1    FishPhyloMaker: An R package to generate phylogenies for ray-finned fishes

2    Authors: Gabriel Nakamura[1,2,*], Aline Richter[1], Bruno E. Soares[3]

3    1 – Universidade Federal do Rio Grande do Sul, Departamento de Ecologia, Bento Gonçalves

4    Avenue, 9500.

5    2 – INCT Ecology, Evolution, and Biodiversity Conservation

6    3 – Universidade Federal do Rio de Janeiro, Programa de Pós-Graduação em Ecologia

7    *correspondence author: gabriel.nakamura.souza@gmail.com

8

9    **Highlights**

10    • We provide the first automated procedure to check species names, construct

11       phylogenetic trees and calculate Darwinian shortfalls for ray-finned fishes

12       (Actinopterygii) by the R package FishPhyloMaker.

13    • This package provides functions to assemble phylogenies through a fast, reliable, and

14       reproducible method, allowing its use and replicability by specialists and non-

15       specialists in fish systematics.

16    • The package also provides an interactive procedure that gives more flexibility to the

17       user when compared with other existing tools that construct phylogenetic trees for

18       other highly speciose groups.

19    • The package includes a new method to compute Darwinian shortfalls for ray-finned

20       fishes, but the rationale of the provided algorithm can be extended in future studies to

21       be used in other groups of organisms

22

23  **Abstract**

24  Phylogenies summarize evolutionary information that is essential in the investigation of

25  ecological and evolutionary causes of diversity patterns. They allow investigating hypotheses

26  from trait evolution to the relationship between evolutionary diversity and ecosystem

27  functioning. However, obtaining a comprehensive phylogenetic hypothesis can be difficult

28  for some groups, especially those with a high number of species, that is the case for fishes,

29  particularly tropical ones. The lack of species in phylogenetic hypotheses, called Darwinian

30  shortfalls, can hinder ecological and evolutionary studies involving this group. To tackle this

31  problem, we developed FishPhyloMaker, an R package that facilitates the generation of

32  phylogenetic trees through a reliable and reproducible procedure, even for a large number of

33  species. The package adopts well-known rules of insertion based on cladistic hierarchy,

34  allowing its use by specialists and non-specialists in fish systematics. We tested the reliability

35  of our algorithm in maintaining important properties of phylogenetic distances running a

36  sensitivity analysis. We also exemplified the use of the FishPhyloMaker package by

37  constructing complete phylogenies for fishes inhabiting the four richest freshwater ecoregions

38  of the world. Furthermore, we proposed a new method to calculate Darwinian shortfalls and

39  mapped this information for the major freshwater drainages of the world. FishPhyloMaker

40  will expand the range of evolutionary and ecological questions that can be addressed using

41  ray-finned fishes as study models, mainly in the field of community phylogenetics, by

42  providing an easy and reliable way to obtain comprehensive phylogenies. Further,

43  FishPhyloMaker presents the potential to be extended to other taxonomic groups that suffer

44  from the same difficulty in the obtention of comprehensive phylogenetic hypothesis.

45  **Keywords**: Phylogenies, community phylogenetics, Darwinian shortfall, gap-analysis.

**Introduction**

46

47    Phylogenies have been widely explored in ecology in the last decades due to the development

48    of theoretical frameworks, numerical methods, and software (*e.g.,* Webb et al. 2008;

49    Felsenstein 1985). The research agenda in ecology and evolution encompasses phylogenetic

50    approaches from organismal to macroecological-scale, including trait evolution, invasion

51    ecology, metacommunity ecology, and ecosystem functioning (Cavender-Bares et al., 2009).

52    Hence, comprehensive phylogenetic trees must be available to address those topics. Large

53    phylogenies were primarily developed by combining source-trees and published-trees (the

54    supertree approach), by concatenating different data matrices of systematic phylogenetic

55    characters to generate a single tree (the supermatrix approach), or by a mix of both

56    approaches (Haeseler, 2012; Smith et al., 2009).

57          Well-established phylogenies for most of the known species are available for some

58    groups, such as terrestrial vertebrates (birds (Jetz et al., 2012), mammals (Upham et al.,

59    2019), amphibians (Jetz and Pyron, 2018), squamates (Tonini et al., 2016), sharks (Stein et

60    al., 2018), and plants (Magallón et al., 2015), which also have powerful tools to generate

61    phylogenetic trees for local/regional pools of species (*e.g.,* Webb & Donoghue 2005 for

62    mammals and plants; Jin & Qian 2019 for plants, to the others see

63    http://vertlife.org/phylosubsets/). Inversely, available phylogenies for bony fishes (Betancur

64    et al., 2017; Rabosky et al., 2018) display issues related to the taxonomic position of some

65    clades (e.g., non-monophyletic groups) and the lack of species representativeness. The latter

66    issue hampers answering some questions on the ecology and evolution of ray-finned fishes by

67    generating inaccuracy in estimates of phylogenetic signal, trait evolution, and phylogenetic

68    diversity (Seger et al., 2013; Boettiger et al., 2012a), or even impeding their calculation.

69          Ray-finned fishes (Actinopterygii) exhibit a complex evolutionary history and high

70    ecological diversity (Albert et al., 2020), making them an interesting group to address

71    questions in the interface of ecology and evolution (*e.g.*, Roa-Fuentes et al. 2019; Nakamura

72    et al. 2020). The difficulty in obtaining phylogenetic information can hinder our efforts to

73    understand fish ecology and evolution. Additionally, the lack of phylogenetic information for

74    species, *i.e.*, Darwinian shortfalls, is currently investigated in a few lineages (*e.g.*, Freitas et

75    al., 2021), which impedes the mapping of the relative demand of additional efforts needed in

76    entire regions or clades to uncover the phylogenetic history of fishes. This problem urges a

77    rapid solution in the context of the accelerated loss of species (Chase et al., 2020).

78        A short-term solution to tackle the Darwinian shortfall for ray-finned fishes would be

79    coupling the phylogenetic information with cladistic classification to produce comprehensive

80    phylogenies (Diniz-Filho et al., 2013). This solution is laborious and lacks reproducibility

81    when adding many species, and the specific steps are not precisely documented when did "by

82    hand" procedures (Webb et al., 2008). An alternative would be using molecular techniques to

83    generate comprehensive phylogenies (e.g. Pie et al., 2021). However, it demands high

84    expertise and high financial investment (Roquet et al., 2013), limiting factors for several

85    institutions. Therefore, automatizing the procedures of constructing comprehensive

86    phylogenies using the information from cladistic hierarchy, as suggested by Diniz-Filho et al

87    (2013), provides a more reliable, accessible, and short-term solution for evolutionary

88    ecologists. The technique produces reliable phylogenetic information for community

89    phylogenetics (Li et al., 2019).

90        In order to tackle the problem of obtention of comprehensive fish phylogenies in a

91    reliable and reproducible way, we developed the FishPhyloMaker. This freely available R

92    package facilitates the obtention of phylogenetic trees for ray-finned fishes. FishPhyloMaker

93    automates the insertion procedure of species in the most comprehensive phylogeny (Rabosky

94    et al., 2018) of ray-finned fishes following their taxonomic hierarchy. We illustrated how the

95    FishPhyloMaker package solves the problem of obtaining comprehensive phylogenies by

96    constructing phylogenetic trees for species inhabiting more than 3000 freshwater basins

97    globally (Tedesco et al., 2017). Further, we developed a new method to quantify the

98    Darwinian shortfalls, which we illustrate by mapping the Darwinian shortfalls for the

99    abovementioned basins. Finally, we performed a sensitivity analysis to evaluate how our

100   method preserves characteristics of the phylogenetic tree (pairwise distances among species

101   and evolutionary distinctiveness), even with a varying number of inserted taxa. Our package

102   overcomes the main problems associated with manually building phylogenies for ray-finned

103   fishes by following a specific and documented procedure and reducing the manual labor in

104   large phylogenies.

