# Defining endemism levels for biodiversity conservation: tree species in the Atlantic Forest hotspot — Source link

Renato A. F. de Lima, Vinicius Castro Souza, Marinez Ferreira de Siqueira, Hans ter Steege

**Institutions:** Naturalis, Escola Superior de Agricultura Luiz de Queiroz

Related papers:

- Defining endemism levels for biodiversity conservation: Tree species in the Atlantic Forest hotspot

- Quantity versus quality: Endemism and protected areas in the temperate forest of South America

- Using endemism to assess representation of protected areas – the family Myrtaceae in the Greater Blue Mountains World Heritage Area

- Species richness and endemism in the native flora of California.

- Angiosperm biodiversity, endemism and conservation in the Neotropics

1 **Defining endemism levels for biodiversity conservation: tree species in**

2 **the Atlantic Forest hotspot**

3

4 **Abstract**

5 Endemic species are important for biodiversity conservation. Yet, quantifying

6 endemism remains challenging because endemism concepts can be too strict (i.e., pure

7 endemism) or too subjective (i.e., near endemism). We propose a data-driven approach

8 to objectively estimate the proportion of records inside a given the target area (i.e.,

9 endemism level) that optimizes the separation of near-endemics from non-endemic

10 species. We apply this approach to the Atlantic Forest tree flora using millions of

11 herbarium records retrieved from multiple sources. We first report an updated checklist

12 of 5044 species for the Atlantic Forest tree flora and then we compare how species-

13 specific endemism levels obtained from herbarium data match species-specific

14 endemism accepted by taxonomists. We show that an endemism level of 90% separates

15 well pure and near-endemic from non-endemic species, which in the Atlantic Forest

16 revealed an overall endemism ratio of 45% for its tree flora. We also found that the

17 diversity of pure and near endemics and of endemics and overall species was congruent

18 in space. Our results for the Atlantic Forest reinforce that pure and near endemic species

19 can be combined to quantify regional endemism and therefore to set conservation

20 priorities taking into account endemic species distribution. We provided general

21 guidelines on how the proposed approach can be used to assess endemism levels of

22 regional biotas in other parts of the world.

23

24 **Keywords**: biodiversity hotspot, endemism centers, endemism ratio, near endemism,

25 occasional species, plant conservation

## 1. Introduction

26

27 One common practice in biodiversity conservation is to focus on species with high

28 conservation value, such as species threatened with extinction (i.e., threatened species)

29 or those exclusive to a given region or habitat (i.e., endemic species). Threatened and

30 endemic species are important for conservation because they have a greater extinction

31 risk than other species (Brooks et al., 2006; Myers et al., 2000; Peterson and Watson,

32 1998). In addition, the spatial patterns of total and endemic species richness can be

33 congruent (Kier et al., 2009; Bonn et al., 2002; Storch et al., 2012), so prioritizing the

34 protection of areas with high-levels of endemism could also safeguard the remaining

35 biodiversity. However, there have been more efforts to delimit threatened species than

36 endemic ones. Threatened species are grouped by clearly-defined categories, enclosed

37 by objective criteria (IUCN, 2018), while species often are classified simply as being

38 endemic or not.

39      There are proposals to divide endemics species based on spatial scale (e.g.,

40 narrow, regional and continental endemics), evolutionary history (e.g., neo and paleo

41 endemics) or habitat specificity (e.g., edaphic endemics; Ferreira and Boldrini, 2011;

42 Kruckeberg and Rabinowitz, 1985; Peterson and Watson, 1998). These proposals,

43 however, implicitly assume that all individuals of a species are confined to a given

44 region or habitat, also known as true or pure endemism (Tyler, 1996). If one record is

45 found outside the target region, the species is to be (re)classified as non-endemic. Since

46 pure endemism is rather strict, the term near-endemism has been used to describe

47 species with few records outside the target region (Matthews et al., 1993; Carbutt and

48 Edwards 2006; Platts et al., 2011; Noroozi et al., 2018). Near-endemics are the result of

49 rare dispersal events, temporary establishment in different habitats or the existence

50 small satellite populations (Matthews et al., 1993; Perera et al., 2011). It is important to

51    emphasize that both types of endemism refer to species restricted to a specific area or

52    habitat, which does not necessarily imply species with small extent of occurrence

53    (<20,000 km$^2$ *sensu* IUCN, 2018) or low local abundance (Rabinowitz, 1981).

54         The differentiation between pure and near endemics is challenging, because it

55    may not be stable in time: near endemics can become pure endemics if habitat loss is

56    higher outside than inside the target region (Carbutt and Edwards, 2006). Conversely,

57    pure endemics may become near endemics with the accumulation of knowledge on their

58    geographical distribution (Werneck et al., 2011). This is particularly true for

59    geographically-restricted species, which often have scarce occurrence data.

60    Furthermore, pure endemics may be classified as near endemics due to species

61    misidentifications (Carbutt and Edwards, 2006) or by a questionable delimitation of the

62    target region (Platts et al., 2011). In practice, conservation aims at protecting as many

63    individuals as possible for a given species (IUCN, 2018). So, the differentiation

64    between pure and near endemism may have little impact to plan conservation actions.

65    Therefore, the question is: how to distinguish both groups of endemic species from non-

66    endemic species? Defining pure endemism is straightforward, but separating near-

67    endemics from non-endemic species can be quite subjective.

68         Here we propose a data-driven approach to objectively separate near-endemic

69    from non-endemic species for conservation purposes. This approach can also be used to

70    separate widespread species from occasional species, i.e., species frequent in other

71    regions but sporadic in a given target region (Barlow et al., 2010). Therefore, its main

72    goal is to classify species occurring inside a target region into pure-endemics, near-

73    endemics, widespread and occasional species, which is done based on their ratio of

74    occurrences inside the target region. As an example, we apply this approach to the

75    Atlantic Forest, a global biodiversity hotspot with abundant knowledge on the

76    taxonomy and distribution of its flora. We focus on the Atlantic Forest arborescent

77    flora, a plant growth form which is well represented in biological collections (Daru et

78    al., 2018). Using millions of carefully curated occurrences from over 500 collections

79    around the world, we evaluate which ratio of occurrences inside the Atlantic Forest

80    match species-specific endemism accepted by taxonomic experts. Finally, we illustrate

81    the implications of the proposed approach to assess endemism ratio, to support on-the-

82    ground conservation actions and to provide additional layer of information to existing

83    tools of spatial prioritization.

84

## 85    2. Material and methods

### 86    2.1 An objective approach to delimit species endemism

87    Here, we formalize the six steps of the proposed approach to objectively classify species

88    endemism levels, which can be applied in respect to any target region based on the

89    distribution of species occurrence records (Figure 1).

90

91    *2.1.1 Define the target area*

92    Species endemism cannot be assessed without defining a geographical area. This area

93    can be a region, domain or a habitat, but endemism is always relative and scale-

94    dependent (Laffan and Crisp, 2003). Although many countries may want to produce

95    their list of endemic species, it is recommended to use natural rather than political

96    boundaries to define the geographical extent of the target area (Ferreira and Boldrini,

97    2011). If the target area is facing changes, such as an increasing loss of natural habitats,

98    specifying the time window over which the target area is being considered may be

99    relevant.

100

101   *2.1.2 Define the target organism(s)*

102   Assessing the endemism of all living species is too time-consuming or data-limited for

103   many taxa. Therefore, one needs to restrict the assessment to one or fewer taxa, that

104   may be chosen according to their taxonomy (e.g., genus, family), or according to their

105   life form, ecology (e.g., ecological guild), function (e.g., trophic levels) or conservation

106   value (e.g., threatened species). Once the target organisms were defined, it is important

107   to build a comprehensive list of names for all species occurring inside the target area.

