1    **HaploVectors: an integrative analytical tool for phylogeography**

2    *Leandro Duarte[1][*], Jacqueline de Souza Lima[1], Renan Maestri[1], Vanderlei Debastiani[1] &*

3    *Rosane Garcia Collevatti[2]*

4    *[1] Departamento de Ecologia, Universidade Federal do Rio Grande do Sul*

5    *Av. Bento Gonçalves 9500 CP 15007, Porto Alegre, RS, 91501-970, Brazil*

6    *[2] Laboratório de Genética & Biodiversidade, Instituto de Ciências Biológicas, Universidade*

7    *Federal de Goiás. 74690-900, Goiânia, GO, Brazil*

8    [*] E-mail corresponding author: duarte.ldas@gmail.com

9    **Running head:** Haplotypic eigenvectors for phylogeography

10    **Abstract**

11    Phylogeographic approaches are commonly used to understand historical-biogeographic

12    patterns in the distribution of haplotypes. However, the emphasis of most tools lies on

13    describing spatial patterns of genetic variation and assess how large are haplotypic differences

14    among populations. An evaluation of the relative influence of environmental factors

15    compared to pure neutral process of haplotypic distribution - a question of great interest for

16    molecular ecologists - is less investigated, in part because appropriate tools are lacking. Here,

17    we introduce HaploVectors, a flexible tool that allows exploring phylogeographical patterns

18    and discriminating biogeographic, neutral and environmental factors acting to shape genetic

19    distribution across space. Haplovectors are variables that summarize the major gradients of

20    haplotypic distribution across a set of localities and allow weighting haplotypic frequencies

21    by the number of mutational steps using a fuzzy weighting approach. HaploVectors is

22    presented as an R package for computing haplotypic eigenvectors and performing null model-

23    based tests. Investigation of HaploVectors using empirical datasets showed that the method is

24    useful to uncover hidden patterns of haplotypic distribution, not easily detected using

25    traditional methods. Using a plant species as study case, we demonstrate by means of

26    HaploVectors that, even though the distribution of plant haplotypes was associated with

27    different biogeographic regions of the Brazilian Cerrado biome, such association was not

28    mediated by evolutionary relationships among haplotypes. The applicability of HaploVectors

29    is broad, ranging from the pure pattern exploration and discrimination of genetic populations,

30    to a hypothesis-testing framework that uses null-models to understand the influence of

31    environmental factors on haplotypic distribution.

34

**Introduction**

A complete understanding of the historical and biogeographic patterns of species distribution

benefits from the connection between micro and macroevolution, a major goal of the field of

phylogeography (Avise, 1987; Avise, 2009). Since the middle 1990's, the number of studies

using molecular markers to understand phylogeographical patterns is increasing at astonishing

rates (Beheregaray, 2008; Turchetto-Zolet, Pinheiro, Salgueiro & Palma-Silva, 2013).

Accordingly, the number of molecular markers used increased from single locus to multiple

genome regions (Freeland, 2014; Blom, Horner & Moritz, 2016). Testing hypotheses about

ecological influences on the genealogical history of populations from a single species, and the

comparison of such patterns across multiple species (comparative phylogeography), has the

potential to shed light on process of species diversification and on the geological and

biogeographic connections of entire regions (Diniz-Filho et al., 2008; Carnaval et al., 2014).

However, hypothesis testing in a rigorous statistical framework were only latter incorporated

into phylogeography, with the advent of statistical phylogeographical approaches (Templeton

et al., 1998; Papadopoulou & Knowles, 2016) that helped moving the field beyond the

essential descriptive nature present in its infancy. Still, the development of analytical tools to

integrate the ecological thinking into phylogeography has lagged behind the ever-increasing

number of molecular loci discovered and the numerous tools focused on spatial genetics (e.g.,

Templeton, 2004; Miller, 2005; Epperson, 2005). In the era of multiple molecular markers

and genomics in phylogeography, new analytical tools are imperative to better understand the

increasingly complex phylogeographical patterns, to compare results from different loci, and

to uncover environmental correlates of genetic distribution in the 'twilight zone' (Diniz-Filho

et al., 2008).

