   trees (left) on the estimation of phylogenetic relationship between species (right, represented as

606   matrices). In each matrix, cells in the lower-left half-matrix represent the existence of phylogenetic

607   information about the relationship between a pair of species, because both species are present in the

608   tree. Conversely, a dash represents the absence of such phylogenetic information, as one or both of

609   the species is missing from the tree. In the upper-right half-matrix a normal distribution represents

610   the variance in the estimated phylogenetic relationship, while zeros represent certainty. (A)

611   Hypothetical, unknown true tree, without missing species or uncertainty in species relationships.

612   (B) Consensus tree with missing species. Relationships are assumed to be known with certainty. (C)

613   Polytomic super-tree. Insertion of missing species in polytomies generates a complete tree, and

614   relationships are assumed to be known with certainty, although the tree differs from the true tree

615   (A). (D) Replication in the use of phylogenetic trees incorporates the uncertainty in the relationship

616   between species, but missing species are ignored. (E) Missing species are inserted in multiple

617   phylogenies, accounting for uncertainty in phylogenetic reconstruction and lack of a complete

618   phylogeny.

619

620  Figure 2: Schematic representation of a workflow using the analytical strategy, proposed here, to

621   account for phylogenetic uncertainty. (A) The true, but unknown, tree for a clade of 8 taxa. (B)

622   Molecular data are available for only six taxa; for taxa B and G no molecular data are available. (C)

623   Two experts, based on their knowledge (e.g. taxonomy, behavior, morphology, geographic

624   distribution, etc.), suggest possible sister taxa of the two taxa that lack molecular data.  (D) Using

625   the available molecular data and phylogenetic reconstruction methods, three backbone phylogenies

626   are proposed. The variation between backbone phylogenies arises from uncertainty in the process of

28

627     phylogenetic reconstruction. Among the backbone phylogenies, the taxa B and G are considered

628     *phylogenetically uncertain taxa* (PUT), because no molecular data are available for them, and

629     therefore they are missing from the backbone phylogenies. Using the information provided by the

630     experts the *most derived consensus clade* (MDCC) is found for each PUT. For PUT B, the MDCC

631     is the clade that necessarily includes taxa C, D, and A (indicated by a bold B in each backbone

632     phylogeny), whereas for PUT G the MDCC must include taxa F and H (indicated by a bold G in

633     each backbone phylogeny). (E) Statistical analyses that use phylogenetic trees as input data should

634     be replicated using samples of operational trees. In each of the operational trees, the PUTs B and G

635     were randomly inserted within their respective MDCCs. The insertion of each PUT was replicated

636     three times (columns) for each backbone tree. The variation among operational trees that use the

637     same backbone tree (each column of the operational trees) arises from variation in placement of the

638     missing taxa (as the backbone tree is not changed by the randomization process), and is caused by

639     uncertainty in the phylogenetic relationships of taxa B and G (both PUTs).

640

641     Figure 3: Hummingbird phylogenies. (A) One backbone phylogeny from McGuire et al. (2007),

642     with red internal branches indicating the *Most Derived Consensus Clades* (MDCCs) used in this

643     analysis. (B) Red taxa (polytomies) indicate PUTs inserted in the original phylogeny at the base of

644     their MDCCs. (C) Red taxa indicate PUTs inserted to yield a fully-resolved phylogeny,

645     randomizing their position within the respective MDCC.

646

647     Figure 4: Moran's I correlogram (top panels) for two hummingbird traits (body size and maximum

648     elevational range limit). The inner 95% confidence interval (C.I.) around each estimate indicates

649     variance due to sampling error, the intermediate 95% C.I. the variance due to missing species, and

650     the outer 95% C.I. the variance due to phylogenetic reconstruction. The bottom panels indicate the

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

651    relative proportion of statistical error caused by each source of uncertainty across evolutionary time:

652    sampling (green), missing species (PUTs, red) and phylogenetic reconstruction (blue).

653

654    Figure 5: Hummingbird assemblages with PSV values significantly different from null expectation.

655    Blue cells indicate significant phylogenetic clustering, whereas red cells indicate significant

656    phylogenetic dispersion. (A) Standard PSV analysis, considering only non-random "sampling" from

657    the phylogeny. (B) Re-analysis of PSV accounting for three sources of error: sampling, missing

658    species, and phylogenetic reconstruction. Notice absence of significant phylogenetic dispersion and

659    decreased areas of phylogenetic clustering when all sources of uncertainty are accounted for.

660

661    Figure 6: Relative statistical error associated with PSV analysis for hummingbird species present in

662    each map cell, partitioned among sources of uncertainty. Each point within the cube has a unique

663    color that represents the relative proportions of uncertainty due to phylogeny, PUTs and sampling

664    error, as shown by the axes. In the map, the prevalence of red cells indicates sampling error as the

665    main source of uncertainty, whereas blue indicates substantial error caused by phylogenetic

666    uncertainty. Notice the absence of green and yellow areas, indicating that phylogenetic uncertainty

667    due to missing species (PUTs) is relatively irrelevant in the phylogenetic structure of assemblages.

For: Evolution

Correspondence to:
Thiago Fernando Rangel
Departamento de Ecologia
Universidade Federal de Goiás
Cx.P: 131, 74001-970
Goiânia, Goiás, Brasil
thiago.rangel@ufg.br

Article:

# Phylogenetic Uncertainty Revisited:

# Implications for Ecological Analyses

Thiago F. Rangel [1]
Robert K. Colwell [1,2,3]
Gary R. Graves [4,5]
Karolina Fučíková [2]
Carsten Rahbek [4]
José Alexandre F. Diniz-Filho [1]

1: Departmento de Ecologia, Universidade Federal de Goiás, CP 131, 74.001-970, Goiânia, GO, Brasil. TFR: thiago.rangel@ufg.br. JAFD-F: diniz@ufg.br.

2: Department of Ecology and Evolutionary Biology, University of Connecticut, 75 North Eagleville Rd., Unit 3043, 06269-3043, Storrs, CT, United States. RKC: robert.colwell@uconn.edu. KF: karolina.fucikova@uconn.edu.

3. University of Colorado Museum of Natural History, Boulder, Colorado 80309, USA

4: Department of Vertebrate Zoology, National Museum of Natural History, Smithsonian Institution, Washington, DC 20013, United States. gravesg@si.edu

5: Center for Macroecology, Evolution and Climate, Department of Biology, University of Copenhagen, Universitetsparken 15, 2100, Copenhagen O, Denmark. crahbek@bio.ku.dk

Counts: 5,953 Words, 0 Tables, 6 Figures

1

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

## 1    **Abstract**

2     Ecologists and biogeographers usually rely on a single phylogenetic tree to study evolutionary

3     processes that affect macroecological patterns. This approach ignores the fact that each

4     phylogenetic tree is a hypothesis about the evolutionary history of a clade, and cannot be directly

5     observed in nature. Also, trees often leave out many extant species, or include missing species as

6     polytomies because of a lack of information on the relationship among taxa. Still, researchers

7     usually do not quantify the effects of phylogenetic uncertainty in ecological analyses. We propose

8     here a novel analytical strategy to maximizes the use of incomplete phylogenetic information, while

9     simultaneously accounting for several sources of phylogenetic uncertainty that may distort

10    statistical inferences about evolutionary processes. We illustrate the approach using a clade-wide

11    analysis of the hummingbirds, evaluating how different sources of uncertainty affect several

12    phylogenetic comparative analyses of trait evolution and biogeographic patterns.  Although no

13    statistical approximation can fully substitute for a complete and robust phylogeny, the method we

14    describe and illustrate enables researchers to broaden the number of clades for which studies

15    informed by evolutionary relationships are possible, while allowing the estimation and control of

16    statistical error that arises from phylogenetic uncertainty. Software tools to carry out the necessary

17    computations are offered.