108   This list that should include synonyms and orthographical variants of the valid species

109   names, to increase changes of occurrence data retrieval. If this a list of names is not

110   available, a list of localities containing the target area (e.g., country names) can be used

111   to generate a list of organisms potentially occurring inside the target area.

112

113   *2.1.3 Obtain species occurrence data*

114   After defining the input list of names to search for species occurrences, it is necessary to

115   define the data sources, which can be primary sources (e.g., personal field collections),

116   secondary sources (e.g., biological collections, floras) or both. In the case of large

117   databases of secondary sources (e.g., GBIF) and/or large number of taxa, the number of

118   occurrences available may be large (thousands to millions). So, the use of automatized

119   tools for data download and documentation may be needed (Chamberlain et al., 2020).

120

121   *2.1.4 Validate occurrence data*

122   Particularly when using data from multiple secondary sources, it is important to validate

123   the information accompanying the occurrences, such as the collector name and number,

124   collection locality and geographical coordinates. In the case of two or more biological

125   collections, the removal of duplicated specimens across collections is advised.

126    Depending on the characteristics of the data sources, one may need to remove spatial

127    duplicates (e.g., records from the same localities) or spatial outliers (probable errors

128    placed too far away from species core distributions). Another important validation step

129    is to define the accepted confidence level of the taxonomic identification of each

130    occurrence (e.g., use only identifications performed by taxonomists).

131

132    *2.1.5 Calculate species endemism levels*

133    Next step simply is the count of the number of valid occurrences inside and outside the

134    target area(s). This can be done by aggregating occurrences by locality names or by

135    crossing a map of the target area with the geographical coordinates of the occurrences.

136    The simplest endemism level metric possible is the number of valid occurrences inside

137    the target area over the total number of occurrences retrieved for each taxon. If there is

138    uncertainty on the delimitation of the boundaries of the target area (e.g., low-resolution

139    map), the occurrences falling close to these boundaries may need a differential

140    treatment to avoid biases on species classifications due to imprecise boundary

141    delimitation (Platts et al., 2011).

142

143    *2.1.6 Classify species for conservation planning*

144    The empirical levels of species endemism calculated in the previous step can be used as

145    a metric of species endemicity in itself or as means to classify species into categories

146    according to their degree of endemicity. For instance, if all records occur inside the

147    target area the species can be classified as pure endemic and if the majority of the

148    records occur outside, the species can be classified as occasional. One needs to assume

149    (or estimate, see example below) thresholds of endemism level to separate near

150    endemics and occasional from other species. Ideally, the classification should be

151    compared to existing classifications of species endemism for validation.

152

153    **2.2 A case study: tree species in the Atlantic Forest**

154    We applied this approach to the arborescent flora of the Atlantic Forest biodiversity

155    hotspot in eastern South America (see Supplementary Material for full details).

156

157    *2.2.1 Target area and organisms.*

158    The Atlantic Forest originally covered ca. 136 million hectares in three different

159    countries, Argentina, Brazil and Paraguay (geographical range: 4–34º S latitude, 35–57º

160    W longitude – Figure S1a). Therefore, we searched for species occurrence data using a

161    list of species names occurring in South America, compiled from different sources

162    (Zuloaga et al., 2008; Oliveira-Filho, 2010; Grandtner and Chevrette 2013; Lima et al.,

163    2015; Zappi et al., 2015; ter Steege et al., 2016). Here we considered only a part of the

164    Atlantic Forest biota, the arborescent species, hereafter referred simply as trees. We

165    considered tree species occurrences in all Atlantic Forest types, which include

166    evergreen, semi-deciduous, deciduous, mixed temperate (locally known as *Araucaria*

167    forests), white-sand ('Restingas' and 'Mussunungas'), alluvial, cloud and swamp

168    forests, as well as in rocky field and inselberg vegetation. Arborescent species are

169    relatively well represented in herbaria (Daru et al., 2018) and they are defined here as

170    species with free-standing stems exceeding 5 cm of diameter at breast height (1.3 m) or

171    4 m in total height, including arborescent palms, cactus, tree ferns, and woody bamboos.

172    Moreover, some tall shrubs and treelets are included here under the term trees. We

173    carefully inspected the input list of names to avoid the inclusion of exotic and non-

174    arborescent species.

175

176 *2.2.2 Retrieval and validation of occurrence data.*

177 The list of South American tree names was used to download occurrence data from

178 multiple secondary sources, namely *species*Link (www.splink.org.br), JABOT

179 (http://jabot.jbrj.gov.br, Silva et al., 2017), 'Portal de Datos de Biodiversidad Argentina'

180 (https://datos.sndb.mincyt.gob.ar) and the Global Biodiversity Information Facility

181 (GBIF.org, 2019). We excluded all occurrences described in the specimen notes as

182 being cultivated or exotic. We checked names for typos, orthographical variants and

183 synonyms in the Brazilian Flora 2020 (BF-2020) project (Filardi et al., 2018; Zappi et

184 al., 2015). Decisions for unresolved names were made by consulting Tropicos

185 (www.tropicos.org) or the World Checklist of Selected Plant Families

186 (http://wcsp.science.kew.org).

187 There was much variation of the notation across herbaria, on the locality details

188 provided and on the precision of the geographical coordinates among the millions of

189 records retrieved (Appendix A). Therefore, we conducted a detailed data cleaning and

190 validation procedure (see Supplementary Material for details). We standardized the

191 notation of different fields (e.g., locality description, collector and identifier names,

192 collection and identification dates), which were then used to (i) search for duplicate

193 specimens among herbaria; (ii) validate the geographical coordinates at country, state

194 and/or county levels and (iii) to assess the confidence level of the identification of each

195 specimen (i.e., 'validated' and 'probably validated' identifications - Appendix B).

196 Moreover, (iv) we cross-validated information of duplicate specimens across herbaria to

197 obtain missing or more precise coordinates and/or valid specimen identifications.

198 Finally, (iv) we removed specimens too distant from their core distributions (i.e., spatial

199    outliers), which are often related to specimens collected from cultivated individuals but

200    that are not declared so by the collectors.

201

202    *2.2.3 Calculating species endemism levels.*

203    We calculated an empirical level of endemism based on the position of records for each

204    species in respect to the Atlantic Forest limits (Olson and Dinerstein, 2002; IBGE,

205    2012). Each record was assigned as being inside, outside or in the transition of the

206    Atlantic Forest to other domains (see details in Figure S1b). Records in the transition

207    were those falling inside the Atlantic Forest limits, but in counties with less than 90% of

208    its area inside the Atlantic Forest or vice-versa. Because of the variable precision of the

209    specimen's coordinates and of the uncertainty of the boundary delimitation at the scale

210    of our target area map (1:5,000,000), records in the transition received half the weight

211    other records to calculated species endemism levels:

212    $$100 \times \left( O_{in} + {O_{ti}}/{2} \right) / \left( O_{in} + {O_{ti}}/{2} + O_{out} + {O_{to}}/{2} \right),$$

213    where, $O_{in}$, $O_{ti}$, $O_{out}$ and $O_{to}$ are the number of specimens inside, inside in the transition,

214    outside and outside in the transition to the Atlantic Forest, respectively. This endemism

215    level is actually a weighted proportion of occurrences inside the Atlantic Forest by the

216    total of valid occurrences found, varying from 0 (no occurrences) to 100% (all

217    occurrences inside the Atlantic Forest).

218         We then obtained the endemism classification derived from the expertise of

219    taxonomists working on the BF-2020 project (Filardi et al., 2018), the best reference

220    currently available for the Atlantic Forest flora. Each species was classified as

221    'endemic' if the BF-2020 field 'phytogeographic domain' contained only the term

222    'Atlantic Rainforest' (equivalent to what we refer here as Atlantic Forest with all of its

223    forest types). Correspondingly, a species was classified as 'occasional' if this field did

224    not include this term. Species with no information on the 'phytogeographic domain'

225    were omitted from this analysis.