58    A principal goal in studies of molecular ecology and phylogeography is to understand

59    to what extent genetic variation of a species, expressed as the spatial distribution of alleles,

60    genotypes and haplotypes, might be affected by different factors, such as dispersal

61    mechanism, history of populations, climate changes or just by isolation by distance (Manel,

62    Schwartz, Luikart & Tarbelet, 2003, Storfer et al., 2007, Avise, 2009). To explore processes

63    underlying genetic diversity distribution, a representation of evolutionary relationships among

64    organisms or populations by means of phylogenetic trees or networks is often performed. In

65    this sense, networks are recognized to be more appropriate than trees, since the former allow

66    visualization of reticulation events, such as hybridization and recombination (Kong, Sánchez-

67    Pacheco & Murphy, 2015). Nonetheless, a network representation per se does not allow

68    robust hypothesis tests to evaluate the interplay between genetic variation across space and

69    alternative explanatory factors. Analysis of molecular variance (AMOVA, Excoffier, Smouse

70    & Quattro, 1992) has been widely used for this purpose (Fitzpatrick, 2009; Maestri et al.,

71    2016; Raffini et al., 2018). AMOVA, which is a permutation procedure akin to approaches

72    often used in ecological studies, such as PERMANOVA (Pillar & Orlóci, 1996; Anderson,

73    2001) allows analyzing whether pairwise genetic dissimilarities between individuals

74    distributed across different sites or regions defined by groups of sites is higher than expected

75    by chance given within site (or region) dissimilarities. Note that AMOVA is based on overall

76    genetic dissimilarities, and therefore does not permit direct inference about effects of spatial,

77    environmental and/or biogeographic factors on the distribution of haplotypes across space.

78    Despite debates about different methods to construct haplotypes networks (see

79    Mardulyn et al., 2012), phylogenetic median-joining network (MJN, Bandelt, Forster & Röhl,

80    1999) is the method normally used in phylogeographical studies. The MJN method is based in

81    the minimum spanning trees and shows the relationships between haplotypes obtained by the

82    distance measured among them (e.g. number of character differences - Hamming distance).

83    The use of MJN has been grown exponentially since its development (Kong et al., 2015).

84    Moreover, the network representation used to explore evolutionary relatedness among

85    haplotypes do not allow neither a clear visualization of haplotype co-occurrence within sites

86    nor general trends in haplotype distribution across space.

87         Phylogenetic eigenvectors have been used to express the variation of phylogenetic

88    beta diversity (or simply phylobetadiversity) among an array of localities (Duarte, 2011;

89    Duarte, Debastiani, Freitas & Pillar, 2016) based on phylogenetic fuzzy weighting (Pillar &

90    Duarte, 2010). Phylogenetic fuzzy weighting allows describing sites by their phylogenetic

91    composition based on species incidences/abundances and the phylogenetic relatedness among

92    species (see also Duarte et al., 2016). The phylogenetic composition of a set of sites can be

93    thereby decomposed into independent phylogenetic eigenvectors using Principal Coordinates

94    of Phylogenetic Structure, or simply PCPS (Duarte, 2011; Duarte et al., 2016). Each

95    eigenvector describes the sites by scores that position them along a phylogenetic gradient

96    expressing a fraction of the total phylogenetic relatedness among the species distributed

97    across the sites. By doing so, PCPS analysis renders single variables that synthesize

98    phylogenetic gradients, and thereby can be used to analyses addressing environmental, spatial

99    or biogeographic determinants of phylogeny-mediated species distribution (Duarte, Bergamin,

100   Marcilio-Silva, Seger & Marques, 2014; Carlucci et al., 2016).

101        In this study we introduce HaploVectors, a new flexible analytical tool for molecular

102   ecologists to disentangle biogeographic or environmental factors driving haplotypic

103   distribution across space. For this, HaploVectors extends the application of phylogenetic

104   fuzzy weighting in order to describe the distribution of haplotypes across sets of localities,

105   allowing flexible exploratory analysis. Moreover, by applying appropriate null models,

5

106    HaploVectors allows to analyze multiple determinants of the frequencies of haplotypes, as

107    well as the number of mutations separating different haplotypes, across sets of localities,

108    providing robust hypotheses tests for phylogeographical studies. We demonstrate the

109    application of HaploVectors in two empirical datasets.

110    **Materials and Methods**

111    *Haplotypic eigenvectors*

112    Haplotypic eigenvector analysis and hypotheses tests based on null models were implemented

113    in the R package *HaploVectors* (freely available at

114    https://github.com/vanderleidebastiani/HaploVectors). The function 'HaploVectors' allows

115    defining the frequency of each haplotype across a set of localities from where the individuals

116    were sampled and to weight those frequencies across localities according to the number of

117    mutational steps between all pairs of haplotypes arranged into a network based on fuzzy set

118    theory (Zadeh, 1965). Further, the function implements null model tests to analyze the

119    influence of environmental, biogeographic or spatial variables on the distribution of

120    haplotypes across sets of localities, as well as to estimate to what extent such influence is

121    mediated by the evolutionary distance between haplotypes. The rationale underlying the

122    HaploVectors approach, including the null model tests, was adapted from phylogenetic fuzzy

123    weighting (Pillar & Duarte, 2010), a method originally developed to analyze multiple

124    determinants of phylogenetic composition across metacommunities based on fuzzy set theory

125    (see details on the method in Duarte et al., 2016).