## Introduction

18

19    Ecological and biogeographical studies often involve hundreds of species, representing both recent

20    and ancient lineages, and both locally endemic and cosmopolitan species. Although the geographic

21    distribution of most groups of terrestrial vertebrates is increasingly well known, species sampling

22    for phylogenetic analysis is rarely complete for larger clades. Even when a comprehensive

23    phylogeny is available for ecological analyses, key sources of phylogenetic uncertainty are seldom

24    taken into account. We begin this paper by outlining the sources of uncertainty in studies that rely

25    on phylogenetic hypotheses to infer evolutionary processes, and how these sources have been

26    considered or ignored in previous studies. With this motivation, we propose here a novel analytical

27    strategy to quantify and account for key sources of phylogenetic uncertainty in any study that uses

28    phylogenetic input data. As an example of the application of our methodology, we implement it to

29    investigate patterns of trait evolution and phylogenetic assemblages in hummingbirds. We contend

30    that any study that infers evolutionary hypotheses should aim to account for all sources of

31    phylogenetic uncertainty. Finally, we offer freely available software tools to help researchers

32    account for phylogenetic uncertainty in their own research.

33    Phylogenetic uncertainty may arise from two distinct sources (FitzJohn et al. 2009, Diniz-

34    Filho et al. 2013): (1) weak, missing, or conflicting empirical support for hypothesized relationships

35    among species in a given clade (e.g. tree topology, branch length estimation and absolute time

36    calibration), which can be expressed in the form of multiple alternative topologies (e.g., resulting

37    from different analysis methods or different data types), polytomic clades, or low branch support

38    values; and (2) incomplete and unrepresentative sampling of known species. Most ecological

39    studies implicitly assume absolute knowledge of phylogenetic history by simply ignoring the

40    uncertainty caused by the lack of empirical support for phylogenetic trees or issues related to branch

41    length estimation (Fig. 1). When contrasting phylogenetic hypotheses are available, ecologists

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

42    usually rely on a consensus tree to "average" phylogenetic information, thus failing to take account

43    of variation among trees (an expression of overall phylogenetic uncertainty). In addition, although

44    ecologists acknowledge that most polytomies are expressions of ambiguous or missing empirical

45    data, highly polytomic trees are nonetheless often used in ecological analyses without proper

46    quantification of phylogenetic uncertainty and without conducting the necessary sensitivity

47    analyses.

48    Ecologists have dealt with incomplete phylogenies in several ways. The first approach is to

49    focus only on clades for which relatively complete phylogenies are available, but this strategy

50    restricts ecological studies to a very small number of groups (Pagel 1999) and can undermine

51    assemblage-level or macroecological studies (Webb and Donoghue 2005). A favored method is to

52    assemble supertrees from smaller, overlapping trees and to fill gaps in phylogenies by placing

53    unsampled species in trees according to their taxonomic classification (Bininda-Emonds 2004,

54    Webb and Donoghue 2005, Hernandez and Vrba 2005, Davies et al. 2004, Ranwez et al. 2007).

55    Nevertheless, even for the best-studied taxa, such as mammals and birds, supertrees that span all

56    species are relatively recent achievements (e.g. Bininda-Emonds et al. 1999, Beck et al. 2006).

57    Unfortunately, supertrees are usually highly polytomic (not fully resolved). Recent approaches

58    based on building megatrees (Roquet et al. 2012) or complex addition of species on backbone trees

59    (Jetz et al. 2012) are still in their infancy. Phylomatic software (Webb and Donoghue 2005), for

60    instance, is a popular tool for assemblage-level ecological studies that constructs customized trees

61    of virtually any size. Although Phylomatic allows additional phylogenetic representations of the

62    clade when compared to the backbone phylogeny, terminal branches are usually based on

63    taxonomic relationships among inserted species and thus are commonly polytomic. Additional

64    drawbacks of supertrees are the absence of information on branch lengths (which require further

65    effort, and assume confidence in the fossil record for time calibration) and, again, the lack of

66    explicit conversion of phylogenetic uncertainty into explicit measures of error associated with

67    statistical inference and parameter estimates.

68        A second and more radical strategy consists of ignoring the species that are absent from the

69    available phylogeny, under the assumption that the species included in the analysis represent an

70    unbiased and representative sample of all species in the clade (FitzJohn et al. 2009). Of course, the

71    full evolutionary history of a clade of substantial age can be described only by a phylogeny with all

72    species (including extinct ones), although this ideal is achievable only in simulated scenarios

73    (Colwell and Rangel 2010). Estimating the degree of bias due to missing species can, in principle,

74    be achieved by replicating the analysis with random sub-samples of species that are present in the

75    phylogeny (rarefaction or "thinning", Davies et al. 2012), but unfortunately this approach is not

76    common in evolutionary or ecological studies.

77        Several other approaches to deal with phylogenetic uncertainty have been recently

78    developed, increasing awareness that species missing from phylogenetic trees may seriously affect

79    statistical inference of evolutionary processes, even when the missing species are inserted in the tree

80    in the form of polytomies (Davies et al. 2012, Diniz-Filho et al. 2013). For example, under the

81    Bayesian framework of molecular phylogenetic reconstruction, missing species can be inserted by

82    assigning empty sequences to missing species and by constraining their insertion using priors on the

83    tree topology (Huelsenbeck and Rannala 2003). Several recent studies have employed different

84    analytical strategies to account for phylogenetic uncertainty. For example, Isaac et al. (2007) dealt

85    with missing species in their analysis of evolutionary distinctiveness of mammals by allocating

86    missing species among their presumed closest relatives using a model of constant rate of speciation

87    and extinction. Day et al. (2008) studied the diversification rates of cichlid fish radiation of Lake

88    Tanganyika, accounting for the potential effects of missing species on the estimates of the timing

89    and rate of diversification. Kuhn et al. (2011) proposed a method and a computational tool that uses

90    a birth-death model of diversification to resolve polytomies and to define not only tree topology,

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

91    but also branch lengths in polytomic trees, and Jetz et al. (2012) built sets of phylogenetic trees for

92    the entire avian clade, combining genetic and taxonomic information and adding a simulation of

93    evolutionary diversification to better establish branch lengths and resolve polytomies. Davies et al.

94    (2012) showed that phylogenetic uncertainty can dramatically inflate estimates of phylogenetic

95    signal and proposed a rarefaction-based method to guide inference. After accounting for

96    phylogenetic uncertainty, Batista et al. (2013) showed that potential loss of evolutionary history

97    caused by extinction of threatened Western Hemisphere anurans would not be significantly higher

98    than that of non-threated anurans. These few examples clearly show the importance of dealing with

99    phylogenetic uncertainty, and demonstrate that systematists and comparative biologists have

100   already begun to develop and employ methods to incorporate estimates of uncertainty in the

101   inference of diversification rates (Day et al. 2008, FitzJohn et al. 2009), character evolution (Losos

102   1994, Huelsenbeck et al. 2000, Housworth and Martins 2001, Huelsenbeck and Rannala 2003, Ives

103   et al. 2007), reconstruction of ancestral states (Ronquist 2004), and tree topology (Felsenstein 1985,

104   Holder and Lewis 2003).

105        Here we propose a unified analytical approach to estimate and account for multiple sources

106   of phylogenetic uncertainty. Our method consists of partitioning variance among estimated

107   parameters of evolutionary processes, using random sampling from the universe of probable

108   phylogenies. We illustrate the approach using a clade-wide analysis of the hummingbirds,

109   evaluating how different sources of uncertainty affect several phylogenetic comparative analyses of

110   trait evolution and biogeographic patterns.