226    The comparison between the empirical classification of species endemism and

227    the reference BF-2020 classification was based on thresholds values varying from 0 to

228    100%, in intervals of 1% (i.e., 0, 1, …, 99, 100%). If a given species had an observed

229    endemism level equal or higher than a given threshold, it was classified as 'endemic'.

230    For each threshold value, we calculated the number of mismatches between the two

231    classifications (i.e., species classified as 'endemic' in the BF-2020 and 'not endemic'

232    from the observed endemism level or vice-versa). The same procedure was used to

233    calculate the number of mismatches for occasional species. We then plotted the number

234    of mismatches against all thresholds and estimated the optimum threshold that

235    minimizes the number of mismatches between classifications. Optimum thresholds were

236    estimated using piecewise regression, allowing up to five segments (i.e., four breaking

237    points). Thus, we provided the breaking point of each curve (and its 95% confidence

238    interval). We compared the results using only taxonomically 'validated' and using both

239    taxonomically 'validated' and 'probably validated' records.

240

241    *2.2.4 Species classification and implications for conservation planning*

242    We used the optimum threshold values obtained above to classify species into pure

243    endemics, near endemics, widespread and occasional species. Because endemic species

244    are not necessarily narrowly distributed and occasional species may be frequent

245    elsewhere, this terminology tried to reflect broad patterns of species occurrence in

246    respect to the target region (pure and near endemics: all or nearly all occurrences within

247    the target region; widespread: species with many occurrences both within and outside

248    the target area; occasional: species with most occurrences outside the target area). We

249 then used this classification to delimit the centers of diversity for each group of species

250 (Laffan and Crisp, 2003). In order to do so, we plotted the valid occurrences of each

251 group of species against a 50×50 km grid covering the Atlantic Forest and surrounding

252 domains. Next, we obtained different diversity metrics for each group of species per

253 grid cell. We selected two metrics with best performance to describe our data (Figures

254 S2 and S3): corrected weighted endemism (WE) and rarefied/extrapolated richness

255 ($S_{RE}$). The WE is the species richness weighted by the inverse of the number of cells

256 where the species is present, divided by cell richness (Crisp et al., 2001). The $S_{RE}$ is the

257 rarefied/extrapolated richness (depending on the observed number of occurrences per

258 cell) for a common number of 100 occurrences, calculated based on the species

259 frequencies per cell (Chao et al., 2014). We also obtained the sample coverage estimate

260 (Chao and Jost, 2012), used here as a proxy of sample completeness. We evaluated the

261 relationship of the diversity of endemic and occasional species with overall species

262 diversity using spatial regression models (i.e., linear regression with spatially correlated

263 errors - Pinheiro and Bates, 2000). Centers of diversity were delimited using ordinary

264 kriging and only the grid cells meeting some minimum criteria of sampling coverage

265 (see Supplementary Material). We used the 80% quantile of predicted distributions to

266 delimit the centers of endemism.

267

## 3. Results

269 The search for occurrence records based on this input list of tree names resulted in a

270 total of 3.11 million records from 543 collections (Appendix A). After the removal of

271 duplicates, spatial outliers and the geographical and taxonomic validation, we retained

272 593,920 valid records (disregarding records with 'probably validated' taxonomy) for the

273 classification of species endemism. We found 252,911 valid records being collected

274    inside the Atlantic Forest limits, which contained a total of 5044 arborescent species

275    (4054 species excluding tall shrubs; Appendix C). If we consider the valid occurrences

276    in the transitions of the Atlantic Forest to other domains, we could add 294 species as

277    probably occurring in the Atlantic Forest (Appendix D). Another 3158 names were

278    retrieved but were finally excluded from the list for different reasons (e.g., synonyms,

279    typos, orthographical variants, species not occurring naturally in the Atlantic Forest,

280    etc.; Appendix E).

281        Based on the valid records retrieved for the Atlantic Forest, we found evidence

282    of pure endemism (i.e., endemism level= 100%) for 1547 tree species (31%; Appendix

283    F). We found that 90.2% of records inside the Atlantic Forest (95% Confidence

284    Interval, CI: 89.3–91.2%) was the threshold of endemism level that best matched the

285    endemism currently accepted by taxonomy experts (Figure 2a). The curve of

286    mismatches between the observed and reference classifications decreases until it reaches

287    a minimum and then it increases again, meaning that more or less restrictive thresholds

288    lead to an increase the number of mismatches. The 90.2% threshold in the Atlantic

289    Forest added 733 near endemic species (15%). Together, pure and near endemics lead to

290    an overall endemism ratio of 45.2% for the Atlantic Forest arborescent flora (Figure 2b)

291    and 1.01 endemic arborescent species per 100 $km^2$ of remaining forest (i.e., 2261.2 $km^2$;

292    Fundación Vida Silvestre Argentina and WWF, 2017). Conversely, we found that 8.7%

293    (95% CI: 8.2–9.3%) was the best threshold for separating occasional from widespread

294    species occurring in the Atlantic Forest (Figure 2a), leading to a total of 639 occasional

295    species (13%). The remaining 42% of the species were classified as widespread

296    (Appendix F). Results using only occurrences with taxonomy flagged as 'validated'

297    were similar (pure endemism: 32%; near endemism: 15%; occasional species: 14%,

298    widespread species: 39% - Figure S4, Appendix F).

299    The diversity of endemic species was strongly correlated with the overall species

300    diversity in the Atlantic Forest (Figure 3). There was also a strong and positive

301    correlation between the number of pure and near endemic species (Figure S5), meaning

302    that the centers of diversity of pure and near endemics are highly congruent in space.

303    The diversity of pure endemics was higher in the rainforests along the coast (Figure 4),

304    corresponding to the rainforests of the Serra do Mar and Bahia Coastal Forests

305    ecoregions (Olson and Dinerstein, 2002). The inclusion of near endemics expanded the

306    diversity of endemic species towards more inland parts of the Atlantic Forest, but

307    spatial patterns remained quite similar (Figure 4 and Figures S6-S8). This expansion

308    was more conspicuous in the colder Araucaria forests in the southern Atlantic Forest,

309    but not to the point of including these forests as centers of diversity (i.e., areas with the

310    80% higher values). On the other hand, occasional species were really rare in the

311    Araucaria forests. Most of the distribution of occasional species was concentrated in the

312    Brazilian Cerrado, but also in the Amazon and slightly less in the Caatinga domain.

313    General patterns were fairly similar when using other diversity measures (Figures S6-

314    S8).

315

316    **4. Discussion**

317    **4.1 Describing species endemism**

318    Near endemism has been used to assess endemism levels of regional floras and faunas.

319    However, such assessments often use loose (Carbutt and Edwards, 2006; Platts et al.,

320    2011) or arbitrary definitions (Perera et al., 2011; Noroozi et al., 2018) of near

321    endemics. Here, we propose and apply an objective approach to find that 90% of the

322    occurrences inside a target region can be used to tell apart endemic species from non-

323    endemic species, a result supported by endemism classifications performed by

13

324      taxonomic experts. This 90% limit has one important implication: the average

325      endemism concept adopted by taxonomic experts implicitly includes the concept of near

326      endemism, at least for the Atlantic Forest. Indeed, the overall endemism ratio found

327      here for pure and near endemics combined (45%) is within the range of 40-50%

328      endemism level previously reported for the flora of this biodiversity hotspot (Myers et

329      al., 2000; Stehmann et al., 2009; Zappi et al., 2015). Thus, we propose that pure and

330      near endemics can be used together to objectively delimit endemism or as two

331      categories of endemism, similarly to what already exists for the categories of species

332      threat (IUCN, 2018). Moreover, conservation funding is not always aligned with the

333      degree of species endemism (Martín-López et al., 2009), despite the civic and scientific

334      awareness of the role of endemics for prioritizing conservation (Myers et al., 2000;

335      Brooks et al., 2006; Meuser et al., 2009; see Scarano, 2009 for a different point of

336      view). Thus, we hope that the quantitative description of endemism proposed here can

337      help to bridge the scarcity of conservation actions using information on species

338      endemicity.