126         The first analytical step implemented in HaploVectors consists of defining haplotypes

127    for a set of samples and computing the frequency (number of individuals) of each haplotype

128    per locality. For this, two input datasets are required: (1) a *.fas file containing aligned

6

129    genetic sequences for each sampled individual, and (2) a matrix describing the incidence of

130    each individual (rows) in a given locality (column). Based on these two datasets, the function

131    'HaploNetDist' extracts the haplotypes for the *.fas file using the function 'haplotype', and

132    computes an haplotypic network using the function 'haploNet', functions originally

133    implemented in the R package *pegas* (Paradis, 2010).

134          Further, 'HaploNetDist' computes the frequency of each haplotype per locality

135    (matrix **W**), and submit **W** to fuzzy weighting as follows: From the matrix **D** describing all

136    pairwise distances between haplotypes based on the number of mutational steps between any

137    pair of haplotypes, remove all distances between haplotypes *not* connected in the network,

138    replacing them by the mutations separating two haplotypes in the network plus one. This

139    procedure generates matrix **D**$_\text{N}$, which reproduces haplotype connections described in the

140    network. Matrix **D**$_\text{N}$ is further converted into a similarity matrix **S** describing the similarities

141    between all pairs of haplotypes $i$ and $j$ ($s_{ij}$, ranging between 0 and 1), as follows:

142

$$s_{ij} = 1 - \left( \frac{d_{ij}}{\max\left(d_{ij}\right)} \right)$$

143          $d_{ij}$ is the number of mutations separating the haplotypes $i$ and $j$ in **D**$_\text{N}$, and max ($d_{ij}$) is

144    the maximum number of mutations between any pair of haplotypes in the network. Matrix **S**

145    is then standardized by marginal total within columns, generating a matrix **Q** that describes

146    fuzzy belongings between haplotypes (Pillar & Duarte, 2010; Duarte et al., 2016). Matrix **Q**,

147    containing haplotypic fuzzy sets are then employed to weight the frequencies of haplotypes

148    per locality described in matrix **W** by their evolutionary relatedness, generating matrix **P**,

149    which describes localities by their haplotypic composition, that is, haplotype frequencies per

150    locality weighted by their evolutionary relatedness. The output of HaploVectors function

151    provides matrices **W**, **D**, **D**$_\text{N}$, **Q** and **P**.

7

152    Performing principal coordinates on **P** generates haplotypic eigenvectors, which

153    decompose the total variation in the haplotypic composition across the set of localities into

154    independent fractions proportional to its respective eigenvalue λ. Those eigenvectors

155    representing the higher amount of variation in **P** can be used to explore major trends in

156    haplotype distribution across the localities. Those localities sharing most haplotypes will

157    show similar scores, and therefore will group to each other in the scatter plot. Thus, this

158    scatter plot allows simultaneously explore evolutionary links among haplotypes and localities.

159    The function 'HaploVectors' also allows analyzing multiple environmental,

160    biogeographic or spatial determinants of haplotype distribution across a set of localities, and

161    therefore is a useful tool for robust hypothesis test in phylogeography. For this, the function

162    implements two null model tests, adapted from Duarte et al. (2016) and designed to test the

163    following hypotheses:

164    Hypothesis 1: *A given environmental, biogeographic or spatial factor **E** affects the*

165    *distribution of haplotypes across a set of localities.* This hypothesis is tested by means of a

166    null model called *site shuffle*, which is a classical permutation-based procedure assuming

167    independence between haplotypes and localities. The test can be performed based on either

168    haplotypic dissimilarities between localities computed based on **P** (hereafter $D_P$) using an

169    appropriate resemblance measure, such as Euclidean distances or Bray-Curtis dissimilarities

170    (Legendre & Legendre, 2012), or directly on single haplotypic eigenvectors. The test consists

171    of 1) computing a $F_{Obs}$ statistic. If $D_P$ is modeled on **E**, the test is based on a dissimilarity

172    regression on distance matrices (hereafter called ADONIS; see McArdle & Anderson, 2001).

173    For single haplotypic eigenvectors modeled on **E**, the test is an ordinary linear squares (OLS)

174    model; 2) freely permuting the localities a number of times (say 999); 3) at each permutation

175    step, computing $F_{null}$; and 4) defining the probability of obtaining the observed statistic by

176  chance ($H_0 = F_{Obs} \leq F_{null}$), as the proportion of permutations in which $F_{null}$ exceeded $F_{Obs}$.