111

## Accounting for multiple sources of phylogenetic uncertainty

113   We propose here an empirical approach, based on simulations, that maximizes the use of

114   incomplete phylogenetic information in ecological studies, while simultaneously accounting for

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

115    several sources of phylogenetic uncertainty that may distort statistical inferences about evolutionary

116    and ecological processes.

117

118

119    *Missing species*

120    Here we define a *phylogenetically uncertain taxon* (henceforth a "PUT" for a single taxon, and

121    "PUTs" for multiple taxa) as a taxonomic unit (e.g. population, species, genus) that is recognized as

122    valid and is accepted as belonging to a particular clade, but is missing from the available

123    phylogenetic tree(s) for that clade. The absence from the tree could be due to multiple causes, such

124    as unavailability of molecular and morphological data (Fig. 2b). However, although the data

125    required to formally reconstruct the phylogenetic relationships of a PUT may be missing,

126    evolutionary information about the PUT is rarely completely non-existent, as no authority would

127    dispute that some other taxon in the clade must be the closest relative of the PUT at some level in

128    the phylogeny. In fact, additional available information, such as taxonomic, morphological or

129    behavioral data, may be useful to define a list of the most likely sister species or lineages of any

130    given PUT (Fig. 2c). Therefore, it makes sense to use all acceptable information available to

131    conservatively define what we will call the *most derived consensus clade* (MDCC), i.e., the node

132    that unequivocally contains each PUT.

133        Thus, an MDCC can be defined as the most recent common ancestor of all unquestioned

134    candidates for being the closest relative of the PUT (Fig. 2d). Of course, the validity of information

135    (e.g. taxonomic, biogeographical, behavioral, morphological, etc) used to define a MDCC can be

136    questioned. However, if two taxonomists, for example, disagree on the proper MDCC for a given

137    PUT, instead of eliminating the PUT from the statistical analysis, thereby ignoring phylogenetic

138    uncertainty, the MDCC should be conservatively redefined as the most recent common ancestor of

7

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

139    the two candidate MDCCs for the PUT, each championed by a different taxonomist (Fig. 2). Thus,

140    it is important to stress that the goal of defining a MDCC to a PUT is not to replace the formal

141    methods of phylogenetic reconstruction, but to conservatively use all available information to allow

142    biologists to make inference of ecological and evolutionary processes using all available

143    information, even with incomplete phylogenies.

144         We propose here a simple process of building a possible phylogeny that includes all extant

145    species. This process relies on the definition of an MDCC for each PUT and the insertion of the

146    PUT in a random position within its MDCC. We begin by randomizing the order in which PUTs are

147    to be added to the tree. Next, within each MDCC, we assign the PUT to a point along one a branch

148    of the clade. The choice of insertion point along a branch can be uniform random if no additional

149    information is available, or it may be guided by a biologically realistic model of diversification

150    (Kuhn et al. 2011, Davies et al. 2012, Jetz et al. 2012). If the baseline tree is ultrametric, the branch

151    length for the inserted PUT is simply the distance from the attachment point to the end of all other

152    tips. If, however, the baseline tree is not ultrametric, the branch length of the PUT can be sampled

153    from a distribution of possible branch length values. Once a PUT has been inserted, its own branch

154    may serve as a potential insertion point for subsequent species assigned from the PUT queue. Our

155    algorithm iterates until each PUT has been sequentially added to the appropriate MDCC, producing

156    a complete phylogeny of known species (Martins et al. 2013, Batista et al. 2013).

157

158    *Polytomies*

159    If all polytomies in a tree are the consequence of phylogenetic uncertainty (i.e. "soft" polytomies),

160    then it is necessary to ensure that such uncertainty is quantified and accounted for in statistical

161    analyses that use such a tree (Lewis et al. 2005). To explore the space of all possible dichotomous

162    trees one must resolve the polytomies, assuming no additional information about the evolutionary

163    history of the clade (Batista et al. 2013), or employing a model of diversification (Martins 1996,

164    Housworth and Martins 2001, Kuhn et al. 2011).

165         Our approach consists in producing fully dichotomous trees by resolving polytomies

166    stochastically. For each node in the tree with three or more branches, we choose two branches at

167    random and reassign them to the same node. Each remaining branch from the original polytomy is

168    then inserted sequentially, in random order, within the clade constructed from the former polytomy,

169    at a randomly chosen position along the length of the existing branches. Dichotomous nodes in the

170    original phylogeny are not changed, thereby preserving the phylogenetic information in the original

171    baseline tree. This process guarantees that the resulting tree is fully resolved and that the

172    phylogenetic uncertainty arising from polytomies is also taken into account when multiple

173    randomized trees are compared (Batista et al. 2013). As with sampling multiple phylogenetic trees

174    (below) or randomizing the position of the PUTs in the phylogeny, not every possible tree topology

175    is examined. However, given a large enough number of replicates, our stochastic procedure ensures

176    that the parameter space will be explored, thereby accounting for phylogenetic uncertainty.

177

178    *Multiple phylogenetic trees*

179    Modern methods of phylogenetic reconstruction and inference are based on searching the universe

180    of possible phylogenetic trees. The search is guided by empirical data among possible trees, which

181    are judged by their ability to predict the observed data, given a model of molecular evolution

182    (Holder and Lewis 2003). However, because of scarce, ambiguous or missing empirical data,

183    searches are often unable to unambiguously rank all trees within an entire group of equally likely

184    trees, therefore yielding a large set of possible trees. Moreover, especially analyses of large,

185    genome-scale data sets may yield multiple yet well-supported topologies depending on the

186    analytical approach or on selected subsets of data (e.g., Xi et al. 2014, examples in Cooper 2014).

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

187    To account for the phylogenetic uncertainty that arises from lack of definitive empirical support for

188    a single phylogenetic hypothesis, one must not limit the statistical analysis to a single sampled tree

189    or to a majority-rule consensus tree that "averages" a larger set of possible trees (Fig. 1). Sampling

190    or averaging among possible trees can potentially improve a point-estimate of a parameter, but does

191    not account for statistical error in evolutionary inference, thereby masking phylogenetic uncertainty

192    associated with tree reconstruction. In practice, variation in the results of ecological analyses that

193    use different phylogenetic trees is an expression of uncertainty in phylogenetic reconstruction.

194    Thus, it is necessary to replicate these analyses over a large number of possible phylogenetic trees,

195    and subsequently study the variation in parameter estimates that arises as a consequence of

196    differences among trees (Fig. 2).

197

198    *Uncertainty and sensitivity analysis*

199    Computer simulations have long been used to estimate the magnitude of statistical error and to

200    quantify the probability of a given outcome when a system is not fully known (Doubilet et al. 1998,

201    Saltelli et al. 2008). Usually run in tandem with uncertainty quantification, a sensitivity analysis

202    allows one to determine the degree to which the sources of uncertainty in model inputs are

203    responsible for the uncertainty in the model output (Pannell 1997). Coupling uncertainty and

204    sensitivity analysis offers us the opportunity to quantify the robustness of models in the presence of

205    phylogenetic uncertainty, account for uncertainty in evolutionary inference, and enhance

206    communication of the magnitude of statistical error in the study.