339          The threshold of 90% found here was also used to assess plant endemism in the

340      Mediterranean Basin biodiversity hotspot (Médail and Baumel, 2018), suggesting that

341      this threshold could be used in the assessment of plant endemism of other species-rich

342      regions. However, we did not find similar assessments in the literature to confirm this

343      suggestion. Thus, although our approach to delimit species endemism is objective and

344      more comprehensive than pure endemism, similar assessments in other parts of the

345      world and for other groups of species are still needed. We provide a workflow to

346      perform such assessments, which would require (*i*) a list of species names, (*ii*) available

347      sources of occurrence data, (*iii*) a data cleaning/validation pipeline, (*iv*) a digitized map

348      of the study area, and (*v*) a classification of endemism based on taxonomists expertise.

349    Online occurrence data sources (e.g., GBIF) and tools to download data (e.g.,

350    Chamberlain et al., 2020) and validate their geographical coordinates (e.g., Zizka et al.,

351    2019) are becoming increasingly available. Here, we propose a simple but efficient way

352    to validate the taxonomic determinations of specimens (see Supplementary Material).

353    The bottleneck for applying this approach remains on the availability of regional lists of

354    species names and on the quantity and accessibility of data from local collections

355    (Boakes et al., 2010). These constraints may become more restrictive in species-rich and

356    less economically developed regions. The Atlantic Forest, used here as a testing ground

357    to our proposed approach, combines one of the largest number of species occurrences

358    available for the tropics (see details below), with one of the most completed national

359    floras (i.e., expert endemism information available – Brazilian Flora project) and

360    herbaria networks (e.g., *species*Link, JABOT).

361

362    **4.2 Implications for conservation**

363    The application of our approach to the tree flora of the Atlantic Forest offers insights on

364    how it can be used for supporting the conservation of local floras or faunas. The first

365    insight is related to the total number of species reported to a given region. The Atlantic

366    Forest is arguably the tropical forest with one of the largest botanical knowledge

367    available, with ca. 680,000 unique specimens of tree species, or 42 specimens per 100

368    $km^2$ – average collection density in the Amazon forest is below 10 per 100 $km^2$ (ter

369    Steege et al., 2016). Nevertheless, we over 700 new valid occurrences of tree species for

370    this biodiversity hotspot, an increase of 21% to the 3343 trees previously reported by

371    the Brazilian Flora 2020 project (Zappi et al., 2015). About 47% of these new records

372    were represented by occasional species, which correspond to 13% of the total richness

373    of the Atlantic Forest tree flora. This result confirms that occasional species, despite of

15

374   their infrequency, make an important contribution to overall biodiversity of regional

375   biotas (Barlow et al., 2010; ter Steege et al., 2019). But more importantly, 53% of the

376   new records correspond to widespread species and endemic species. An increase of 16%

377   in the total richness was also observed for the Espírito Santo state flora compared to the

378   reported in the Brazilian Flora (Dutra et al., 2015). The Brazilian Flora 2020 project is

379   permanently being improved and is of utmost importance for the understanding of the

380   Brazilian flora (Zappi et al., 2015; Filardi et al., 2018), the richest in the world (Ulloa et

381   al., 2017). Here, we provide products that can be readily integrated into the Brazilian

382   Flora project (e.g., more refined endemism filters), illustrating how data-driven

383   approaches as the one proposed here can help to refine the knowledge of regional floras,

384   even in regions with a great knowledge about its flora, promoting the accumulation of

385   critical knowledge to support biodiversity conservation.

386           Another possible application of the approach is the detection of centers of

387   endemic species diversity. In the Atlantic Forest example provided here, the centers

388   detected were congruent with previous proposals, which suggested areas of high

389   endemism in the moist and rain forests between the Brazilian states of São Paulo and

390   Rio de Janeiro and between Espírito Santo and Bahia states (Thomas et al., 1998;

391   Murray-Smith et al., 2009). However, our results provided evidence that the coastal

392   lowland forests in the states of Paraná and Santa Catarina (PR-SC) should also be

393   included as important centers of tree endemism for the Atlantic Forest. In accordance to

394   Murray-Smith et al. (2009), we found no strong support for the existence of an area of

395   endemism along the coastal and '*brejo de altitude*' forests in Paraíba, Pernambuco and

396   Alagoas states (Thomas et al., 1998), at least not at the spatial scale used here (50×50

397   km). The Atlantic Forests of northeast Brazil are closer or are surrounded by seasonally

398   dry vegetation (i.e., *Caatinga*) and they share many floristic elements with Amazon

16

399   forests (Santos et al., 2007), which could lead to the lower endemism levels found for

400   the species occurring in this part of the Atlantic Forest.

401         The provision of lists of species along with their degree of endemicity can

402   support the selection of species for conservation projects (Martín-López et al., 2009;

403   Meuser et al., 2009). These projects could be related to on-the-ground actions targeting

404   individual species (e.g., Martins, 2014 or www.saveourspecies.org) or to restoration

405   plans aiming at the maximization of biodiversity conservation outcomes while restoring

406   ecosystem services (Brancalion et al., 2018). Moreover, since range-restricted endemics

407   are probably also threatened, existing initiatives such as the Brazilian Alliance for

408   Extinction Zero (www.biodiversitas.org.br/baze) could incorporate the information on

409   degree of endemicity in their species selection methods. It is important to emphasize

410   that not only the degree of endemicity should be taken into account in the selection of

411   species for conservation projects. Widespread species may play important functional

412   roles in natural ecosystems, so they should be included in conservation projects as well

413   (Scarano, 2009).

414         The delimitation of centers of endemic diversity also has direct implications for

415   conservation planning. For instance, they can assist the identification of Important Plant

416   Areas (IPA), provided by the Target 5 of the Global Strategy for Plant Conservation

417   (www.cbd.int/gspc), or of Key Biodiversity Areas (KBA -

418   www.keybiodiversityareas.org). Although the delimitation of IPAs and KBAs predicts

419   the use of endemic species, their definition is mainly based on the presence of

420   threatened species. Also, IPAs are highly concentrated non-tropical regions of the

421   northern hemisphere (www.plantlifeipa.org). Our data driven approach, based on

422   careful data curation, proved to be efficient to identify areas of high endemicity in one

423   of the richest tropical floras of the world and could be used to expand the IPA and KBA

424    programs. In the specific case of the Atlantic Forest, which has less than 20% of its

425    original forest cover, conservation actions are urgently needed. When combined with

426    other layers of information (e.g., socio-economic), maps of endemic species diversity

427    can be used as an additional layer of biodiversity information in existing tools of spatial

428    prioritization (e.g., Brancalion et al., 2019; Strassburg et al., 2019), aiming to pinpoint

429    remaining natural areas that should be protected or degraded lands that could be

430    prioritized in restoration actions. This suggestion is reinforced by the spatial congruence

431    found between the diversity of endemic and non-endemic tree species, meaning that

432    conservation of areas with high-levels of endemism could also safeguard a great deal of

433    the remaining Atlantic Forest tree flora (Kier et al., 2009; Bonn et al., 2002). Thus,

434    considering that defining threatened and endemic species have the same constraints

435    related to data availability and to the time and spatial scale considered (Ferreira and

436    Boldrini, 2011), the detection of endemics is more straightforward than threatened

437    species, which could speed up the decision-making process for conservation in rich

438    tropical biotas around the world.

439

440    **Data Availability**

441    All data providers and their citations are given in Appendix A. GBIF data used in the

442    analysis is also provided in the references.