105

106   **Methods**

107   **Inside the Fish(PhyloMaker): an overview of the package**

108        A stable version of FishPhyloMaker can be downloaded from the CRAN repository

109   (https://cran.r-project.org/web/packages/FishPhyloMaker/index.html), and a development

110   version is available at the GitHub repository

111   (https://github.com/GabrielNakamura/FishPhyloMaker). All analyses shown here were

112   performed using the development version of FishPhyloMaker.

113   FishPhyloMaker is a freely available R package containing three main functions,

114   *FishTaxaMaker*, *FishPhyloMaker,* and *PD_deficit*. Below, we describe the functions to

115   generate phylogenetic trees, highlighting the input data, intermediate steps, and output

116   objects. Brief descriptions of the package functions are available in Table 1.

117

118   *FishTaxaMaker*

119   The *FishTaxaMaker* function checks the validity of species names provided by the user and

120   prepares a formatted data frame for the *FishPhyloMaker* function. The input data must be a

121    string vector or a data frame containing a list of species from the regional pool or an

122    occurrence matrix (sites x species). The genus and specific epithet (or subspecies) must be

123    separated by underline (e.g., *Genus_epithet*). The function first classifies the provided species

124    names as valid or synonymies based on Fishbase (Froese & Pauly, 2000) using the *rfishbase*

125    package (Boettiger et al., 2012b). A new column summarizes names initially valid and the

126    current valid names substituting identified synonymies. Unknown species to Fishbase are

127    printed in the command line, and the user must manually inform the Family of these species.

128    If the user types a Family not recognized in the FishBase, the user is asked to check the

129    spelling and type the Order of this family. The output of the function is a list containing three

130    elements: 1) a data frame displaying the taxonomic information (Valid name, Subfamily,

131    Family, Order, Class, and SuperClass) for each species; 2) a data frame displaying the

132    taxonomic information (Species, Family, and Order), only for the valid species; 3) a character

133    vector displaying the species names not found in Fishbase.

134

135    Table 1: Functions presented in the package FishPhyloMaker and their descriptions.

| Function | Description |
| --- | --- |
| *FishTaxaMaker*() | Checks species names according to Fishbase and prepares the species list for the other functions in the package. |
| *whichFishAdd*() | Identifies the species already included in the backbone tree and in which taxonomic level each remaining species will be inserted. |
| *FishPhyloMaker*() | Builds the phylogeny and may return a data frame identifying step-by-step the performed insertions. |
| *PD_deficit*() | Calculates the Darwinian shortfall for the |

provided species list through a Phylogenetic

Diversity (PD Faith (1992)) ratio

136

137     *FishPhyloMaker*

138     *FishPhyloMaker* is the core function of the package. This function builds a phylogenetic tree

139     for the provided species list by inserting in and pruning species from the Rabosky et al.,

140     (2018) phylogenetic tree (Figure 1) downloaded by the fishtreeoflife R package (Chang et al.

141     2019). This phylogeny is the most up-to-date and comprehensive phylogenetic hypothesis for
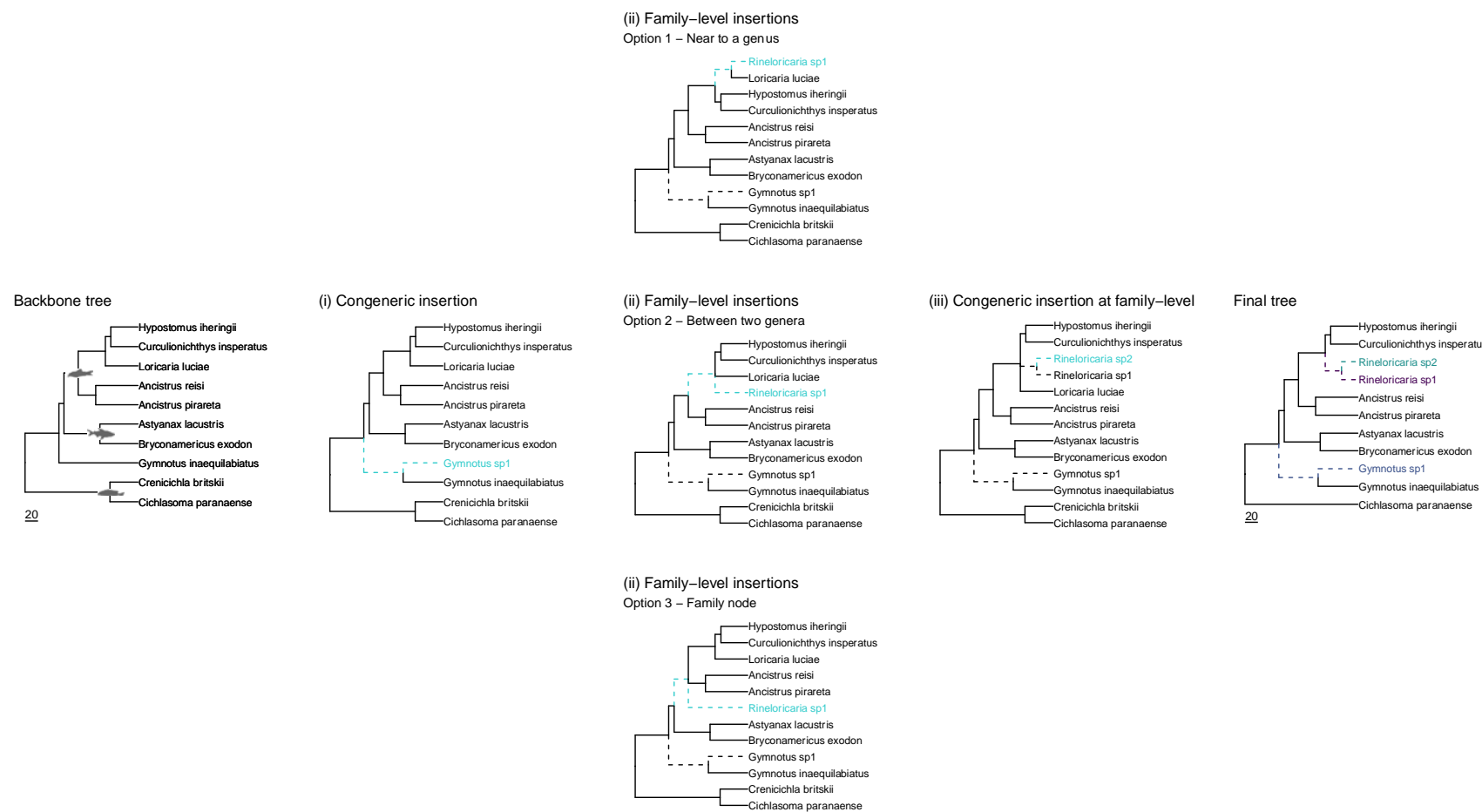
142     ray-finned fishes.