207         Martins (1996) proposed a way to carry out phylogenetic comparative studies when the

208    phylogenetic relationships among species are unknown. In her method, a large sample of trees is

209    generated by randomly resolving a fully polytomic tree ("star phylogeny"), using models of

210    phenotypic evolution and diversification rates (so that uncertainty in branch lengths is also

10

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

211    implicitly incorporated). Species traits are then analyzed on each of the possible trees. Although the

212    mean results of such analyses will converge to a non-phylogenetic analysis (Abouheif 1998), the

213    approach of Martins (1996) represents a landmark in phylogenetic inference under uncertainty

214    present in an unresolved phylogeny (see also Housworth and Martins 2001).  The mean of the

215    squared standard error of the calculated evolutionary statistic (e.g. the correlation between two

216    species' traits), known as $V_s$, estimates the true variance of the statistic, which is due to sample

217    variance. The variance of the evolutionary statistic calculated among the randomly generated trees,

218    known as $V_p$, estimates the variance due to phylogenetic uncertainty. Thus, inferences based on

219    parameter estimation must account for both sources of error, which can be done by computing

220    confidence intervals or $P$-values based on the sum of $V_p$ and $V_s$.

221        We expand Martins' (1996) strategy for incorporating phylogenetic uncertainty in parameter

222    estimation further to accommodate multiple sources of phylogenetic uncertainty. We treat each

223    source of uncertainty as a factor in an experimental design, while parameter estimates are treated as

224    the response variable under study. Thus, we can ask how different phylogenetic trees, alternative

225    resolutions of polytomies, and/or probable configurations of PUTs would affect parameter

226    estimates. We partition the amount of variance in a parameter estimate arising from each source of

227    phylogenetic uncertainty with an Analysis of Variance (ANOVA), which not only isolates variance

228    among sources but also calculate the magnitude of standard error in the parameter estimate (a

229    separate Martins' (1996) $V_p$ for each source of uncertainty).

230        We view the sources of phylogenetic uncertainty we have outlined here as a hierarchy to be

231    employed in the design of an experiment used to assess uncertainty and sensitivity. If multiple

232    empirical phylogenies are available, then this source of uncertainty can be regarded as the highest

233    level of uncertainty in the analysis. Because each empirical phylogeny is unique, its polytomies

234    must be treated individually. Thus, polytomies can be regarded as the second level of phylogenetic

235    uncertainty, immediately below the multiple empirical phylogenies. Of course, for each available

11

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

236  empirical phylogeny there is a virtually unlimited number of alternative ways to resolve polytomies

237  randomly or by following some diversification model. Finally, for each combination of original

238  phylogeny and polytomy resolution, PUTs can be inserted according to the procedure described

239  above. Thus, insertion of PUTs is the lowest level in the hierarchy of sources of phylogenetic

240  uncertainty. Because sources of uncertainty are nested within a higher level of classification, and

241  groups representing subordinate levels are randomly chosen, a nested (hierarchic) analysis of

242  variance is ideal for the analysis of sensitivity and uncertainty.

243

244  **Application example: trait evolution and phylogenetic assemblage patterns in**

245  **hummingbirds**

246  We illustrate our method by applying it to ecological hypotheses for the hummingbirds

247  (Trochilidae), a large, monophyletic clade (~330 species) with a rich natural history literature and

248  ongoing phylogenetic, morphological, behavioral, and ecological research. McGuire et al. (2007)

249  published a multilocus molecular phylogeny for the hummingbirds that includes only 146 species,

250  but encompasses all the higher-level trochilid diversity (73 of the approximately 104 recognized

251  genera, Schuchmann 1999). Although a more complete hummingbird phylogeny is now available

252  (McGuire et al. 2014), we use the earlier phylogeny(McGuire et al. 2007) to demonstrate our

253  methods because it is typical of current level of phylogenetic knowledge for many comparably

254  diverse taxa.

255      Replicating the phylogenetic reconstruction methodology of McGuire et al. (2007), we

256  sampled 25,000 trees from the posterior distribution (hereafter called "backbone" trees), after

257  discarding a burn-in of 5,000 steps. We relied upon the taxonomic classification of Schuchmann

258  (1999) to designate the most likely MDCC for each PUT (Fig. 3).

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

259       From the total of about 330 described species of hummingbirds, we compiled geographical

260    distributions and ecologically important morphological data for 304 species. Species endemic to

261    islands (the West Indies and the Juan Fernández Archipelago) were not included in the analysis. We

262    compiled estimates of average body mass (intersexual mean) for each of the 304 species from

263    Dunning (2007) and Schuchmann (1999). Although large intersexual and geographical variation in

264    hummingbird body size is well documented (Colwell 2000), the averages used here are useful for

265    broad taxonomic analyses. Based on the literature, we also recorded the maximum known

266    elevational range limits for each species (see Appendix A for references). Body mass (as a measure

267    of body size) and maximum elevational range limit were log-transformed to conform to a normal

268    distribution.

269       We extracted distributional data for 304 species from an updated version (16 July 2010) of

270    the comprehensive database for all land and fresh-water birds known to have breeding populations

271    in the Western Hemisphere, complied by Rahbek and Graves (2000, 2001), and mapped the

272    geographical range of each species on a gridded map at a resolution of 1° x 1° (latitude-longitude).

273    These maps represent a conservative extent-of-occurrence estimate of the breeding range based on

274    museum specimens, published sight records, and spatial distribution of habitats based on

275    documented records for South America.

276

277    *Uncertainty quantification*

278    Variability in input data is a fundamental source of uncertainty. To estimate the degree of

279    phylogenetic uncertainty in the hummingbird data due to phylogeny reconstruction and missing

280    species we randomly sampled 100 trees from the posterior (here after called "backbone" trees), and

281    for each tree we replicated the insertion of PUTs 100 times, thereby generating a total of 10,000

282    fully resolved phylogenetic trees. Next, we calculated the tree-to-tree pairwise distance matrix using

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

283     the weighted Robinson-Foulds (wRF) metric, which measures the minimum number of internal

284     branches that must be collapsed or expanded to make two trees identical. The weighted version of

285     Robinson-Foulds distance accounts for differences in branch lengths between trees (Steel and

286     Penny 1993). Distance matrices of phylogenetic trees have been used in phylogenetic reconstruction

287     to calculate consensus trees (Swofford 1991), produce supertrees (Bansal et al. 2010), and visualize

288     multi-dimensional space of possible trees (Hillis et al. 2005).

289          The average wRF pair-wise distance between the 10,000 trees is 2.7423, with a standard

290     deviation of 0.136. We used a multivariate analysis of variance (PERMANOVA, Anderson 2001)

291     to partition variance among phylogenetic trees through the pair-wise wRF distance matrix. The

292     estimated components of variance indicated that 17.02% of total variance among trees is

293     attributable to phylogenetic reconstruction (i.e. differences among backbone trees), whereas 82.97%

294     of total variance is due to phylogenetic uncertainty of missing species (PUTs).