443

444    **CRediT authorship contribution statement**

445    Renato A. F. de Lima: Conceptualization, Methodology, Formal analysis, Data curation,

446    Funding acquisition, Writing - original draft. Vinicius C. Souza: Validation, Data

447    curation, Writing - review & editing. Marinez F. Siqueira: Methodology, Writing -

448 review & editing. Hans ter Steege: Methodology, Funding acquisition, Writing - review

449 & editing.

450

**Declaration of competing interest**

452 The authors declare that they have no known competing financial interests or personal

453 relationships that could have appeared to influence the work reported in this paper.

454

**Acknowledgments**

461

**References**

463 Barlow, J., Gardner, T.A., Louzada, J., Peres, C.A., 2010. Measuring the conservation

464     value of tropical primary forests: the effect of occasional species on estimates of

465     biodiversity uniqueness. PLoS One 5, e9609.

466 Boakes, E.H., McGowan, P.J.K., Fuller, R.A., Chang-Qing, D., Clark, N.E., O'Connor,

467     K., Mace, G.M., 2010. Distorted views of biodiversity: spatial and temporal bias in

468     species occurrence data. PLoS Biol. 8, e1000385.

469 Bonn, A., Rodrigues, A.S.L., Gaston, K.J., 2002. Threatened and endemic species: are

470     they good indicators of patterns of biodiversity on a national scale? Ecol. Lett. 5,

471     733–741.

472 Brancalion, P.H.S., Bello, C., Chazdon, R.L., Galetti, M., Jordano, P., Lima, R.A.F.,

473    Medina, A., Pizo, M.A., Reid, J.L., 2018. Maximizing biodiversity conservation

474    and carbon stocking in restored tropical forests. Conserv. Lett. e12454.

475    Brancalion, P.H.S., Niamir, A., Broadbent, E., Crouzeilles, R., Barros, F.S.M.,

476    Zambrano, A.M.A., Baccini, A., Aronson, J., Goetz, S., Leighton Reid, J.,

477    Strassburg, B.B.N., Wilson, S., Chazdon, R.L., 2019. Global restoration

478    opportunities in tropical rainforest landscapes. Sci. Adv. 5, 1–12.

479    Brooks, T.M., Mittermeier, R.A., Da Fonseca, G.A.B., Gerlach, J., Hoffmann, M.,

480    Lamoreux, J.F., Mittermeier, C.G., Pilgrim, J.D., Rodrigues, A.S.L., 2006. Global

481    biodiversity conservation priorities. Science 313, 58–61.

482    Carbutt, C., Edwards, T.J., 2006. The endemic and near-endemic angiosperms of the

483    Drakensberg Alpine Centre. South African J. Bot. 72, 105–132.

484    Chamberlain, S., Barve, V., Mcglinn, D., Oldoni, D., Desmet, P., Geffert, L., Ram, K.,

485    2020. rgbif: Interface to the Global Biodiversity Information Facility API.

486    https://cran.r-project.org/package=rgbif.

487    Chao, A., Jost, L., 2012. Coverage-based rarefaction and extrapolation : standardizing

488    samples by completeness rather than size. Ecology 93, 2533–2547.

489    Chao, A., Gotelli, N.J., Hsieh, T.C., Sander, E.L., Ma, K.H., Colwell, R.K., Ellison,

490    A.M., 2014. Rarefaction and extrapolation with Hill numbers: A framework for

491    sampling and estimation in species diversity studies. Ecol. Monogr. 84, 45–67.

492    Crisp, M.D., Laffan, S., Linder, H.P., Monro, A., 2001. Endemism in the Australian

493    flora. J. Biogeogr. 28, 183–198.

494    Daru, B.H., Park, D.S., Primack, R.B., Willis, C.G., Barrington, D.S., Whitfeld, T.J.S.,

495    Seidler, T.G., Sweeney, P.W., Foster, D.R., Ellison, A.M., Davis, C.C., 2018.

496    Widespread sampling biases in herbaria revealed from large-scale digitization.

497    New Phytol. 217, 939–955.

498    Dutra, V.F., Alves-Araújo, A., Carrijo, T.T., 2015. Angiosperm checklist of Espírito

499        Santo: Using electronic tools to improve the knowledge of an Atlantic Forest

500        biodiversity hotspot. Rodriguésia 66, 1145–1152.

501    Ferreira, P.M.A., Boldrini, I.I., 2011. Potential Reflection of Distinct Ecological Units

502        in Plant Endemism Categories. Conserv. Biol. 25, 672–679.

503    Filardi, F.L.R., De Barros, F., Baumgratz, J.F.A., Bicudo, C.E.M., Cavalcanti, T.B.,

504        Nadruz Coelho, M.A., Costa, A.F., Costa, D.P., Goldenberg, R., Labiak, P.H.,

505        Lanna, J.M., Leitman, P., Lohmann, L.G., Costa Maia, L., Mansano, V.F., Morim,

506        M.P., Peralta, D.F., Pirani, J.R., Prado, J., Roque, N., Secco, R.S., Stehmann, J.R.,

507        Sylvestre, L.S., Viana, P.L., Walter, B.M.T., Zimbrão, G., Forzza, R.C. et al.,

508        2018. Brazilian Flora 2020: Innovation and collaboration to meet Target 1 of the

509        Global Strategy for Plant Conservation (GSPC). Rodriguésia 69, 1513–1527.

510    GBIF.org. 2019. GBIF Occurrence Download. https://doi.org/10.15468/dl.mzmat2.

511    Grandtner, M.M., Chevrette, J., 2013. Dictionary of Trees, Volume 2: South America:

512        Nomenclature, Taxonomy and Ecology. Academic Press, Amsterdam.

513    IBGE, 2012. Mapa da Área de Aplicação da Lei no 11.428 de 2006. Brasília, Brazil.

514    IUCN, 2018. The IUCN Red List of Threatened Species. http://www.iucnredlist.org.

515    Kier, G., Kreft, H., Lee, T.M., Jetz, W., Ibisch, P.L., Nowicki, C., Mutke, J., Barthlott,

516        W., 2009. A global assessment of endemism and species richness across island and

517        mainland regions. Proc. Natl. Acad. Sci. 106, 9322–9327.

518    Kruckeberg, A.R., Rabinowitz, D., 1985. Biological aspects of endemism in higher

519        plants. Annu. Rev. Ecol. Syst. 16, 447–479.

520    Laffan, S.W., Crisp, M.D., 2003. Assessing endemism at multiple spatial scales, with an

521        example from the Australian vascular flora. J. Biogeogr. 30, 511–520.

522    Lima, R.A.F. de, Mori, D.P., Pitta, G., Melito, M.O., Bello, C., Magnago, L.F.,

523    Zwiener, V.P., Saraiva, D.D., Marques, M.C., Oliveira, A.A. de, Prado, P.I., 2015.

524    How much do we know about the endangered Atlantic Forest? Reviewing nearly

525    70 years of information on tree community surveys. Biodivers. Conserv. 24, 2135–

526    2148.

527    Martín-López, B., Montes, C., Ramírez, L., Benayas, J., 2009. What drives policy

528    decision-making related to species conservation? Biol. Conserv. 142, 1370–1380.

529    Martins, E.M., Fernandes, F.M., Maurenza, D., Pougy, N., Loyola, R., Martinelli, G.,

530    2014. Plano de ação nacional para a conservação do Faveiro-de-wilson

531    (*Dimorphandra wilsonii* Rizzini). Jardim Botânico do Rio de Janeiro, Rio de

532    Janeiro, 52 p.

533    Matthews, W.S., van Wyk, A.E., Bredenkamp, G.J., 1993. Endemic flora of the north-

534    eastern Transvaal Escarpment, South Africa. Biol. Conserv. 63, 83–94.