143         The input for the *FishPhyloMaker* function can be the second element in the list

144     returned by the *FishTaxaMaker* (Taxon_data_FishPhyloMaker) function or a manually

145     constructed data frame with the same configuration (species, family, and order names for

146     each taxon). The function also contains three logical arguments: insert.base.node,

147     return.insertions and progress.bar. These three arguments are set by default as FALSE,

148     TRUE, and TRUE, respectively, and allow the user to choose if the species must be at the

149     base node of families/orders, if the insertions made by each species must be shown in the

150     output and if a progress bar must be shown in the console.

151         The function works sequentially, first identifying which of the provided species are in

152     the backbone phylogenetic tree (Rabosky et al., 2018). If all of them are already present in

153     the backbone tree, the function returns a pruned one. If any of the provided species is not in

154     the backbone tree, the function performs a four-level insertion routine. First, species from

155     genera already included in the backbone tree are inserted as polytomies at the most recent

156     ancestral node that links all congeneric species or as the sister species of the only species

157     representing a genus in the backbone tree, as shown in *i* in Figure 1. In the case of *i* in Figure

158     1 the branch length is divided at half of its length and the species is inserted. Second, species

159    not inserted in the previous step are then inserted at the family level by an interactive

160    procedure using a returned list of all the genera within the same family of the target species.

161    The user has the option to insert the target species as a sister taxon to a genus (*ii* in Figure 1,

162    option 1, near to *Loricaria* genus), between two genera (*ii* in Figure 1, option 2, between

163    genus *Loricaria* and *Hypostomus*), or at the node of the family (*ii* in Figure 1, option 3). If the

164    user enters a single genus from the list, the function splits its branch and inserts the target as a

165    sister taxon of this genus (option 1). If the user enters two genera separated by a blank space,

166    the function inserts the target species as a polytomy at the most recent node that links the

167    selected genera (option 2). If the user enters the family name, the function attaches the target

168    species at the family node as a polytomy (option 3). Third, if any remaining species can now

169    be inserted at the genus level, the function repeats the first procedure but records it as a

170    Congeneric family-level insertion by splitting the branch length of the congeneric species at

171    half of its length (*iii* in Figure 1). Fourth, remnant species are inserted at the order level

172    following similar to the second step, by an interactive procedure using a returned list of all

173    the families within the order of the target species. Hence, the user may specify a family to

174    insert the target species as sister taxon (option 1), two families to insert it as a polytomy at the

175    most recent node linking them (option 2), or the order to insert it as a sister taxon (option 3).

176    The function will not perform insertions steps beyond the order level because it would add

177    too much uncertainty to the phylogenetic tree.

178        Setting the argument `insert.base.node` as TRUE automatically inserts the target

179    species from the second and fourth steps in the family and order nodes, respectively. This

180    option facilitates the insertion of a large number of species or species with the unknown

181    phylogenetic position. The default output is a list with two objects: (i) the pruned tree

182    including only the provided species list (Final tree in Figure 1); (ii) a data frame identifying if

183    each provided species was initially present in the backbone tree, in which step it was inserted,

184     or not inserted at all. This data frame will flag each species with one of the six classification

185     based on the insertion procedure: 1 – Present in tree will indicate species that were already

186     present in the backbone tree; 2 – Congeneric insertion will indicate species that present at

187     least one species of the same genus in the backbone tree and was inserted as congeneric of

188     this species; 3 – Family insertion will indicate inserted species that did not present any

189     congeneric species at backbone tree, but had at least one species of the same family in

190     backbone tree; 4 – Congeneric at Family-level will indicate species that was added as

191     congeneric after another species of the same genus was inserted at the Family level; 5 – Order

192     insertion will indicate inserted species that did not presented any species of the same family

193     in the backbone tree and must be inserted near to an extant family or in node corresponding

194     to the order root in the backbone tree; 6 – not inserted will indicate species that did not

195     present any species of the same order in the backbone tree, therefore was not inserted due

196     their high uncertainty in the phylogenetic position.

7

8    Figure 1: Schematic representation of insertion and subsetting procedure performed by the FishPhyloMaker() function. Here we used a hypothetical phylogeny

9    containing ten species and four families (silhouettes inside the tree) as the backbone phylogeny. Step (i) represents the congeneric level of insertion. Step (ii)

0    represents the three options that the user may choose in the Family-level round of insertions (Option 1 – near to a genus; Option 2 – between two genera;

1    Option 3 – at the family node). (iii) represents the congeneric insertions at the family level and, finally, the final pruned tree containing only the species of

2    interest.

203     *PD_deficit*

204         The *PD_deficit* function calculates a measure of Darwinian shortfalls following

205     Equation 1:

206 $$PD_{inserted} \Big/ PD_{inserted} + PD_{present\ in\ tree} \qquad (1)$$

207

208     In this function, $PD_{inserted}$ is the sum of the branch lengths of species in the phylogenetic tree

209     before the insertion procedure. $PD_{present\ in\ tree}$ is the sum of branch lengths of the species

210     inserted in the tree. Therefore, the Darwinian deficit ranges from 0 (all species already

211     present in the backbone tree before the insertion procedure) to 1 (all the species in the

212     phylogenetic tree were inserted and were not presented in backbone phylogeny). PD_deficit

213     function returns a vector with three values, the Darwinian shortfall (Equation 1), the total

214     phylogenetic diversity calculated as the sum of branch lengths of the tree ($PD_{total}$) with all

215     species provided by the user, the sum of branch lengths inserted ($PD_{inserted}$) in the tree and

216     that was already present in the backbone tree ($PD_{present\ in\ tree}$). It is worth noting that the sum

217     of $PD_{inserted}$ and $PD_{present}$ are complementary, summing up to $PD_{total}$. To calculate the

218     Darwinian shortfall through the *PD_deficit* function, the user must provide a phylogenetic

219     tree and a table of insertions, both obtained from the *FishPhyloMaker* function.

220

221     *Sensitivity analysis*

222     We performed a sensitivity analysis to assess how the insertion procedure implemented

223     herein and the amount of inserted species affect two characteristics of phylogenetic trees: the

224     mean pairwise distance among species and the phylogenetic distinctiveness.

225         We 1) randomly change the name of a subsample of species within Rabosky's

226     phylogeny. Then, 2) we built a phylogeny for the species sampled with changed names in the

227     previous step using the FishPhyloMaker function. Finally, we computed: 3) the matrix

228    correlation (Pearson correlation) between the cophenetic distances of the subsampled species

229    in Rabosky's phylogeny and the FishPhyloMaker phylogeny; and 4) the Pearson correlation

230    between the phylogenetic distinctness values for the Rabosky's and FishPhyloMaker

231    phylogenies. The evolutionary distinctness was calculated as the equal splits measure that is

232    the sum of the contribution of all branches of a given lineage divided among its daughter

233    branches (Redding and Mooers, 2006). Evolutionary distinctness measure was calculated

234    using the phyloregion package (Daru et al., 2020).

235         The abovementioned steps (1, 2 and 3) were repeated 100 times for eleven different

236    quantities (10%, 15%, 20%, 25%, 30%, 35%, 40%, 50%, 55%, 60%) of subsampled species

237    from Rabosky's phylogeny and inserted by the FishPhyloMaker function.