177  Using this procedure, the test simultaneously addresses the influence of **E** on the distribution

178  of haplotypes across the localities (the number of haplotypes shared between pairs of

179  localities) and to what extent such influence is mediated by the evolutionary relatedness

180  between the haplotypes (the number of mutational steps between the pairs of haplotypes).

181  Therefore, even if this first hypothesis is corroborated ($H_0$ is rejected), such test does not

182  allow us to conclude that the influence of **E** on haplotype distribution is or is not dependent

183  on the evolutionary relatedness among haplotypes. To accomplish that, it is necessary to test a

184  second hypothesis, which is conditioned on the validity of the first one:

185  Hypothesis 2: *The influence of **E** on the distribution of haplotypes across a set of*

186  *localities depends on the evolutionary relatedness among them.* To test this hypothesis, a

187  second round of permutations (*network shuffle*) is needed in order to compute $F_{null}$.

188  Accordingly, the frequency of haplotypes in **W** is kept constant across the localities while

189  evolutionary relatedness between them is permuted by haplotype label shuffling (Kembel et

190  al., 2010). After computing $F_{Obs}$ (step 1, as described for site shuffle), the procedure consists

191  of 2) freely permuting haplotype labels in the network to generate random evolutionary

192  relatedness among haplotypes and computing null matrices **D**, **D**$_N$, **Q** and **P**. The procedure is

193  repeated 999 times; 3) at each permutation procedure, computing null D$_P$ and, if necessary,

194  null haplotypic eigenvectors. In this later case, null haplotypic eigenvectors are submitted to

195  procrustean adjustment (Jackson, 1995) and fitted values between observed and null

196  eigenvectors are obtained; 4) take null D$_P$ or selected adjusted null eigenvectors as response

197  variable in ADONIS or OLS, respectively, using **E** as predictor, and compute $F_{null}$ values; 5)

198  generating a set of $F_{null}$ to get the probability under the null hypothesis ($H_0 = F_{Obs} \leq F_{null}$); 6)

199  defining a probability under the null hypothesis.

200    By performing both null model tests, two probability values are generated. Previous

201    analyses using simulated data demonstrated that both null models show appropriate type I

202    error and statistical power (Duarte et al., 2016). Using site shuffle, whenever the null

203    hypothesis is rejected, we conclude that **E** affects the distribution of haplotypes across a set of

204    localities (hypothesis 1). Then we proceed to test for the hypothesis 2 (via network shuffle). If

205    the null hypothesis is rejected, we conclude that the influence of **E** on the distribution of

206    haplotypes across the localities depends on the evolutionary relatedness among them. In Fig.

207    1 we illustrate the expected distribution of haplotypes and the respective probabilities

208    generated under site and network shuffle models.

209    *Application using empirical datasets*

210    We demonstrate the application of HaploVectors in phylogeographical analyses through two

211    empirical data sets: a set of cpDNA sequences from 333 adult individuals of *Eugenia*

212    *dysenterica*, a tree species from Myrtaceae family, sampled from 23 localities (Lima, Telles,

213    Chaves, Lima-Ribeiro & Collevatti, 2017), and a set of cpDNA sequences from 257 adult

214    individuals of *Mauritia flexuosa*, a palm tree sampled from 26 localities (Lima, Lima-Ribeiro,

215    Tinoco, Terribile & Collevatti, 2014). Both data sets are available at GenBank (accession

216    numbers: MF752706- MF753038 and KC527837-KC528609, respectively).

217    Dataset 1: *Eugenia dysenterica*

218    Using a phylogeographical approach, Lima et al. (2017) investigated the spatial pattern of

219    genetic diversity on *E. dysenterica*, a widely distributed tree species in the Brazilian savanna.

220    Four regions of the chloroplast were sequenced from individuals sampled at 23 populations

221    throughout its distribution (for details about sampling and genetic data see Lima et al., 2017).

222    The authors inferred the phylogenetic relationships among haplotypes using median-joining

10

223    network analysis. Furthermore, Analysis of Molecular Variance (AMOVA) was used to

224    analyze spatial patterns of genetic variation among biogeographic regions of the Cerrado

225    biome (Table 1). Although AMOVA pointed out genetic differentiation among sites located at

226    three different Cerrado regions (Central, Northeast and Southeastern; $P < 0.001$), the results

227    of the network analysis visually suggested that the phylogenetic relationships among

228    haplotypes did not match the geographical distribution of the lineages (Fig. 2a).