295

296     *Phylogenetic autocorrelation in species traits*

297     We used Moran's I correlograms to estimate the magnitude of phylogenetic autocorrelation in body

298     size and elevational range limits among hummingbird species (see Gittleman and Kot 1990; Diniz-

299     Filho 2001; Pavoine and Ricota 2012). Moran's I is larger when species within a given phylogenetic

300     distance interval have similar trait values and smaller when species within the same distance

301     interval have very different trait values. To account for uncertainty in phylogenetic reconstruction

302     and missing species, the calculation of correlograms was repeated 1,000,000 times (i.e., replicating

303     the insertion of PUTs 1,000 times in each of the 1,000 randomly sampled backbone phylogenies).

304     We used a nested ANOVA to partition sources of uncertainty for each of the 12 distance classes in

305     the correlogram. Applying our method, the average Moran's I within a distance class is the

14

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

306    parameter estimate, whereas additive confidence intervals are calculated around the estimate for

307    each source of uncertainty.

308         Results of these analyses showed that phylogenetic autocorrelation in body size and

309    maximum elevational range limit both decline with phylogenetic distance, although at different

310    rates (Fig. 4, top panel in each plot). Closely related species tend to have similar body sizes but less

311    similar elevational range limits. Accounting for phylogenetic uncertainty widens confidence

312    intervals by at least a factor of two, regardless of trait (Fig. 4, confidence intervals in upper panel;

313    see the caption for details). For all distance classes, standard confidence intervals due to sampling

314    error are always narrower than confidence intervals due to phylogenetic uncertainty. In addition, the

315    magnitude of error attributed to each source of uncertainty is not constant over phylogenetic

316    distance (Fig. 4, bottom panels). The relative proportion of error caused by missing species tends to

317    be higher in short distance classes, as taxonomic information has been used to assign PUTs to

318    relatively derived positions in the phylogeny (MDCCs). In contrast, large phylogenetic confidence

319    intervals at intermediate and long distance classes arise from uncertainty in the empirical phylogeny

320    located at the base of the coquettes and brilliants clades, which together form the Andean clade

321    (McGuire et al. 2007).

322

323    *Phylogenetic signal and evolutionary models in species traits*

324    We used Blomberg et al.'s (2003) *K* statistic (hereafter *K*) to estimate the magnitude of deviation

325    from Brownian motion in body size and elevational range limits. *K* measures the degree of

326    similarity in species' traits in relation to the similarity expected under a Brownian motion model of

327    phenotypic evolution, given a phylogenetic hypothesis. *K* is based on the ratio between two

328    measures of distance: $MSE_0$, the squared distance between trait values and the phylogenetically

329    corrected mean trait value, and MSE, the squared distance between trait values estimated from a

15

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

330    variance-covariance matrix derived from a phylogenetic hypothesis. Thus, large values of

331    $MSE_0/MSE$ indicate a strong phylogenetic signal. To allow comparisons between different traits

332    and trees, observed $MSE_0/MSE$ is standardized by expected $MSE_0/MSE$, assuming that the trait

333    evolves under a Brownian motion model of phenotypic evolution. $K$ values less than one indicate

334    that species are less similar for a given trait than expected under Brownian motion evolution,

335    whereas $K$ values greater than one indicate that species are more similar for a given trait than

336    expected under Brownian motion evolution.

337         To account for uncertainty in phylogenetic reconstruction and missing species, we again

338    calculated 10,000 possible $K$ values, replicating the insertion of PUTs 100 times in each of the 100

339    randomly sampled phylogenies. To estimate the accuracy (standard error) of $K$ within each tree we

340    used jackknife permutation for each combination of PUT insertion and sampled phylogeny (Efron

341    and Tibshirani 1994). Finally, we used a nested ANOVA to partition error in estimated $K$ among

342    different sources of uncertainty.

343         Our results revealed a surprisingly low phylogenetic signal in body size among

344    hummingbirds ($\overline{K}$ = 0.0597); thus body size evolution can hardly be explained by a simple model

345    of Brownian evolution. However, the standard error of $K$ due to sample size ($s_{KSamp}$=0.0011) is

346    relatively small (2.23%) compared to the error arising from phylogenetic uncertainty due to missing

347    species ($s_{KSamp+PUT}$=0.0494), as this source of uncertainty represents 97.1% of the total error in $K$.

348    Finally, uncertainty due to phylogenetic reconstruction represents only 0.67% of total error of the

349    evolutionary inference ($s_{KSamp+PUT+Phy}$=0.0497).

350         The relative importance of sources of uncertainty shift discordantly when maximum

351    elevational range limit is considered, although this trait also has very little phylogenetic signal ($\overline{K}$ =

352    0.0221). Standard error of $K$ due to sampling size contributes 59.7% of total error ($s_{KSamp}$ = 0.0179),

16

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

353 whereas error due to missing species represents 40.1% of total error ($s_{KSamp+PUT}$= 0.0300), and error

354 due to phylogenetic reconstruction represents only 0.2% of total error ($s_{KSamp+PUT+Phy}$=0.03001).

355 To test the hypothesis that estimated *K* values are not significantly different from random

356 expectation, we calculated 100 *Ks* for each combination of PUT insertion and sampled phylogeny,

357 randomizing trait values among species (Blomberg et al. 2003, Revell et al. 2008). We employed

358 the same nested ANOVA design to estimate standard error of the null distribution due to each

359 source of uncertainty. Finally, we used Welch's *t*-test to evaluate the hypothesis that estimated *Ks*

360 do not differ significantly from the null expectation:

$$t = \frac{\overline{K}_{trait} - \overline{K}_{H0}}{\sqrt{\frac{s^2_{K_{trait}}}{n_{K_{trait}}} + \frac{s^2_{K_{H0}}}{n_{K_{H0}}}}}$$

361

362 Estimated *K* for hummingbird body mass is significantly larger than expected under the null

363 expectation ($\overline{K}_{H0}$= 0.0237), and accounting for phylogenetic uncertainty does not affect the

364 hypothesis test ($t_{Samp}$= 2589.92, $t_{Samp+PUT}$ = 62.5, $t_{Samp+PUT+Phy}$ = 76.72, all *P*-values < 0.001).

365 Conversely, phylogenetic signal in maximum elevational range limit is significantly smaller than

366 expected under the null hypothesis ($\overline{K}_{H0}$= 0.0333). Accounting for phylogenetic uncertainty also

367 does not change the inference of statistical significance in hummingbird maximum elevational

368 range limit ($t_{Samp}$ = -275.36, $t_{Samp+PUT}$ = -34.74, $t_{Samp+PUT+Phy}$ = -44.47, all *P*-values < 0.001).

369 These results indicate that, although body size evolution deviates significantly from

370 Brownian motion, it carries a small phylogenetic signal. Thus, evolutionary constraints and niche

371 conservatism can be invoked for body size (for instance, evolution under an O-U process – see

372 Hansen et al. 2008), whereas elevational range has less phylogenetic signal than expected by

373 chance.

17

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

374

375    *Phylogenetic community structure at the macroecological scale*

376    Phylogenetic species variability (PSV, Helmus et al. 2007) of an assemblage is maximized (PSV =

377    1) when the assemblage is composed of the least related species in a clade, and minimized (PSV =

378    0) when the most related species coexist in an assemblage. We calculated PSV for the 1979

379    assemblages in the Western Hemisphere based on mapped species ranges with at least two

380    hummingbird species, and tested statistical significance of each PSV value using 300 permutations

381    of species identities. However, because phylogenetic distance between hummingbird species is not

382    known with certainty, we generated 300 possible phylogenies with different PUT insertions, and

383    replicated the calculation of PSV for each of the 1979 assemblages using the each of the 300

384    randomly selected phylogenetic trees.