238

239    *Illustrating the use of FishPhyloMaker package*

240    We provide an example of the usage of the *FishPhyloMaker* package by creating a

241    phylogenetic tree using a global dataset of freshwater fishes inhabiting 3,119 freshwater

242    drainage basins that cover more than 80% of the Earth surface and 14886 species (Tedesco et

243    al., 2017). This dataset allowed in-depth investigation on the global patterns of species

244    distribution and their evolutionary determinants (*e.g.*, Miller & Román-Palácios, 2021). We

245    built a phylogenetic for all species presented in Tedesco´s et al. dataset and mapped all the

246    insertions realized. Moreover, we used this same dataset to demonstrate how to map

247    Darwinian shortfalls, calculated following Equation 1 through *PD_deficit* function for all the

248    drainage basins in the Tedesco et al. (2017) dataset. All the analyses were performed using

249    the development version of the FishPhyloMaker package, which can be downloaded using
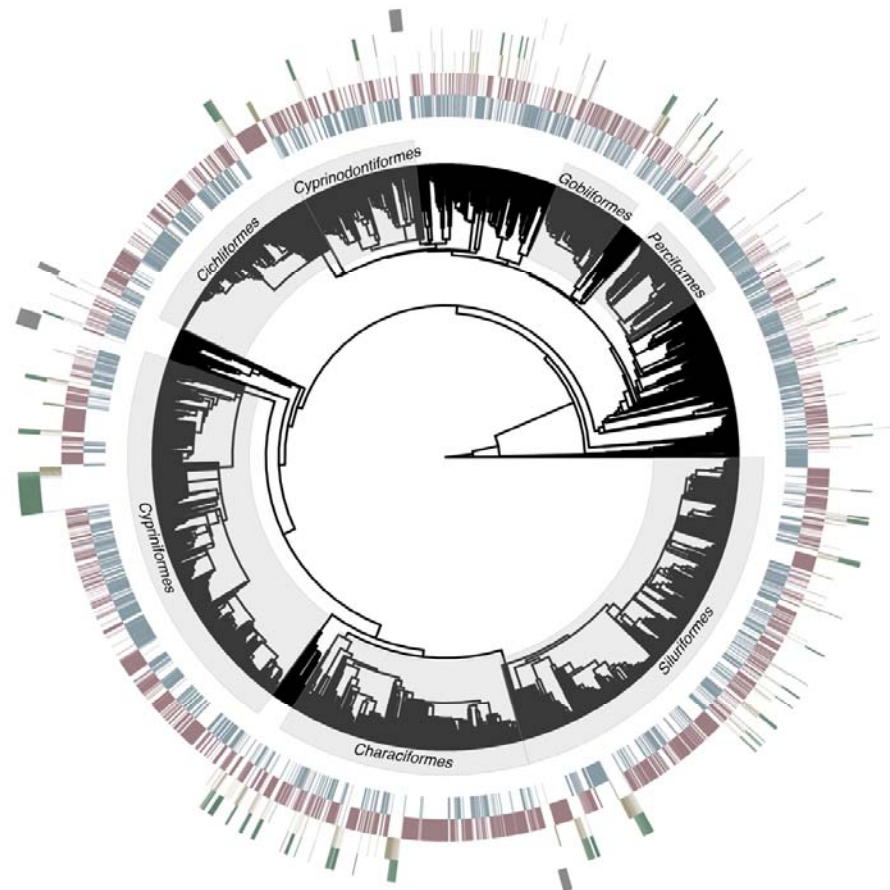
250    the following command line:

251    `devtools::install_github("GabrielNakamura/FishPhyloMaker", ref = "main",`

252    `build_vignettes = TRUE)`

253   We recommend that the user updates all the requested packages to avoid errors related to

254   packages versions. We first prepared the fish occurrence by checking the validity of its names

255   by using the function *FishTaxaMaker*. The occurrence matrix encompassed 14,886 species,

256   from which 13,992 were valid names. The remaining 961 names were substituted by their

257   corresponding valid names according to FishBase. We applied the *FishPhyloMaker* function

258   to build a phylogenetic tree containing all the 14,886 species with valid names retrieved from

259   *FishTaxaMaker* (Figure 2). For simplicity and reproducibility, we set the argument

260   `insert.base.node` as TRUE, thus, inserting all species at the base node of its corresponding

261   family and order when needed. We also set the argument `return.insertions = TRUE` for

262   retrieving the insertion information of each species. Then, we applied the *PD_deficit* function

263   to calculate the Darwinian shortfall for all the freshwater basins of the world harboring at

264   least two species (Tedesco et al. 2017). The *PD_deficit* function was calculated considering

265   congeneric insertions and insertions at the family level, however, the function may also

266   include other levels of phylogenetic insertion, like order insertions. All the codes need to

267   fully reproduce these analyses are provided at the GitHub repository

268   (GabrielNakamura/MS_FishPhyloMaker). For further explanations and examples illustrating

269   the usage of functions in the FishPhyloMaker package, the user can assess the package

270   website https://gabrielnakamura.github.io/FishPhyloMaker/index.html and see the Articles

271   section.

272   **Results**

273   The entire insertion procedure lasted approximately three hours using one core from a

274   computer machine with an i5 processor. A total of 11,569 species were inserted, 6,418

275   species were already present in the backbone phylogeny, and 181 were not inserted at all,

276   resulting in a phylogenetic tree containing 14,705 species (Figure 2). We also showed in

277   Figure 2 all the insertions realized through the FishPhyloMaker function and the seven orders

278     of ray-finned fishes with that present the highest number of species. We can see in Figure two

279     that the insertions are evenly distributed throughout the phylogenetic tree.



280

281     Figure 2: Phylogenetic tree obtained from FishPhyloMaker, containing 14,705 finned-ray

282     species with their respective insertions. We also highlight in the gray rectangles the seven

283     most speciose Orders.

284

285        We also depicted all the insertions made by FishPhyloMaker for all freshwater

286    Ecoregion of the world. This was only possible because FishPhyloMaker flags all the

287    insertions made during the insertion procedure. Figure 3 shows that Neotropics and

288    Afrotropics regions exhibited the largest number of species inserted. On the contrary, despite

289    the great area and number of basins, the Nearctic Ecoregion presented the smallest percentage

290    of insertions, most of them congeneric. All Ecoregions and the percentage of species