229         We analyzed the variation in the distribution of haplotypes across the biogeographic

230    regions of the Cerrado using HaploVectors approach. Our hypotheses propose that (1) the

231    spatial distribution of haplotypes varies among the different biogeographic regions of the

232    Cerrado biome, and that (2) the biogeographic distribution of haplotypes across the Cerrado

233    regions depends on the evolutionary relationships among them. We tested these hypotheses

234    using ADONIS, based on haplotypic dissimilarities between localities ($D_P$), and OLS using

235    haplotypic eigenvectors. For ADONIS, we computed matrix **P** using log-transformed

236    frequencies of haplotypes per site (matrix **W**) and square-rooted Hamming distances between

237    haplotypes (matrix $D_N$). Haplotypic dissimilarities between sites were computed using

238    Euclidean distances. The same three biogeographic regions of Cerrado were taken as a

239    categorical predictor in the analysis (matrix **E**). For OLS, we first computed haplotypic

240    eigenvectors (haplovectors) based on $D_P$. Haplovectors containing more than 5% of total

241    information in **P** were taken as response variables in linear models, while **E** was used as

242    predictor.

243    Dataset 2: *Mauritia flexuosa*

244    *Mauritia flexuosa* is a dioecious palm species distributed widely across South America,

245    occurring in Brazilian savannas and Amazonia (Lima et al., 2014). Because its occurrence is

246    associated with wetlands, Lima et al. (2014) investigated spatial patterns in chloroplast

247    genome regions among populations of *M. flexuosa* occurring across four different river basins

248    (Amazon, Araguaia/Tocantins, Paraná and São Francisco). For this purpose, the authors

249    estimated phylogenetic relationships among haplotypes using the median-joining network

250    analysis and performed a hierarchical analysis of molecular variance (AMOVA) to analyze

251    the genetic differentiation among populations from different river basins. The haplotypic

252    network (Fig. 2c) did not allow a clear congruence between the geographic distribution of

253    phylogenetic lineages and river basins, although AMOVA found significant genetic

254    differentiation among river basins (Table 1). For this dataset we performed similar analyses as

255    described for *E. dysenterica*.

256    **Results**

257    For analyses performed using *E. dysenterica* dataset, ADONIS indicated that sites occurring

258    at the same biogeographic region of the Cerrado biome share more haplotypes with each other

259    than with sites located at different regions ($P_{\text{site shuffle}} = 0.001$, Table 1); nonetheless, such

260    difference in haplotype composition is not related with evolutionary relatedness among

261    haplotypes ($P_{\text{network shuffle}} = 0.973$). The two first haplovectors (Fig. 2b), containing 32% and

262    17% of total variation in haplotype composition of sites, respectively, indicated association

263    between haplotype distribution across sites and biogeographic regions ($P_{\text{site shuffle}} < 0.02$). The

264    first haplovector shows a separation of populations from Southeast Brazil (green circles) from

265    other regions (Fig. 2b), which relies only on haplotype composition, but not evolutionary

266    relatedness among haplotypes ($P_{\text{network shuffle}} = 0.936$). On the other hand, the second

267    haplovector discriminated some populations from Northeast Brazil from other regions (Fig.

268    2b). In this case, an evolutionary signal in the association between haplotype composition of

269    sites and biogeographic regions was detected ($P_{\text{network shuffle}} = 0.024$), based mostly on the

270    evolutionary path from haplotype one to 19 to 11 (Figs. 2a, 2b).

12

271     For the *Mauritia flexuosa* dataset, ADONIS rejected null hypothesis for both site

272     shuffle and network shuffle null models (Table 1), indicating that haplotype composition

273     differed between river basins ($P_{\text{site shuffle}} = 0.01$), and that such difference was mediated by

274     evolutionary relatedness between haplotypes ($P_{\text{network shuffle}} = 0.005$). For this dataset, only the

275     first haplovector (Fig. 2d), containing 42% of total variation in haplotype composition of

276     sites, indicated association between haplotype distribution across sites and river basins ($P_{\text{site}}$

277     $_{\text{shuffle}} = 0.021$), which was mediated by evolutionary relatedness between haplotypes ($P_{\text{network}}$

278     $_{\text{shuffle}} = 0.007$). The first haplovector showed a clear separation between populations located in

279     Amazon (left side of the plot in Fig. 2d), which were related to haplotypes seven and eight

280     (Fig. 2c), and Araguaia/Tocantins basins (right side of the plot in Fig. 2d), mostly associated

281     with haplotypes one to five.