385         Some hummingbird assemblages are composed of species with non-random phylogenetic

386    relationships (Fig. 5). However, identifying any significant departure from randomness requires

387    accounting for all sources of uncertainty. Standard null models to test phylogenetic structure of

388    communities (e.g. Graham et al. 2009), which randomize species identities while preserving species

389    richness, are designed to account only for non-random "sampling" from the phylogeny. For

390    hummingbird assemblages across the Western Hemisphere, an analysis with such a null model

391    would lead to the inference of significant phylogenetic dispersion along the middle and upper

392    Andes (contrary to Graham et al. 2009), and significant phylogenetic clustering across the Pacific

393    coast of Central and North America (Fig 5A). However, uncertainty caused by missing species and

394    phylogenetic reconstruction may seriously affect pattern detectability. Because of substantial

395    uncertainty in the phylogenetic relationship of the two Andean clades (Coquettes and Brilliants),

396    when phylogenetic uncertainty is taken into account no significant phylogenetic dispersion in

397    Andean assemblages is detected. Notice that this result is concordant with the lack of phylogenetic

18

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

398    signal and lack of phylogenetic autocorrelation for elevational range at the species level, as we

399    previously discussed. Moreover, once phylogenetic uncertainly is taken into account, many

400    assemblages in Central America can no longer be considered phylogenetically clustered (Fig. 5B).

401         Because species are neither randomly distributed in the phylogeny nor in geographic space,

402    the magnitude of error in the analysis of phylogenetic structure of assemblages tends to be strongly

403    spatially autocorrelated. In addition, sampling effect is usually higher in species-poor assemblages.

404    On the other hand, assemblages with a higher proportion of species that are members of clades

405    characterized by phylogenetic uncertainty in deep (basal) nodes are subject to a higher proportion of

406    error due to phylogenetic reconstruction. Figure 6 depicts the relative contribution of each source of

407    error in analysis the analysis of phylogenetic structure of hummingbird assemblages. Because of

408    phylogenetic uncertainty in the reconstruction of the relationship between the two, basal, Andean

409    clades (Coquetes and Brilliants), statistical error is substantially higher than the other two sources of

410    error for Andean assemblages. Conversely, statistical error in North American assemblages is

411    mostly due to sampling, as species richness is very low. Uncertainty due to missing species (PUTs)

412    is not particularly pronounced in any assemblage, as no assemblage is composed by more than 17%

413    PUTs.

414

415    **Concluding Remarks**

416    Phylogenetic uncertainty is not evenly distributed across time and space. As a consequence, the

417    effects of phylogenetic uncertainty in the statistical analysis of phylogenetic data cannot be

418    estimated without a thorough sensitivity analysis on a case-by-case basis. To illustrate the new

419    methods we propose to account for phylogenetic uncertainty, in this paper we used phylogenetic

420    hypotheses derived from molecular data to analyze the hummingbird clade, a well-studied

421    taxonomic group widely accepted as monophyletic. Partition of variance among sources of

19

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

422  uncertainty reveals that variance among simulated phylogenies is caused primarily by missing

423  species, in this case, even using the best information available to assign phylogenetically uncertain

424  taxa (PUTs). Conversely, variance among backbone phylogenies, which are estimated through

425  molecular data, is substantially smaller, indicating that efforts to improve the knowledge of

426  evolutionary history of hummingbirds (and other clades that share similar patterns of uncertainty)

427  should be concentrated on gathering molecular data for additional species (e.g. McGuire et al.

428  2014).

429  Because phylogenetic uncertainty is expected to be relatively more concentrated within

430  some clades of a phylogeny than within others, statistical analyses of species assemblages or

431  temporal segments of the phylogeny will be affected differently. Thus, statistical analyses of

432  different traits for the same group of species, using the same phylogenetic information, may be

433  differently affected by phylogenetic uncertainty, both in intensity and direction of bias. Of course,

434  comparative analyses among multiple taxonomic groups, based on independently built phylogenetic

435  hypotheses, requires additional caution, as the heterogeneity of variance among groups may

436  seriously distort results in unpredictable directions.

437  Because the effects of phylogenetic uncertainty in ecological and evolutionary analyses are

438  not subject to generalization, quantifying and accounting for phylogenetic uncertainty through

439  sensitivity analysis is required in all ecological studies. Modern techniques of sensitivity analyses

440  typically involve the application of Monte Carlo methods. Although general simulation methods,

441  such as the one used in this study, could be modified to account for phylogenetic uncertainty in

442  most evolutionary and ecological analyses, the framework for uncertainty quantification and

443  sensitivity analysis should be tailored to the purpose of the study (e.g. estimate of diversification

444  rates, community phylogenetics, comparative analysis of species traits).

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

445     Finally, the approach proposed here can be used to quantify the full spectrum of components

446     of phylogenetic uncertainty, guiding sampling strategies for future studies and allowing more

447     reliable interpretations of the relative magnitude of historical and phylogenetic components of

448     biodiversity patterns.

449

450     **Software Tools**

451     We provide two software toolkits to enable the application of the analytical strategy

452     proposed here. The first software is SUNPLIN (Martins et al. 2013;

453     https://sourceforge.net/projects/sunplin), which is capable of generating randomized phylogenies

454     after the insertion PUTs into backbone trees, with MDCCs assigned. The generated trees can then

455     be applied in any analysis that requires phylogenies as input data. SUNPLIN can be used as an

456     online web service (http://wsmartins.net/sunplin/), as a library that connects through APIs to any

457     compiled software, or directly integrated into R (http://www.ecoevol.ufg.br/pam).

458     The second software toolkit, PAM (*Phylogenetic Analysis in Macroecology*,

459     http://www.ecoevol.ufg.br/pam), is a compiled computational platform for inference of ecological

460     and evolutionary processes in a spatially explicit context. In PAM, users can not only generate

461     replicates of phylogenetic trees to be used in other software applications, but can also run several

462     statistical analyses commonly used in biodiversity analysis, while estimating and accounting for

463     multiple sources of uncertainty using the analytical framework proposed here. PAM is a work in

464     progress and will be continuously expanded in the future.

465

466     **Acknowledgements**

21

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

478

## 479    **References**

480    Abouheif, E. 1998. Random trees and the comparative method: a cautionary tale. Evolution

481        52:1197-1204.

482    Anderson, M.J. 2001. A new method for non-parametric multivariate analysis of variance. Austral

483        Ecology 26:32-46.

484    Bansal, M. S., J. G. Burleigh, O. Eulenstein, and D. Fernández-Baca. 2010. Robinson-Foulds

485        Supertrees. Algorithms for Molecular Biology 5:18.

486    Batista, M. C. G, S. F. Gouveia, D. L. Silvano and T. F. Rangel. 2013. Spatially explicit analyses

487        highlight idiosyncrasies: species extinctions and the loss of evolutionary history. Diversity

488        and Distributions (*in press*).

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

489    Beck, R. M. D., O. R. P. Bininda-Emonds, M. Cardillo, F.-G., R. Liu and A. Purvis. 2006. A

490        higher-level MRP supertree of placental mammals. Evolutionary Biology 6:93-107.

491    Bininda-Emonds, O. R., J. L. Gittleman and A. Purvis. 1999. Building large trees by combining

492        phylogenetic information: a complete phylogeny of the extant Carnivora (Mammalia).