535    Médail, F., Baumel, A., 2018. Using phylogeography to define conservation priorities:

536    The case of narrow endemic plants in the Mediterranean Basin hotspot. Biol.

537    Conserv. 224, 258–266.

538    Meuser, E., Harshaw, H.W., Mooers, A.Ø., 2009. Public preference for endemism over

539    other conservation-related species attributes. Conserv. Biol. 23, 1041–1046.

540    Murray-Smith, C., Brummitt, N.A., Oliveira-Filho, A.T., Bachman, S., Moat, J.,

541    Lughadha, E.M.N., Lucas, E.J., 2009. Plant diversity hotspots in the Atlantic

542    Coastal Forests of Brazil. Conserv. Biol. 23, 151–163.

543    Myers, N., Mittermeier, R.A., Mittermeier, C.G., Fonseca, G.A.B., Kent, J., 2000.

544    Biodiversity hotspots for conservation priorities. Nature 403, 858–863.

545    Noroozi, J., Talebi, A., Doostmohammadi, M., Rumpf, S.B., Linder, H.P., Schneeweiss,

546    G.M., 2018. Hotspots within a global biodiversity hotspot-areas of endemism are

547    associated with high mountain ranges. Sci. Rep. 8, 1–10.

548    Oliveira-Filho, A.T., 2010. TreeAtlan 2.0, Flora arbórea da América do Sul cisandina

549        tropical e subtropical: Um banco de dados envolvendo biogeografia, diversidade e

550        conservação. Universidade Federal de Minas Gerais. www. cb.ufmg.br/treeatlan

551    Olson, D.M., Dinerstein, E., 2002. The Global 200: Priority ecoregions for global

552        conservation. Ann. Missouri Bot. Gard. 89, 199–224.

553    Perera, S.J., Ratnayake-Perera, D., Procheş, Ş., 2011. Vertebrate distributions indicate a

554        greater Maputaland-Pondoland-Albany region of endemism. S. Afr. J. Sci. 107,

555        68–71.

556    Peterson, A.T., Watson, D.M., 1998. Problems with areal definitions of endemism: the

557        effects of spatial scaling. Divers. Distrib. 4, 189–194.

558    Pinheiro, J., Bates, D., 2000. Linear mixed-effects models: basic concepts and

559        examples. Mix. Model. S S-Plus.

560    Platts, P.J., Burgess, N.D., Gereau, R.E., Lovett, J.C., Marshall, A.R., McClean, C.J.,

561        Pellikka, P.K.E., Swetnam, R.D., Marchant, R., 2011. Delimiting tropical

562        mountain ecoregions for conservation. Environ. Conserv. 38, 312–324.

563    Rabinowitz, D., 1981. Seven forms of rarity, in: The biological aspects of rare plant

564        conservation. John Wiley and Sons, Chichester, pp. 205–217.

565    Santos, A.M.M., Cavalcanti, D.R., 2007. Biogeographical relationships among tropical

566        forests in north-eastern Brazil. J. Biogeogr 34, 437–446.

567    Scarano, F.R., 2009. Plant communities at the periphery of the Atlantic rain forest:

568        Rare-species bias and its risks for conservation. Biol. Conserv. 142, 1201–1208.

569    Silva, L.A.E. da, Fraga, C.N. de, Almeida, T.M.H. de, Gonzalez, M., Lima, R.O.,

570        Rocha, M.S. da, Bellon, E., Ribeiro, R. da S., Oliveira, F.A. de, Clemente, L. da S.,

571        Magdalena, U.R., Medeiros, E. von S., Forzza, R.C., 2017. Jabot - Sistema de

572        Gerenciamento de Coleções Botânicas: a experiência de uma década de

573        desenvolvimento e avanços. Rodriguésia 68, 391–410.

574    ter Steege, H., Vaessen, R.W., Cárdenas-López, D., Sabatier, D., Antonelli, A., de

575        Oliveira, S.M., Pitman, N.C.A., Jørgensen, P.M., Salomão, R.P., 2016. The

576        discovery of the Amazonian tree flora with an updated checklist of all known tree

577        taxa. Sci. Rep. 6, 29549.

578    ter Steege, H., Oliveira, S.M., Pitman, N.C.A., Sabatier, D., Antonelli, A., Andino,

579        J.E.G., Aymard, G.A., Salomão, R.P., 2019. Towards a dynamic list of Amazonian

580        tree species. Sci. Rep. 9, 1–5.

581    Stehmann, J.R., Forzza, R.C., Salino, A., Sobral, M., Pinheiro, D., Kamino, H.Y., 2009.

582        Plantas da Floresta Atlântica. Jardim Botânico do Rio de Janeiro, Rio de Janeiro.

583    Storch, D., Keil, P., Jetz, W., 2012. Universal species–area and endemics–area

584        relationships at continental scales. Nature 488, 78–81.

585    Strassburg, B.B.N., Beyer, H.L., Crouzeilles, R., Iribarrem, A., Barros, F., de Siqueira,

586        M.F., Sánchez-Tapia, A., Balmford, A., Sansevero, J.B.B., Brancalion, P.H.S.,

587        Broadbent, E.N., Chazdon, R.L., Oliveira Filho, A., Gardner, T.A., Gordon, A.,

588        Latawiec, A., Loyola, R., Metzger, J.P., Mills, M., Possingham, H.P., Rodrigues,

589        R.R., Scaramuzza, C.A. de M., Scarano, F.R., Tambosi, L., Uriarte, M., 2019.

590        Strategic approaches to restoring ecosystems can triple conservation gains and

591        halve costs. Nat. Ecol. Evol. 3, 62–70.

592    Thomas, W.W., Carvalho, A.M.V. de, Amorim, A.M.A., Garrison, J. & Arbeláez, A.L.,

593        1998. Plant endemism in two forests in southern Bahia, Brazil. Biodivers. Conserv.

594        7, 311–322.

595    Tyler, P.A., 1996. Endemism in freshwater algae, in: Biogeography of Freshwater

596        Algae. Springer Netherlands, Dordrecht, pp. 127–135.

597    Ulloa, C.U., Acevedo-Rodríguez, P., Beck, S., Belgrano, M.J., Bernal, R., Berry, P.E.,

598    Brako, L., Celis, M., Davidse, G., Forzza, R.C., Gradstein, S.R., Hokche, O., León,

599    B., León-Yánez, S., Magill, R.E., Neill, D.A., Nee, M., Raven, P.H., Stimmel, H.,

600    Strong, M.T., Villaseñor, J.L., Zarucchi, J.L., Zuloaga, F.O., Jørgensen, P.M.,

601    2017. An integrated assessment of the vascular plant species of the Americas.

602    Science 358, 1614–1617.

603    Werneck, M. de S., Sobral, M.E.G., Rocha, C.T.V., Landau, E.C., Stehmann, J.R.,

604    2011. Distribution and endemism of angiosperms in the Atlantic Forest. Nat.

605    Conserv. 9, 188–193.

606    Zappi, D.C., Ranzato Filardi, F.L., Leitman, P., Souza, V.C., Walter, B.M.T., Pirani,

607    J.R., Morim, M.P., Queiroz, L.P., Cavalcanti, T.B., Mansano, V.F., Forzza, R.C.,

608    2015. Growing knowledge: An overview of Seed Plant diversity in Brazil.

609    Rodriguésia 66, 1085–1113.

610    Zizka, A., Silvestro, D., Andermann, T., Azevedo, J., Duarte Ritter, C., Edler, D.,

611    Farooq, H., Herdean, A., Ariza, M., Scharn, R., Svantesson, S., Wengström, N.,

612    Zizka, V., Antonelli, A., 2019. CoordinateCleaner: Standardized cleaning of

613    occurrence records from biological collection databases. Methods Ecol. Evol. 10,

614    744–751.