282     **Discussion**

283     Current implemented methods for phylogeographical analyses treat haplotypic frequency

284     across localities over a given environmental gradient/factor and phylogenetic relationships

285     among those haplotypes in disconnected manners, lacking a clear conceptual framework to

286     integrate both. Haplovectors provide such integration, allowing disentangling the

287     environmental or biogeographic influence on haplotypic distribution and assessing whether

288     that distribution is resulting from the evolutionary relationship among haplotypes. In cases

289     where AMOVA (Excoffier et al., 1992) and analyses using haplotype networks (e.g.,

290     Templeton, 1998) reveal contradicting results of haplotypic distributions, we propose that

291     HaploVectors can elucidate the conundrum.

292     In the first example (*E. dysenterica* dataset) we found that the influence of

293     biogeographic regions at structuring haplotypes is independent from the evolutionary

294     relatedness among them (i.e. the number of mutational steps separating haplotypes). This

13

295    shows that biogeographic regions are indeed structuring haplotypes (i.e. different haplotypes

296    can be found in different regions, implying few haplotypes occurring in more than one

297    region), however, inside any given biogeographic region, haplotypes are not the closely

298    related to each other based on the phylogenetic relationships among haplotypes. This means

299    that phylogenetic closely related haplotypes occur in distinct biogeographic regions, and each

300    biogeographic region comprises exclusive haplotypes from multiple evolutionary origins.

301          Different from that observed in *E. dysenterica* dataset, for the *M. flexuosa* dataset we

302    found that haplotype composition differed between river basins and this difference was

303    associated with the evolutionary relatedness among haplotypes. This reveals that haplotypes

304    that occur in the same river basins are more phylogenetically related than those that occur in

305    different river basins. These interpretations provided by HaploVectors solve the apparent

306    paradox found in the results of previous analyses in both cases: the AMOVA found

307    haplotypic differences among biogeographic regions and river basins; however, the

308    haplotypic network failed to revel a clear structured haplotypic distribution over the same

309    regions. The null model tests implemented in HaploVectors permit treating the haplotype

310    frequency across localities on a given environmental factor, independently from the

311    phylogenetic similarities among haplotypes. The combination of both tests in a joint approach

312    allows for tracing a complete picture of the evolutionary history of populations.

313          We hope that this approach will be useful in all cases where the distribution of

314    haplotypes is hypothesized to be under the influence of an environmental or biogeographic

315    factor. These questions are likely to be encountered with increasingly frequency by molecular

316    ecologists and phylogeographers.

317    **Acknowledgements**

14

325    **References**
326    Anderson, M. J. (2001). A new method for non-parametric multivariate analysis of variance.
327        *Austral Ecology, 26*, 32-46. doi: 10.1111/j.1442-9993.2001.01070.pp.x
328    Avise, J. C., Arnold, J. & Ball, R.M., Bermingham, E., Lamb, T., Neigel, J.E., … Saunders,
329        N.C. (1987). Intraspecific phylogeography: the mitochondrial DNA bridge between
330        population genetics and systematics. *Annual Review of Ecology and Systematics*, *18*, 489-
331        522. doi: 10.1146/annurev.es.18.110187.002421
332    Avise, J. C. (2009). Phylogeography: retrospect and prospect. *Journal of Biogeography*, *36*,
333        3-15. doi: 10.1111/j.1365-2699.2008.02032.x
334    Bandelt, H. J., Forster, P. & Röhl, A. (1999). Median-joining networks for inferring
335        intraspecific phylogenies. *Molecular Biology and Evolution*, *16*, 37-48. doi:
336        10.1093/oxfordjournals.molbev.a026036
337    Beheregaray, L. B. (2008). Twenty years of phylogeography: the state of the field and the
338        challenges for the Southern Hemisphere. *Molecular Ecology*, *17*, 3754-3774. doi:
339        10.1111/j.1365-294X.2008.03857.x
340    Blom, M. P. K., Horner, P. & Moritz, C. (2016). Convergence across a continent: adaptive
341        diversification in a recent radiation of Australian lizards. *Proceedings of the Royal
342        Society B: Biological Sciences*, *283*, 20160181. doi: 10.1098/rspb.2016.0181
343    Carlucci, M. B., Seger, G. D. S., Sheil, D., Amaral, I. L., Chuyong, G. B., Ferreira, L. V., …
344        Duarte, L. S. (2016). Phylogenetic composition and structure of tree communities shed
345        light on historical processes influencing tropical rainforest diversity. Ecography, *40*, 521-
346        530. doi: 10.1111/ecog.02104
347    Carnaval, A. C., Waltari, E., Rodrigues, M. T., Rosauer, D., VanDerWal, J., Damasceno, R.,
348        ... Moritz, C. (2014) Prediction of phylogeographic endemism in an environmentally
349        complex biome. *Proceedings of the Royal Society B: Biological Sciences*, *281*, 20141461.
350        doi: 10.1098/rspb.2014.1461
351    Diniz-Filho, J. A. F., Telles, M. P. C., Bonatto, S. L., Eizirik, E., Freitas, T. R. O., De Marco
352        Jr, P., … Soares, T. N. (2008). Mapping the evolutionary twilight zone: molecular
353        markers, populations and geography. *Journal of Biogeography*, *35*, 753-763. doi:
354        10.1111/j.1365-2699.2008.01912.x
355    Duarte, L. D. S. (2011). Phylogenetic habitat filtering influences forest nucleation in
356        grasslands. *Oikos*, 120, 208-215. doi: 10.1111/j.1600-0706.2010.18898.x