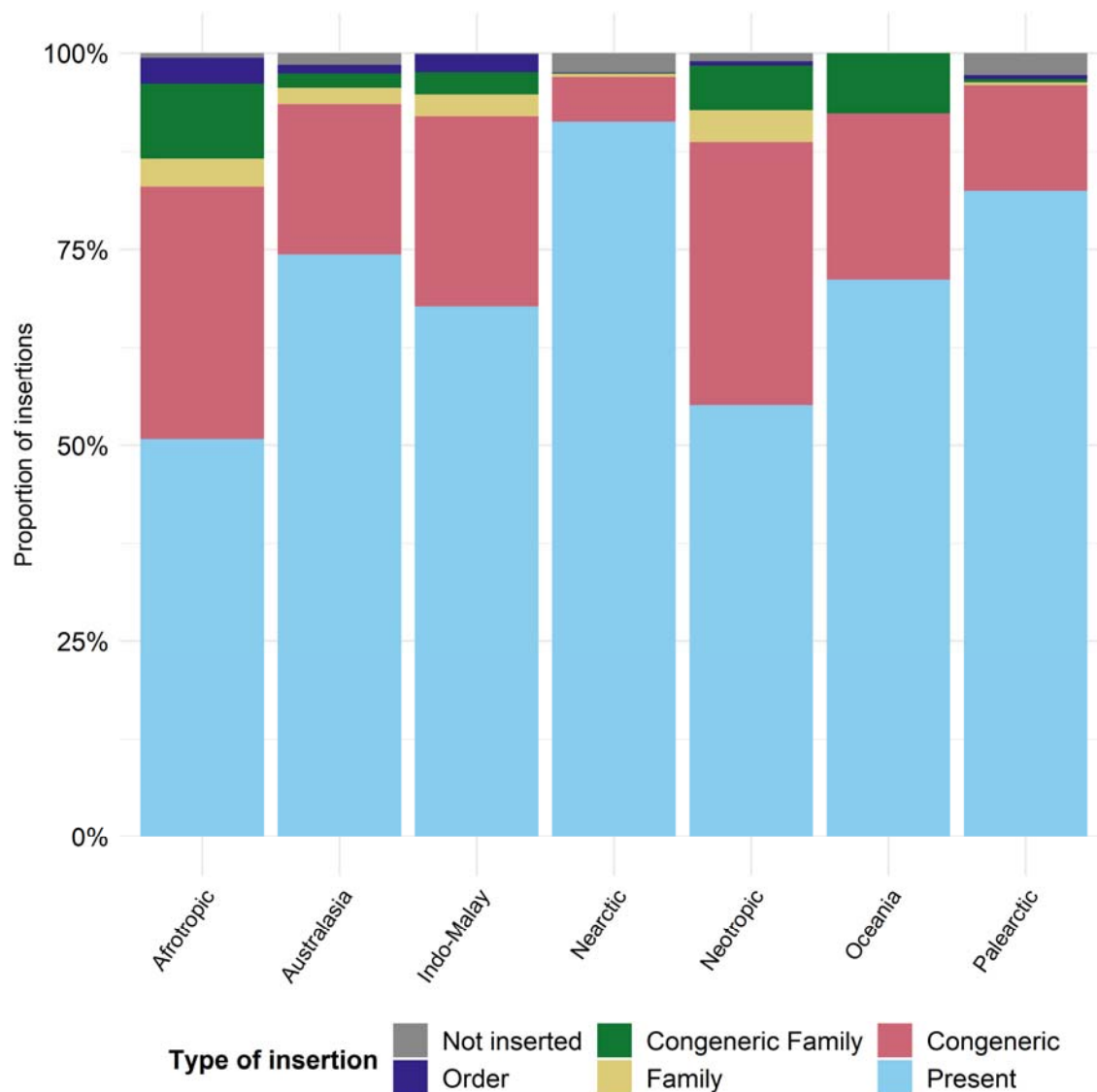291    insertions per level are shown in Figure 3.

Figure 3: Barplot showing the percentage of species inserted in each one of the seven freshwater ecoregions of the world and their respective type of insertions mapped by FishPhyloMaker package.

We spatialized the Darwinian shortfalls per basin and observed that tropical regions exhibited the highest shortfalls, while northern sites had the lowest (Figure 4). The highest values of Darwinian shortfalls were found in Afrotropics and Neotropics, as some drainages did not harbor any (or only a few) species in the Rabosky's phylogeny. The grey areas correspond to sites that do not present species occurrences accordingly to Tedesco et al. (2017) or presented less than two occurrences for the Order considered. We also depicted the

302    Darwinian shortfalls for the four major orders in terms of species richness (bottom maps in

303    Figure 4). For all the groups, the highest values of Darwinian shortfalls were found in the

304    neotropical region, except for Cypriniformes, the group responsible for the highest values of

305    Darwinian shortfalls in the watersheds in Asia and some basins in North America.
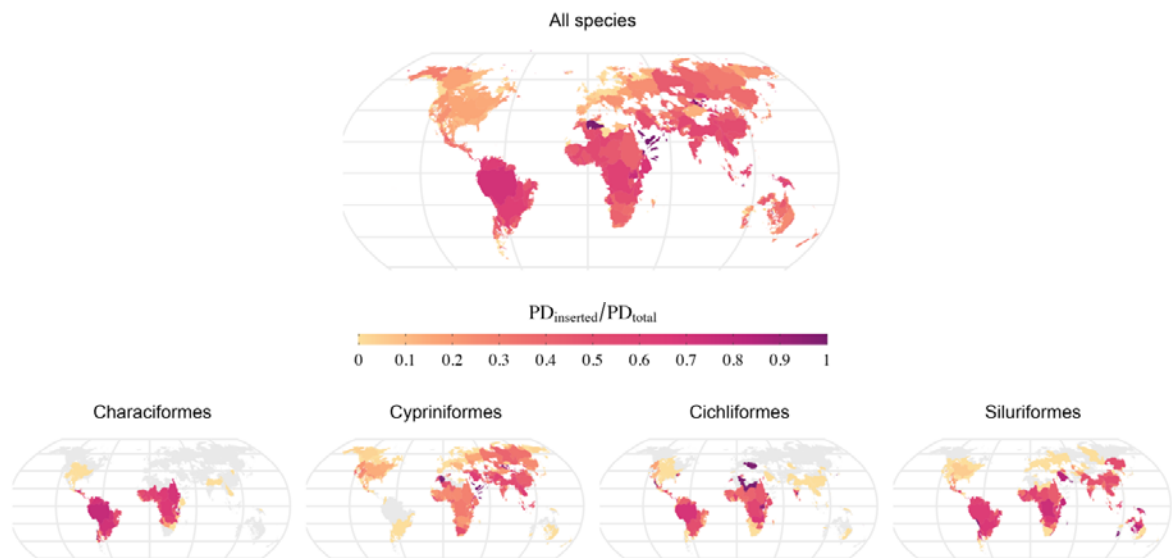
306



308    Figure 4: Global distribution of the Darwinian Shortfalls for ray-finned fishes, based on

309    freshwater species occurrences in more than 3000 basins. Values near to 1 indicate a high

310    Darwinian shortfall (a large number of congeneric insertions), while values near zero indicate

311    low shortfalls. We depicted the Darwinian shortfall for the four major orders in terms of

312    species richness (Characiformes, Cypriniformes, Cichliformes, and Siluriformes). Gray color

313    indicates areas with no occurrence of species for a given order.

314

315    The sensitivity analysis highlights the strong correlation between the cophenetic distances of

316    Rabosky's and FishPhyloMaker phylogenies (Figure 5 B) even in varying levels of taxa

317    insertions. Inversely, an increasing number of insertions on Rabosky's phylogeny reduced the

318    correlation between phylogenetic distinctness in the original phylogeny and that assembled