615    Zuloaga, F., Morrone, O., Belgrano, M., 2008. Catálogo de las plantas vasculares del

616    cono sur (Argentina, southern Brazil, Chile, Paraguay y Uruguay). Monographs in

617    Systematic Botany from the Missouri Botanical Garden 107.

25

## Appendices

618

**Appendix A**: List of collections and data providers used for data compilation.

The numbers of records retrieved per collection correspond to overall sum of records

before data validation, thus including both valid and invalid records.

**Appendix B**: List of names of taxonomists per family used for taxonomical validation.

The 'tdwg.name' represents the taxonomist name following the standard notation of the

Biodiversity Information Standards (https://www.tdwg.org), which includes different

variants of notation found for the same taxonomist name.

**Appendix C**: Updated, taxonomically vetted checklist of the Atlantic Forest tree flora.

For each name included in the checklist we provide the life form, the status of the name

in respect to the Brazilian Flora 2020 project, the number of records found inside the

Atlantic Forest (both 'validated' and 'probably validated' taxonomy) and a list of up to

30 vouchers (only specimens with 'validated' taxonomy), giving priority to type

specimens. We also indicate which species were regarded as being taxa of low

taxonomic complexity (TBC) or taxa commonly cultivated outside its original range.

**Appendix D**: List of species with probable occurrence in the Atlantic Forest.

We present all names with valid records found only in the transition of the Atlantic

Forest to other domains and those names cited in the Brazilian Flora 2020 project as

being an Atlantic Forest species, but for which we did not find any valid records. Again,

we present for each name the life form, the number of records found and a list of up to

30 vouchers.

643    **Appendix E**: List of names excluded from the final Atlantic Forest checklist.

644    For each name on the list we provide the life form and the reason why the name was

645    excluded. For synonyms, orthographical variants, common typos we also provide the

646    corresponding valid name used in this study.

647

648    **Appendix F**: Endemism levels for the Atlantic Forest tree flora and the corresponding

649    classification into pure endemic, near endemic, widespread and occasional species.

650    For each species name, we provide the number of valid records outside the Atlantic

651    Forest, outside but in the transition to the Atlantic Forest, inside the Atlantic Forest but

652    in the transition to other domains, and inside the Atlantic Forest. We present the

653    endemism levels and species classifications using only records with validated taxonomy

654    and using records with validated and probably validated taxonomy. Finally, we present

655    the endemism classification currently accepted in the Brazilian Flora 2020 in respect to

656    the Atlantic Forest.

657

658    **Appendix G**: Shapefiles delimiting the centers of the endemic and occasional species

659    diversity in the Atlantic Forest for pure endemics, near endemics, pure + near endemics

660    and occasional species.

661    Each shapefile contains the isoclines corresponding to the 75%, 80%, 85%, 90% and

662    95% quantiles of the distribution of rarefied/extrapolated richness for 100 specimens,
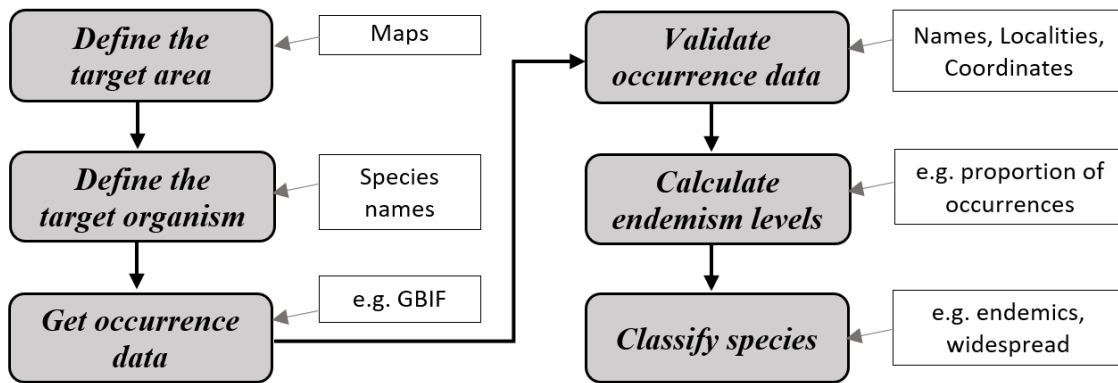
663    predicted using ordinary kriging.

# Figures



**Figure 1**. Flow chart showing the six steps of the proposed approach to classify species based on their endemism levels.
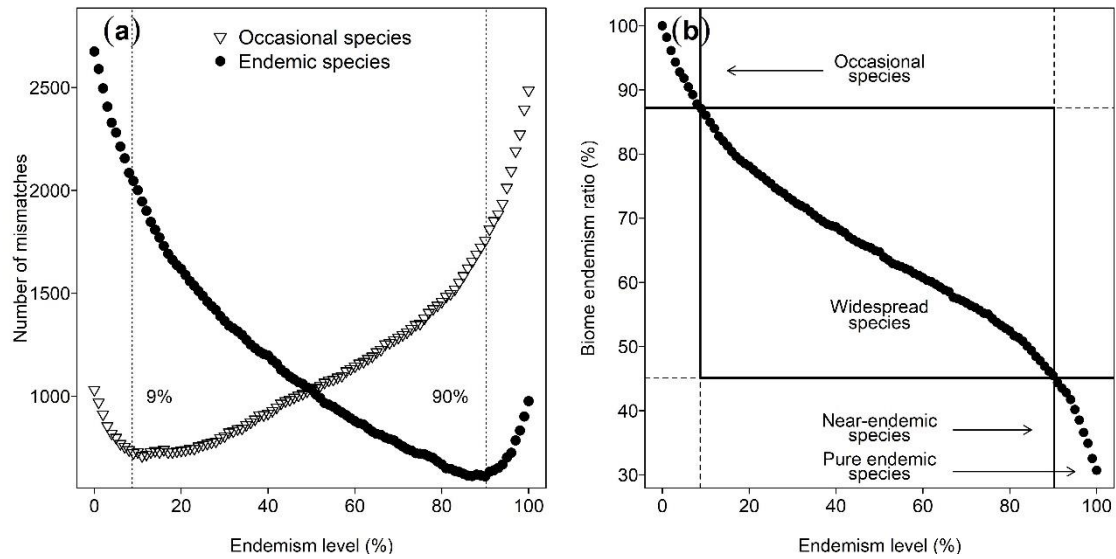
664