357 Duarte, L. D. S., Bergamin, R. S., Marcilio-Silva, V., Seger, G. D. S. & Marques, M. C. M.
358     (2014). Phylobetadiversity among forest types in the Brazilian Atlantic Forest complex.
359     PLoS ONE 9:e105043. doi: 10.1371/journal.pone.0105043
360 Duarte, L. D. S., Debastiani, V. J., Freitas, A. V. L. & Pillar, V. D. (2016). Dissecting
361     phylogenetic fuzzy weighting: theory and application in metacommunity phylogenetics.
362     *Methods in Ecology and Evolution*, 7, 937-946. doi: 10.1111/2041-210X.12547
363 Epperson, B. K. (2005). Estimating dispersal from short distance spatial autocorrelation.
364     *Heredity*, *95*, 7-15. doi: 10.1038/sj.hdy.6800680
365 Excoffier, L., Smouse, P. E. & Quattro, J. M. (1992). Analysis of molecular variance inferred
366     from metric distances among DNA haplotypes: Application to human mitochondrial
367     DNA restriction data. *Genetics*, *131*, 479-491.
368 Fitzpatrick, B. M. (2009). Power and sample size for nested analysis of molecular variance.
369     *Molecular Ecology*, *18*, 3961-6. doi: 10.1111/j.1365-294X.2009.04314.x
370 Freeland, J. R. (2014). *Molecular Ecology*. In: eLS. John Wiley & Sons, Ltd: Chichester. doi:
371     10.1002/9780470015902.a0003268.pub2
372 Jackson, D. A. (1995). Protest: A Procrustean Randomization test of community environment
373     concordance. *Ecoscience.*, *2*, 297-303. doi: 10.1080/11956860.1995.11682297
374 Kembel, S. W., Cowan, P. D., Helmus, M. R., Cornwell, W. K., Morlon, H., Ackerly, D. D.,
375     … Webb, C. O. (2010). Picante: R tools for integrating phylogenies and ecology.
376     *Bioinformatics*, *26*, 1463-4. doi: 10.1093/bioinformatics/btq166
377 Kong, S., Sánchez-Pacheco, S. J. & Murphy, R. W. (2015). On the use of median-joining
378     networks in evolutionary biology. *Cladistics*, *32*, 691-699. doi: 10.1111/cla.12147
379 Legendre, P. & Legendre, L. (2012) *Numerical Ecology* (3rd ed.). Elsevier, Amsterdam.
380 Lima, N. E., Lima-Ribeiro, M. S., Tinoco, C. F., Terribile, L. C. & Collevatti, R. G. (2014).
381     Phylogeography and ecological niche modelling, coupled with the fossil pollen record,
382     unravel the demographic history of a Neotropical swamp palm through the Quaternary.
383     *Journal of Biogeography*. doi: 10.1111/jbi.12269
384 Lima, J. S., Telles, M. P. C., Chaves, L. C., Lima-Ribeiro, M. S. & Collevatti, R. G. (2017).
385     Demographic stability and high historical connectivity explain the diversity of a savanna
386     tree species in the Quaternary. *Annals of Botany*, *119*, 645-657. doi:
387     10.1093/aob/mcw257
388 Maestri, R., Fornel, R., Gonçalves, G. L., Geise, L., Freitas, T. R. O. & Carnaval, A. C.
389     (2016). Predictors of intraspecific morphological variability in a tropical hotspot:
390     comparing the influence of random and non-random factors. *Journal of Biogeography*,
391     *43*, 2160-2172. doi: 10.1111/jbi.12815
392 Mardulyn, P. (2012). Trees and */*or networks to display intraspecific DNA sequence variation?
393     *Molecular Ecology*, 21, 3385-3390. doi: 10.1111/j.1365-294X.2012.05622.x
394 Manel, M., Scwartz, M. K., Luikart, G. & Taberlet, P. (2003). Landscape genetics: combining
395     landscape ecology and population genetics. *Trends in Ecology and Evolution*, *18*, 189-
396     197. doi: 10.1016/S0169-5347(03)00008-9
397 McArdle, B. H. & Anderson, M. J. (2001). Fitting multivariate models to community data: a
398     comment on distance-based redundancy analysis. *Ecology*, *82*, 290-297. doi:
399     10.1890/0012-9658(2001)082[0290:FMMTCD]2.0.CO;2
400 Miller, M. P. (2005). Alleles in Space: computer software for the joint analysis of
401     interindividual spatial and genetic information. *Journal of Heredity*, *96*, 722-724. doi:
402     10.1093/jhered/esi119