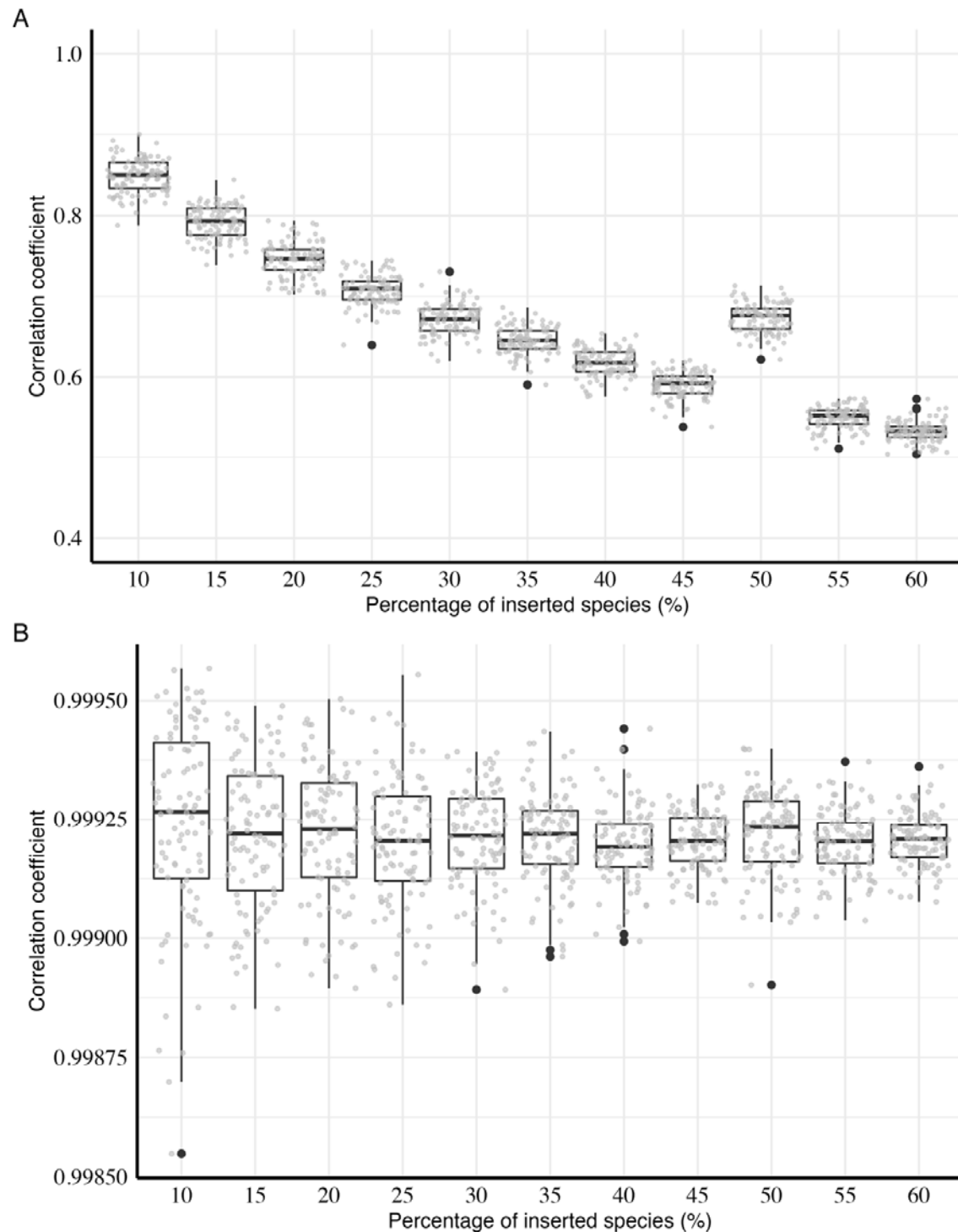319    by FishPhyloMaker (Figure 5 A).

320

321

Figure 5: Barplots showing the correlation of evolutionary distinctness values (A) and

between cophenetic distances (B) calculated from original phylogeny and inserted

phylogenies with varying percentages of species inserted in the original phylogeny. Grey dots

325     represent individual correlation values. The lower and upper hinges in boxplots represent the

326     first and third quantiles while the middle hinge represents the median.

327     **Discussion**

328     We provided a user-friendly, fast, reliable, and reproducible way to construct phylogenetic

329     trees for a megadiverse group (Actinopterygii). The FishPhyloMaker package is in line with

330     tools developed for plants, such as Phylomatic (C++ application) and V.PhyloMaker (R

331     package) (Jin and Qian, 2019; Webb and Donoghue, 2005), but includes different features.

332     These features include new options for inserting species through an interactive procedure in

333     phylogenies and recording insertions. The latter feature allows a better systematization of

334     building supertrees and calculating the first, to our knowledge, quantitative measure of the

335     Darwinian shortfall.

336         Whereas Phylomatic allows the insertion of absent species only as congeneric or at

337     the node corresponding to the family of the focal species (Webb and Donoghue, 2005), the

338     FishPhyloMaker package delivers options through an interactive procedure of insertion. The

339     performed insertions can be easily recorded in an R script, providing flexibility and the same

340     level of reproducibility as other algorithms designed for similar purposes (*e.g.*, Jin and Qian,

341     2019). This interactive option is a novelty when compared to similar insertion algorithms

342     (*e.g.*, Phylomatic).

343         The spatial distribution of the Darwinian shortfall is paramount to guide our future

344     efforts to understand the history of life. The phylogenetic gaps in the knowledge of ray-

345     finned fishes are geographically biased, with tropical basins presenting higher Darwinian

346     shortfalls levels, as evidenced in this study. This gap in evolutionary knowledge could lead to

347     a bias in evaluating the effects of evolutionary history and the interpretation of

348     macroecological patterns for fish assemblages in these regions, which can affect conservation

349     decisions based on the phylogenetic dimension of diversity (Assis, 2018).

350    Several biological and sociological factors can explain the observed bias in Darwinian

351    shortfalls. First, the regions exhibiting the most significant Darwinian gaps also exhibit the

352    largest freshwater fish diversity, which we can not describe at the same speed as less

353    biologically rich areas (Hortal et al., 2015). Second, on-ground accessibility, human

354    occupation, and economic development constrain investments in biodiversity research

355    (Moura et al., 2018; Moura and Jetz, 2021), which is probably more pronounced in tropical

356    regions than temperate ones, which may hamper field sampling and phylogenetic analyses.

357    Despite being more simple when compared with other insertion methods (e.g., Pearse

358    and Purvis, 2013), FishPhyloMaker provided reliable results by preserving important

359    characteristics of the phylogenetic tree, as we showed through the sensitivity analysis.

360    Commonly used measures of phylogenetic diversity are based on the pairwise distance of

361    species from a phylogenetic tree (e.g., Kraft et al., 2007; Webb et al., 2002), and we showed

362    that the algorithm implemented in FishPhyloMaker successfully preserve the distances

363    among species in the phylogenetic tree even for a great number of insertions.

364

365    *Limitations and possible applications*

366    Future developments of the package should consider the Catalog of Fishes (van der

367    Laan et al., 2021) to improve the nomenclature checking procedures. Despite Fishbase being

368    a widely used database to check for the taxonomic classification of fishes, it may present

369    delays in updating taxonomic information because it is not its primary purpose. Inversely, the

370    Catalog of Fishes is an authoritative taxonomic list frequently updated.

371    An inherent limitation of the phylogenetic hypothesis produced by FishPhyloMaker is

372    the large number of polytomies resulting from the insertion procedures. We recommend that

373    users directly assess how the phylogenetic uncertainty affects further analysis when not using

374    a fully solved phylogenetic tree (Martins et al., 2013). Furthermore, we recommend caution

375 in the use of FishPhyloMaker phylogenies to compute measures that depend on speciation

376 events (e.g., evolutionary distinctiveness and other split-based metrics) since the insertion

377 procedure modifies the split events in the tree as shown in the sensitivity analysis.