**Figure 2**. Defining near endemic and occasional tree species using herbarium records for the Atlantic Forest biodiversity hotspot. For both endemic (black circles) and occasional species (triangles), we present (a) the optimum endemism levels (vertical dashed lines) estimated from the distribution of mismatches between the empirical and the Brazilian Flora 2020 classifications and (b) the overall endemism ratio of the Atlantic Forest in intervals of 1% (*x*-axis in both panels).
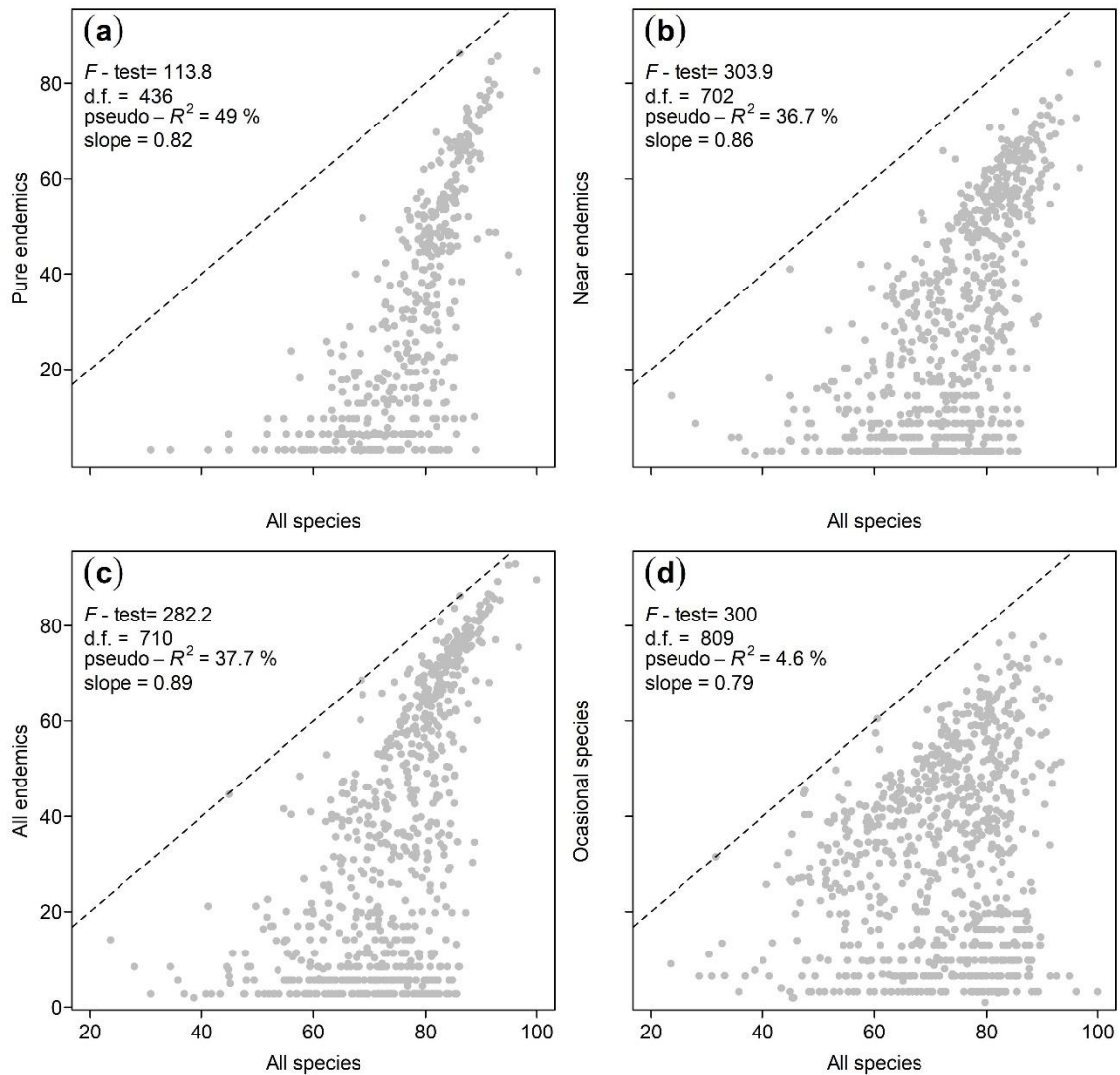
**Figure 3**. Relationship between the number of rarefied/extrapolated richness per 50×50 km grid cell and the same diversity metric obtained for (a) pure endemics, (b) near endemics, (c) all endemics (pure + near endemics) and (d) occasional species. For each group of species, we present the summary statistics of each spatial regression model (top left; d.f.= degrees of freedom), including the predicted slope of the regression prediction. The spatial regression analysis was performed only for grid cells meeting some minimum criteria of sampling coverage (see Supplementary Methods). The dashed line represents the 1:1 line. All *p*-values are below 0.001.
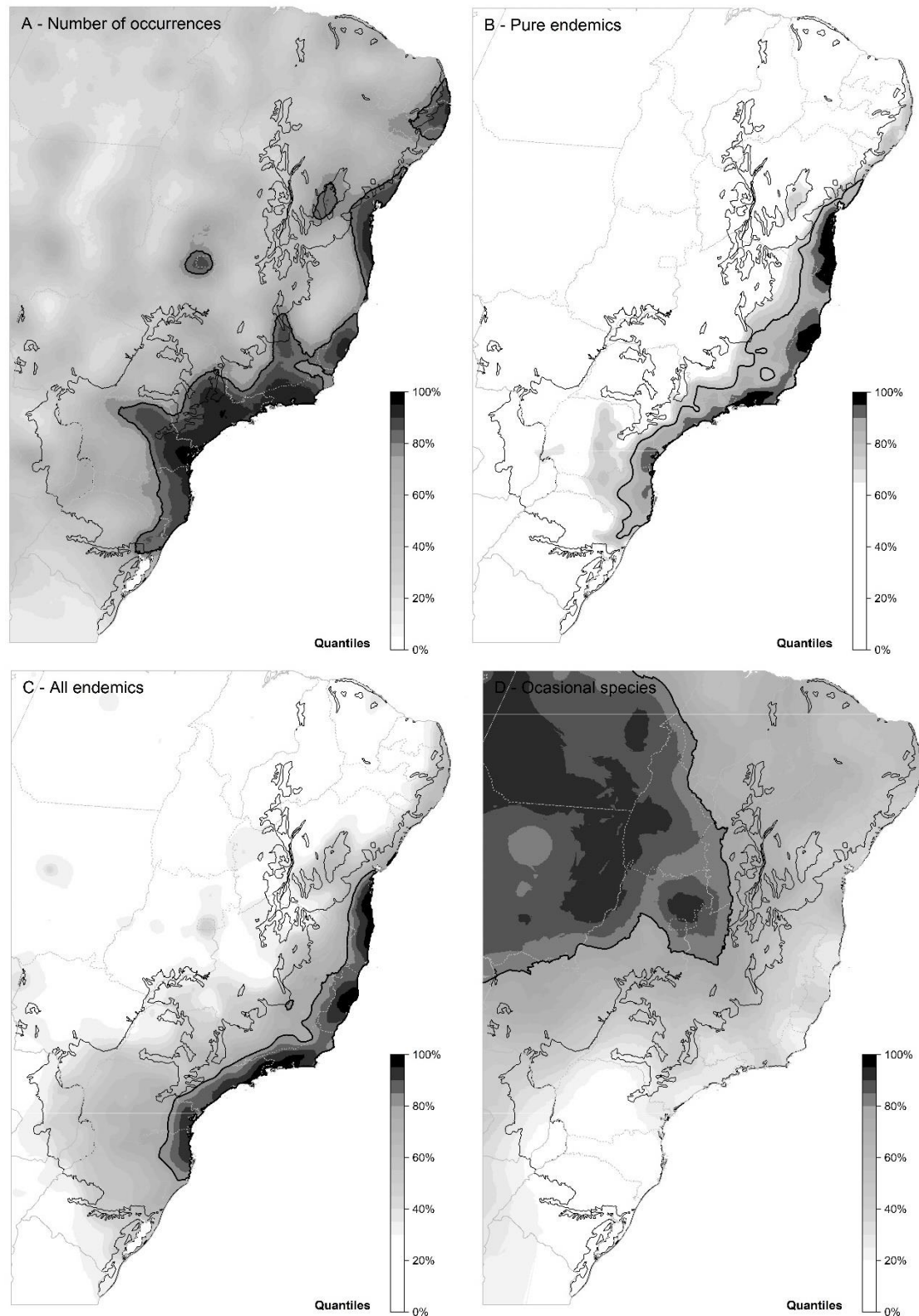
**Figure 4**. The spatial distribution of (A) the number of occurrences retrieved for the species occurring in the Atlantic Forest, and the centers of diversity of (B) pure endemics, (C) all endemics (pure + near) and (D) occasional species. Maps were

produced using ordinary kriging based on rarefied/extrapolated species richness obtained for a common number of 100 records per grid cell. The color scale represents the 5% quantiles of the metrics distribution, from 0-5% (white) to 95-100% (black). Bold black lines are the area containing the 80% higher richness values. The black line marks the limits of the Atlantic Forest, while the solid and dashed grey lines mark the limits of South American countries and of the Brazilian states, respectively.