493        Biological Reviews 74: 143-175.

494    Bininda-Emonds, O. R. P. 2004. The evolution of supertrees. Trends in Ecology and Evolution

495        19:315-322.

496    Blomberg, S. P., T. Garland, and A. R. Ives. 2003. Testing for phylogenetic signal in comparative

497        data: behavioral traits are move labile. Evolution 57:717-745.

498    Colwell, R. K. 2000. Rensch's Rule crosses the line: Convergent allometry of sexual size

499        dimorphism in hummingbirds and flower mites. American Naturalist 156:495-510.

500    Colwell, R. K. and T. F. Rangel. 2010. Modelling quaternary range shifts and richness on tropical

501        elevational gradients. Phylosofical Transactions of the Royal Society, B 365: 3695-3707.

502    Cooper E.D. 2014. Overly simplistic substitution models obscure green plant phylogeny. Trends in

503        Plant Science 19: 576-582.

504    Davies, T. J., N. J. B. Kraft, N. Salamin and E. M. Wolkovich. 2012. Incompletely resolved

505        phylogenetic trees inflate estimates of phylogenetic conservatism. Ecology 93:242-247

506    Davies, T.J., T.G. Barraclough, M.W. Chase, P.S. Soltis, D.E. Soltis and V. Savolainen. 2004

507        Darwin's abominable mystery: insights from a supertree of the angiosperms. Proceedings of

508        the National Academy of Sciences 107, 1904–1909.

509    Day, J. D., J. A. Cotton and T. G. Barraclough. 2008. Tempo and Mode of Diversification of Lake

510        Tanganyika Cichlid Fishes. PLoS One 3:e1730

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

511     Diniz-Filho, J. A. F. 2001. Phylogenetic autocorrelation under distinct evolutionary processes.

512          Evolution 55:1104-1109.

513     Diniz-Filho, J. A. F., R. D. Loyola, P. Raia, A. O. Mooers and L. M. Bini. 2014. Darwinian

514          shortfalls in biodiversity conservation. Trends in Ecology and Evolution, *in press*.

515     Doubilet P., C.B. Begg, M. C. Weinstein, P. Braun, B. J. McNeil. 1985. Probabilistic sensitivity

516          analysis using Monte Carlo simulation. A practical approach. Medical Decision Making

517          5:157-177.

518     Dunning, J. B. 2007. CRC handbook of avian body masses. 2nd ed. CRC press.

519     Efron, B. and R. J. Tibshirani. 1994. An introduction to the bootstrap. Chapman & Hall.

520     Felsensetein, J. 1985. Phylogenies and the comparative method. American Naturalist 125:1-15.

521     FitzJohn, R., W. Maddison, and S. Otto. 2009. Estimating trait-dependent speciation and extinction

522          rates from incompletely resolved phylogenies. Systematic Biology 58:595–611.

523     Gittleman, J. L. and M. Kot. 1990. Adaptation: Statistics and a null model for estimating

524          phylogenetic effects. Systematic Zoology 39:227-241.

525     Graham, C. H., J. L. Parra, C. Rahbek, and J. A. McGuire. 2009. Phylogenetic structure in tropical

526          hummingbird communities. Proceedings of the National Academy of Sciences of the United

527          States of America 106:19673-19678.

528     Helmus, M. R., T. J. Bland, C. K. Williams, and A. R. Ives. 2007. Phylogenetic measures of

529          biodiversity. American Naturalist 169:E68-E83.

530     Hernandez, F.M. and E. S. Vrba. 2005. A complete estimate of the phylogenetic relationships in

531          Ruminantia: a dated species-level supertree of the extant ruminants. Biological Reviews 80:

532          269–302.

533    Hillis, D. M., T. Heath and K. St. John. 2005. Analysis and Visualization of Tree Space. Systematic

534          Biology 54:1-12.

535    Housworth, E. A. and E. P. Martins. 2001. Random sampling of constrained phylogenies:

536          conducting phylogenetic analyses when the phylogeny is partially known. Systematic

537          Biology 50:628-639.

538    Holder, M. and P. O. Lewis. 2003. Phylogeny Estimation: Traditional and Bayesian Approaches.

539          Nature Reviews Genetics 4:275-284.

540    Huelsenbeck, J. P. and B. Rannala. 2003. Detecting correlation between characters in a comparative

541          analysis with uncertain phylogeny. Evolution 57:1237-1247.

542    Huelsenbeck, J. P., B. Rannala, and J. P. Masly. 2000. Accommodating Phylogenetic Uncertainty in

543          Evolutionary Studies. Science 288:2349-2350.

544    Ives, A. R., P. E. Midford and T. Garland Jr. 2007. Within-species variation and measuremente

545          Error in Phylogenetic Comparative Methods. Systematic Biology 56:252-270.

546    Isaac, N. J. B., S. T. Turvey, B. Collen, C. Waterman, and J. E. M. Baillie. 2007. Mammals on the

547          EDGE: Conservation Priorities Based on Threat and Phylogeny. PLoS ONE 2:e296.

548    Kuhn, T. S., A. O. Mooers, and G. H. Thomas. 2011. A simple polytomy resolver for dated

549          phylogenis. Methods in Ecology and Evolution 2:427-436.

550    Lewis, P. O., M. T. Holder and K. E. Holsinger. 2005. Polytomies and Bayesian Phylogenetic

551          Inference. Systematic Biology 54:241-253.

552    Losos, J. B. 1994. An approach to the Analysis of Comparative Data When a Phylogeny is

553          Unavailable or Incomplete. Systematic Biology 43:117-123.

554    Martins, E. P. 1996. Conducting phylogenetic comparative studies when the phylogeny is not

555          known. Evolution 50:12-22.

556 Martins, W.S., W. C. Carmo, H. J. Longo, T. C. Rosa and T. F. Rangel. 2013. SUNPLIN:

557      Simulation with Uncertainty for Phylogenetic Investigations. BMC BioInformatics 14:324-

558      335.

559 McGuire, J. A., C. C. Witt, D. L. Altshuler, and J. V. Remsen Jr. 2007. Phylogenetic systematics

560      and biogeography of hummingbirds: Bayesian and maximum likelihood analyses of

561      partitioned data and selection of an appropriate partitioning strategy. Systematic Biology

562      56:837-856.

563 McGuire, J. A., C. C. Witt, J. Remsen Jr, A. Corl, D. L. Rabosky, D. L. Altshuler, and R. Dudley.

564      2014. Molecular phylogenetics and the diversification of hummingbirds. Current Biology

565      24:910-916.

566 Jetz, W., G. H. Thomas, J. B. Joy, K. Hartmann and A. O. Mooers. 2013. The global diversity of

567      birds in space and time. Nature 491:444-448.

568 Pagel, M. D. 1999. Inferring the historical patterns of biological evolution. Nature 401:877-884.

569 Pannell, D.J. 1997. Sensitivity analysis of normative economic models: Theoretical framework and

570      practical strategies. Agricultural Economics 16: 139-152.

571 Pavoine, S. and C. Ricota. 2012. Testing for phylogenetic signal in biological traits: the ubiquity of

572      cross-product statistics. Evolution, 67:828-840.

573 Rahbek, C. and G. R. Graves. 2000. Detection of macro-ecological patterns in South American

574      hummingbirds is affected by spatial scale. Proceedings of the Royal Society of London, B

575      267:2259-2265.