378 These limitations do not preclude the package applicability for studies in phylogenetic

379 community ecology since synthesis phylogenies do not significantly impact phylogenetic

380 diversity indices as showed by previous studies (Li et al., 2019) and confirmed in ours

381 (through sensitivity analysis). Moreover, this is the only automated tool able to provide a

382 complete phylogenetic tree that can easily handle large datasets. FishPhyloMaker can be

383 relevant for addressing several critical questions in ecology and evolution by facilitating the

384 obtention of phylogenetic hypotheses for local pools of ray-finned fishes. This facilitation can

385 be essential for regions with a large gap in the phylogenetic knowledge of fishes, such as the

386 Neotropical region (Albert et al., 2020). Such phylogenetic hypotheses allow understanding

387 how ecological traits evolved or how the current and past environmental conditions selected

388 the lineages in different areas.

389 Biogeographical studies are usually restricted to one or a few lineages at larger scales

390 due to the availability of molecular phylogenies (e.g. García-Andrade et al., 2021) or with

391 phylogenies with a considerable number of absent species (Miller, 2021). The

392 FishPhyloMaker package facilitates large-scale investigations on the biogeographic history of

393 the most diverse group of vertebrates on Earth, the Actinopterygians, helping us understand

394 the processes that drive this high diversity. Finally, we can map where the lack of

395 phylogenetic information is the most critical once the function returns the insertion level of

396 species. This information can directly elucidate the patterns of the Darwinian shortfalls for

397 ray-finned fishes, contributing not only to direct sampling and studying efforts but also to

398 evidence the need for increased efforts to decolonize science (Trisos et al., 2021). Therefore,

399 we expect that the FishPhyloMaker package reduces the gaps and barriers to addressing

400   ecological and evolutionary questions due to the difficulty or lack of a reliable phylogenetic

401   hypothesis for local and regional pools of ray-finned fishes.

402

403   **Contributions**

404   GN Conceptualization; Data curation; Formal Analysis; Methodology; Software; Writing –

405   original draft. AR Data curation; Methodology; Software, Writing – review, and editing. BES

406   Writing – original draft; Methodology.

407

**Acknowledgments**

GN is a member of the National Institutes for Science and Technology (INCT) in

Ecology, Evolution, and Biodiversity Conservation, supported by MCTIC/CNPq (proc.

465610/2014-5). BES and AR are grateful to FAPERJ and CAPES for their postdoctoral and

doctoral grants, respectively. The authors also thank valuable suggestions made by LDS

Duarte and other two anonymous Reviewers. Brazilian science resists.

**References**

Albert, J.S., Tagliacollo, V.A., Dagosta, F., 2020. Diversification of Neotropical Freshwater Fishes. Annu. Rev. Ecol. Evol. Syst. 51, 27–53. https://doi.org/10.1146/annurev-ecolsys-011620-031032

Assis, L.C.S., 2018. Revisiting the Darwinian shortfall in biodiversity conservation. Biodivers. Conserv. 27, 2859–2875. https://doi.org/10.1007/s10531-018-1573-3

Betancur, R.R., Wiley, E.O., Arratia, G., Acero, A., Bailly, N., Miya, M., Lecointre, G., Ortí, G., 2017. Phylogenetic classification of bony fishes. BMC Evol. Biol. 17, 1–40. https://doi.org/10.1186/s12862-017-0958-3

Boettiger, C., Coop, G., Ralph, P., 2012. Is your phylogeny informative? Measuring the power of comparative methods. Evolution (N. Y). 66, 2240–2251. https://doi.org/10.1111/j.1558-5646.2011.01574.x

Cavender-Bares, J., Kozak, K.H., Fine, P.V.A., Kembel, S.W., 2009. The merging of community ecology and phylogenetic biology. Ecol. Lett. 12, 693–715. https://doi.org/10.1111/j.1461-0248.2009.01314.x

Chang, J., Rabosky, D.L., Smith, S.A., Alfaro, M.E., 2019. An r package and online resource for macroevolutionary studies using the ray-finned fish tree of life. Methods Ecol. Evol. 10, 1118–1124. https://doi.org/10.1111/2041-210X.13182

Chase, J.M., Blowes, S.A., Knight, T.M., Gerstner, K., May, F., 2020. Ecosystem decay exacerbates biodiversity loss with habitat loss. Nature 584, 238–243. https://doi.org/10.1038/s41586-020-2531-2

Daru, B.H., Karunarathne, P., Schliep, K., 2020. phyloregion: R package for biogeographical regionalization and macroecology. Methods Ecol. Evol. 11, 1483–1491. https://doi.org/10.1111/2041-210X.13478

Diniz-Filho, J.A.F., Loyola, R.D., Raia, P., Mooers, A.O., Bini, L.M., 2013. Darwinian

440    shortfalls in biodiversity conservation. Trends Ecol. Evol. 28, 689–695.

441    https://doi.org/10.1016/j.tree.2013.09.003

442 Faith, D.P., 1992. Conservation evaluation and phylogenetic diversity. Biol. Conserv. 61, 1–

443    10. https://doi.org/10.1016/0006-3207(92)91201-3

444 Felsenstein, J., 1985. Phylogenies and the comparative method. Am. Nat. 125, 1–15.

445    https://doi.org/0003-0147/85/2501-0001

446 Freitas, T.M. da S., Stropp, J., Calegari, B.B., Calatayud, J., De Marco, P., Montag, L.F. de

447    A., Hortal, J., 2021. Quantifying shortfalls in the knowledge on Neotropical

448    Auchenipteridae fishes. Fish Fish. 22, 87–104. https://doi.org/10.1111/faf.12507

449 van der Laan, R., Fricke, R. & Eschmeyer, W. N. (eds) 2021. ESCHMEYER'S CATALOG

450    OF FISHES: CLASSIFICATION. (http://www.calacademy.org/scientists/catalog-of-

451    fishes-classification/).Electronic version accessed dd mmm 2021.

452 García-Andrade, A.B., Carvajal-Quintero, J.D., Tedesco, P.A., Villalobos, F., 2021.

453    Evolutionary and environmental drivers of species richness in poeciliid fishes across the

454    Americas. Glob. Ecol. Biogeogr. 30, 1245–1257. https://doi.org/10.1111/geb.13299

455 Haeseler, A. V., 2012. Do we still need supertrees? BMC Biol. 10, 2–5.

456    https://doi.org/10.1186/1741-7007-10-13

457 Hortal, J., De Bello, F., Diniz-Filho, J.A.F., Lewinsohn, T.M., Lobo, J.M., Ladle, R.J., 2015.

458    Seven Shortfalls that Beset Large-Scale Knowledge of Biodiversity. Annu. Rev. Ecol.

459    Evol. Syst. 46, 523–549. https://doi.org/10.1146/annurev-ecolsys-112414-054400

460 Jetz, W., Pyron, R.A., 2018. The interplay of past diversification and evolutionary isolation

461    with present imperilment across the amphibian tree of life. Nat. Ecol. Evol. 2, 850–858.

462    https://doi.org/10.1038/s41559-018-0515-5

463 Jetz, W., Thomas, G.H., Joy, J.B., Hartmann, K., Mooers, A.O., 2012. The global diversity of

464    birds in space and time. Nature 491, 444–448. https://doi.org/10.1038/nature11631

465    Jin, Y., Qian, H., 2019. V.PhyloMaker: an R package that can generate very large

466        phylogenies for vascular plants. Ecography (Cop.). 42, 1353–1359.