576 Rahbek, C. and G. R. Graves. 2001. Multiscale assessment of patterns of avian species richness.

577      Proceedings of the National Academy of Sciences of the USA 98:4534-4539.

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

578    Ranwez, V., V. Berry, A. Criscuolo, P. H. Fabre, S. Guillemot, C. Scornavacca and E. J. P.

579         Douzery. 2007. PhySIC: a veto supertree method with desirable properties. Systematic

580         Biology 56:798–817.

581    Revell, L. J., L. J. Harmon, and D. C. Collar. 2008. Phylogenetic signal, evolutionary process, and

582         rate. Systematic Biology 57:591-601.

583    Ronquist, F. 2004. Bayesian Inference of Character Evolution. Trends in Ecology and Evolution

584         19:475-481.

585    Roquet, C., W. Thuiller and S. Lavergne. 2012. Building megaphylogenies for macroecology:

586         taking up the challenge. Ecography 36:13-26.

587    Saltelli, A., M. Ratto, T. Andres, F. Campolongo, J. Cariboni, D. Gatelli, M. Saisana, and S.

588         Tarantola. 2008. Global Sensitivity Analysis: The Primer. John Wiley & Sons.

589    Schuchmann, K.L. 1999. Family Trochilidae. Pages 468-680 in J. del Hoyo, A. Elliott, and J.

590         Sargatal, editors. Handbook of the Birds of the World. Vol. 5. Barn-owls to Hummingbirds.

591         Lynx Edicions, Barcelona.

592    Steel, M. A. and D. Penny. 1993. Distributions of tree comparison metrics - some new results.

593         Systematic Biology 42:126-141.

594    Swofford, D. L. 1991. When are phylogeny estimates from molecular and morphological data

595         incongruent? Pages 295-333 in M. M. Miyamoto and J. Cracraft, editors. Phylogenetic

596         analysis of DNA sequences. Oxford Univ. Press, New York.

597    Xi Z., L. Liu, J.S. Rest and C.S. Davis. 2014. Coalescent versus concatenation methods and the

598         placement of *Amborella* as sister to water lilies. Syst. Biol. (*in press*).

599    Webb, C. O. and M. J. Donoghue. 2005. Phylomatic: tree retrieval for applied phylogenetics.

600         Molecular Ecology Notes 5:181-183.

601                                                27

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

602 **Figure Legends:**

603

604    Figure 1: A conceptual example of the effect of different sources of uncertainty in phylogenetic

605 trees (left) on the estimation of phylogenetic relationship between species (right, represented as

606 matrices). In each matrix, cells in the lower-left half-matrix represent the existence of phylogenetic

607 information about the relationship between a pair of species, because both species are present in the

608 tree. Conversely, a dash represents the absence of such phylogenetic information, as one or both of

609 the species is missing from the tree. In the upper-right half-matrix a normal distribution represents

610 the variance in the estimated phylogenetic relationship, while zeros represent certainty. (A)

611 Hypothetical, unknown true tree, without missing species or uncertainty in species relationships.

612 (B) Consensus tree with missing species. Relationships are assumed to be known with certainty. (C)

613 Polytomic super-tree. Insertion of missing species in polytomies generates a complete tree, and

614 relationships are assumed to be known with certainty, although the tree differs from the true tree

615 (A). (D) Replication in the use of phylogenetic trees incorporates the uncertainty in the relationship

616 between species, but missing species are ignored. (E) Missing species are inserted in multiple

617 phylogenies, accounting for uncertainty in phylogenetic reconstruction and lack of a complete

618 phylogeny.

619

620 Figure 2: Schematic representation of a workflow using the analytical strategy, proposed here, to

621 account for phylogenetic uncertainty. (A) The true, but unknown, tree for a clade of 8 taxa. (B)

622 Molecular data are available for only six taxa; for taxa B and G no molecular data are available. (C)

623 Two experts, based on their knowledge (e.g. taxonomy, behavior, morphology, geographic

624 distribution, etc.), suggest possible sister taxa of the two taxa that lack molecular data.  (D) Using

625 the available molecular data and phylogenetic reconstruction methods, three backbone phylogenies

626 are proposed. The variation between backbone phylogenies arises from uncertainty in the process of

28

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

627    phylogenetic reconstruction. Among the backbone phylogenies, the taxa B and G are considered

628    *phylogenetically uncertain taxa* (PUT), because no molecular data are available for them, and

629    therefore they are missing from the backbone phylogenies. Using the information provided by the

630    experts the *most derived consensus clade* (MDCC) is found for each PUT. For PUT B, the MDCC

631    is the clade that necessarily includes taxa C, D, and A (indicated by a bold B in each backbone

632    phylogeny), whereas for PUT G the MDCC must include taxa F and H (indicated by a bold G in

633    each backbone phylogeny). (E) Statistical analyses that use phylogenetic trees as input data should

634    be replicated using samples of operational trees. In each of the operational trees, the PUTs B and G

635    were randomly inserted within their respective MDCCs. The insertion of each PUT was replicated

636    three times (columns) for each backbone tree. The variation among operational trees that use the

637    same backbone tree (each column of the operational trees) arises from variation in placement of the

638    missing taxa (as the backbone tree is not changed by the randomization process), and is caused by

639    uncertainty in the phylogenetic relationships of taxa B and G (both PUTs).

640

641    Figure 3: Hummingbird phylogenies. (A) One backbone phylogeny from McGuire et al. (2007),

642    with red internal branches indicating the *Most Derived Consensus Clades* (MDCCs) used in this

643    analysis. (B) Red taxa (polytomies) indicate PUTs inserted in the original phylogeny at the base of

644    their MDCCs. (C) Red taxa indicate PUTs inserted to yield a fully-resolved phylogeny,

645    randomizing their position within the respective MDCC.

646

647    Figure 4: Moran's I correlogram (top panels) for two hummingbird traits (body size and maximum

648    elevational range limit). The inner 95% confidence interval (C.I.) around each estimate indicates

649    variance due to sampling error, the intermediate 95% C.I. the variance due to missing species, and

650    the outer 95% C.I. the variance due to phylogenetic reconstruction. The bottom panels indicate the

29

Rangel et al., Phylogenetic Uncertainty in Evolutionary Inference

651    relative proportion of statistical error caused by each source of uncertainty across evolutionary time:

652    sampling (green), missing species (PUTs, red) and phylogenetic reconstruction (blue).

653

654    Figure 5: Hummingbird assemblages with PSV values significantly different from null expectation.

655    Blue cells indicate significant phylogenetic clustering, whereas red cells indicate significant

656    phylogenetic dispersion. (A) Standard PSV analysis, considering only non-random "sampling" from

657    the phylogeny. (B) Re-analysis of PSV accounting for three sources of error: sampling, missing

658    species, and phylogenetic reconstruction. Notice absence of significant phylogenetic dispersion and

659    decreased areas of phylogenetic clustering when all sources of uncertainty are accounted for.

660

661    Figure 6: Relative statistical error associated with PSV analysis for hummingbird species present in

662    each map cell, partitioned among sources of uncertainty. Each point within the cube has a unique

663    color that represents the relative proportions of uncertainty due to phylogeny, PUTs and sampling

664    error, as shown by the axes. In the map, the prevalence of red cells indicates sampling error as the

665    main source of uncertainty, whereas blue indicates substantial error caused by phylogenetic

666    uncertainty. Notice the absence of green and yellow areas, indicating that phylogenetic uncertainty

667    due to missing species (PUTs) is relatively irrelevant in the phylogenetic structure of assemblages.