467        https://doi.org/10.1111/ecog.04434

468    Kraft, N.J.B., Cornwell, W.K., Webb, C.O., Ackerly, D.D., 2007. Trait evolution, community

469        assembly, and the phylogenetic structure of ecological communities. Am. Nat. 170, 271–

470        283. https://doi.org/10.1086/519400

471    Li, D., Trotta, L., Marx, H.E., Allen, J.M., Sun, M., Soltis, D.E., Soltis, P.S., Guralnick, R.P.,

472        Baiser, B., 2019. For common community phylogenetic analyses, go ahead and use

473        synthesis phylogenies. Ecology 100, 1–15. https://doi.org/10.1002/ecy.2788

474    Magallón, S., Gómez-Acevedo, S., Sánchez-Reyes, L.L., Hernández-Hernández, T., 2015. A

475        metacalibrated time-tree documents the early rise of flowering plant phylogenetic

476        diversity. New Phytol. 207, 437–453. https://doi.org/10.1111/nph.13264

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

478        Simulation with Uncertainty for Phylogenetic Investigations. BMC Bioinformatics 14.

479        https://doi.org/10.1186/1471-2105-14-324

480    Miller, E.C., 2021. Comparing diversification rates in lakes , rivers , and the sea 1–19.

481        https://doi.org/10.1111/evo.14295

482    Moura, M.R., Costa, H.C., Peixoto, M.A., Carvalho, A.L.G., Santana, D.J., Vasconcelos,

483        H.L., 2018. Geographical and socioeconomic determinants of species discovery trends

484        in a biodiversity hotspot. Biol. Conserv. 220, 237–244.

485        https://doi.org/10.1016/j.biocon.2018.01.024

486    Moura, M.R., Jetz, W., 2021. Shortfalls and opportunities in terrestrial vertebrate species

487        discovery. Nat. Ecol. Evol. 5, 631–639. https://doi.org/10.1038/s41559-021-01411-5

488    Nakamura, G., Vicentin, W., Súarez, Y.R., Duarte, L., 2020. A multifaceted approach to

489        analyzing taxonomic, functional, and phylogenetic β diversity. Ecology.

490     https://doi.org/10.1002/ecy.3122

491  Pearse, W.D., Purvis, A., 2013. phyloGenerator: An automated phylogeny generation tool for

492     ecologists. Methods Ecol. Evol. 4, 692–698. https://doi.org/10.1111/2041-210X.12055

493  Pie, M.R., Carrijo, T.F., Caron, F.S., 2021. The diversification of termites: Inferences from a

494     complete species-level phylogeny. Zool. Scr. 1–11. https://doi.org/10.1111/zsc.12502

495  Rabosky, D.L., Chang, J., Title, P.O., Cowman, P.F., Sallan, L., Friedman, M., Kaschner, K.,

496     Garilao, C., Near, T.J., Coll, M., Alfaro, M.E., 2018. An inverse latitudinal gradient in

497     speciation rate for marine fishes. Nature 559, 392–395. https://doi.org/10.1038/s41586-

498     018-0273-1

499  Redding, D.W., Mooers, A.O., 2006. Incorporating evolutionary measures into conservation

500     prioritization. Conserv. Biol. 20, 1670–1678. https://doi.org/10.1111/j.1523-

501     1739.2006.00555.x

502  Roa-Fuentes, C.A., Heino, J., Cianciaruso, M. V., Ferraz, S., Zeni, J.O., Casatti, L., 2019.

503     Taxonomic, functional, and phylogenetic β-diversity patterns of stream fish assemblages

504     in tropical agroecosystems. Freshw. Biol. 64, 447–460.

505     https://doi.org/10.1111/fwb.13233

506  Roa-Fuentes, C.A., Heino, J., Zeni, J.O., Ferraz, S., Cianciaruso, M.V., Casatti, L., 2020.

507     Importance of local and landscape variables on multiple facets of stream fish

508     biodiversity in a Neotropical agroecosystem. Hydrobiologia 7.

509     https://doi.org/10.1007/s10750-020-04396-7

510  Roquet, C., Thuiller, W., Lavergne, S., 2013. Building megaphylogenies for macroecology:

511     Taking up the challenge. Ecography (Cop.). 36, 13–26. https://doi.org/10.1111/j.1600-

512     0587.2012.07773.x

513  Seger, G.D.S., Duarte, L.D.S., Debastiani, V.J., Kindel, A., Jarenkow, J.A., 2013.

514     Discriminating the effects of phylogenetic hypothesis, tree resolution and clade age

515    estimates on phylogenetic signal measurements. Plant Biol. 15, 858–867.

516    https://doi.org/10.1111/j.1438-8677.2012.00699.x

517  Smith, S.A., Beaulieu, J.M., Donoghue, M.J., 2009. Mega-phylogeny approach for

518    comparative biology: An alternative to supertree and supermatrix approaches. BMC

519    Evol. Biol. 9, 1–12. https://doi.org/10.1186/1471-2148-9-37

520  Stein, R.W., Mull, C.G., Kuhn, T.S., Aschliman, N.C., Davidson, L.N.K., Joy, J.B., Smith,

521    G.J., Dulvy, N.K., Mooers, A.O., 2018. Global priorities for conserving the evolutionary

522    history of sharks, rays and chimaeras. Nat. Ecol. Evol. 2, 288–298.

523    https://doi.org/10.1038/s41559-017-0448-4

524  Tedesco, P.A., Beauchard, O., Bigorne, R., Blanchet, S., Buisson, L., Conti, L., Cornu, J.F.,

525    Dias, M.S., Grenouillet, G., Hugueny, B., Jézéquel, C., Leprieur, F., Brosse, S.,

526    Oberdorff, T., 2017. Data Descriptor: A global database on freshwater fish species

527    occurrence in drainage basins. Sci. Data 4, 1–6. https://doi.org/10.1038/sdata.2017.141

528  Tonini, J.F.R., Beard, K.H., Ferreira, R.B., Jetz, W., Pyron, R.A., 2016. Fully-sampled

529    phylogenies of squamates reveal evolutionary patterns in threat status. Biol. Conserv.

530    204, 23–31. https://doi.org/10.1016/j.biocon.2016.03.039

531  Trisos, C.H., Auerbach, J., Katti, M., 2021. Decoloniality and anti-oppressive practices for a

532    more ethical ecology. Nat. Ecol. Evol. https://doi.org/10.1038/s41559-021-01460-w

533  Upham, N.S., Esselstyn, J.A., Jetz, W., 2019. Inferring the mammal tree: Species-level sets of

534    phylogenies for questions in ecology, evolution, and conservation. PLOS Biol. 17,

535    e3000494. https://doi.org/10.1371/journal.pbio.3000494

536  Webb, C.O., Ackerly, D.D., Kembel, S.W., 2008. Phylocom: Software for the analysis of

537    phylogenetic community structure and trait evolution. Bioinformatics 24, 2098–2100.

538    https://doi.org/10.1093/bioinformatics/btn358

539  Webb, C.O., Ackerly, D.D., McPeek, M. a., Donoghue, M.J., 2002. Phylogenies and

540     Community Ecology. Annu. Rev. Ecol. Syst. 33, 475–505.

541     https://doi.org/10.1146/annurev.ecolsys.33.010802.150448

542  Webb, C.O., Donoghue, M.J., 2005. Phylomatic: Tree assembly for applied phylogenetics.

543     Mol. Ecol. Notes 5, 181–183. https://doi.org/10.1111/j.1471-8286.2004.00829.x

544