16

403 Papadopoulou, A. & Knowles, L. L. (2016). Toward a paradigm shift in comparative
404     phylogeography driven by trait-based hypotheses. *PNAS*, *113*, 8018-8024. doi:
405     10.1073/pnas.1601069113
406 Paradis, E. (2010). pegas: an R package for population genetics with an integrated modular
407     approach. *Bioinformatics, 26*, 419-420. doi:10.1093/bioinformatics/btp696
408 Pillar, V. & Orlóci, L. (1996). On randomization testing in vegetation science: multifactor
409     comparisons of relevé groups. *Journal of Vegetation Science*, *7*, 585-592. doi:
410     10.2307/3236308
411 Pillar, V. & Duarte, L. S. (2010). A framework for metacommunity analysis of phylogenetic
412     structure. *Ecology Letters*, 13, 587–596. doi: 10.1111/j.1461-0248.2010.01456.x
413 Raffini, F., Fruciano, C. & Meyer, A. (2018). Morphological and genetic correlates in the
414     left–right asymmetric scale-eating cichlid fish of Lake Tanganyika. *Biological Journal of*
415     *the Linnean Society*, *124*, 67-84. doi: 10.1093/biolinnean/bly024
416 Storfer, A., Murphy, M. A., Evans, J. S., Goldberg, C. S., Robinson, S., Spear, S. F., …Waits,
417     L. P. (2007). Putting the 'landscape' in landscape genetics. *Heredity*, *98*, 128-142. doi:
418     10.1038/sj.hdy.6800917
419 Templeton, A. R. (1998). Nested clade analysis of phylogeographical data: testing hypotheses
420     about gene flow and population history. *Molecular Ecology*, *7*, 381-397. doi:
421     10.1046/j.1365-294x.1998.00308.x
422 Templeton, A. R. (2004). Statistical phylogeography: methods of evaluating and minimizing
423     inference errors. *Molecular Ecology*, *13*, 789-809. doi: 10.1046/j.1365-
424     294X.2003.02041.x
425 Turchetto-Zolet, A. C., Pinheiro, F., Salgueiro, F. & Palma-Silva, C. (2013).
426     Phylogeographical patterns shed light on evolutionary process in South America.
427     *Molecular Ecology*, *22*, 1193-1213. doi: 10.1111/mec.12164
428 Zadeh, L. A. (1965). Fuzzy sets. *Information and Control*, *8*, 338-353. doi: 10.1016/S0019-
429     9958(65)90241-X
430

431    **Data accessibility**

432    The data are archived in GenBank (accession numbers: MF752706-MF753038 and

433    KC527837-KC528609).

434    **Author contributions**

435    LD designed research; All authors performed research; LD and JL analyzed data; all authors
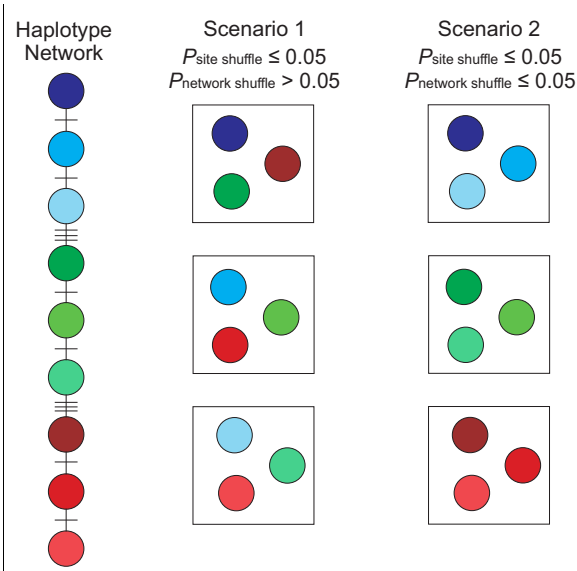
436    wrote the paper.

437

438 **Table 1.** Comparison between results obtained by Analysis of Molecular Variance (AMOVA) and null model
439 tests implemented in HaploVectors. $D_P$ - haplotypic dissimilarities between localities computed based on matrix
440 **P**.

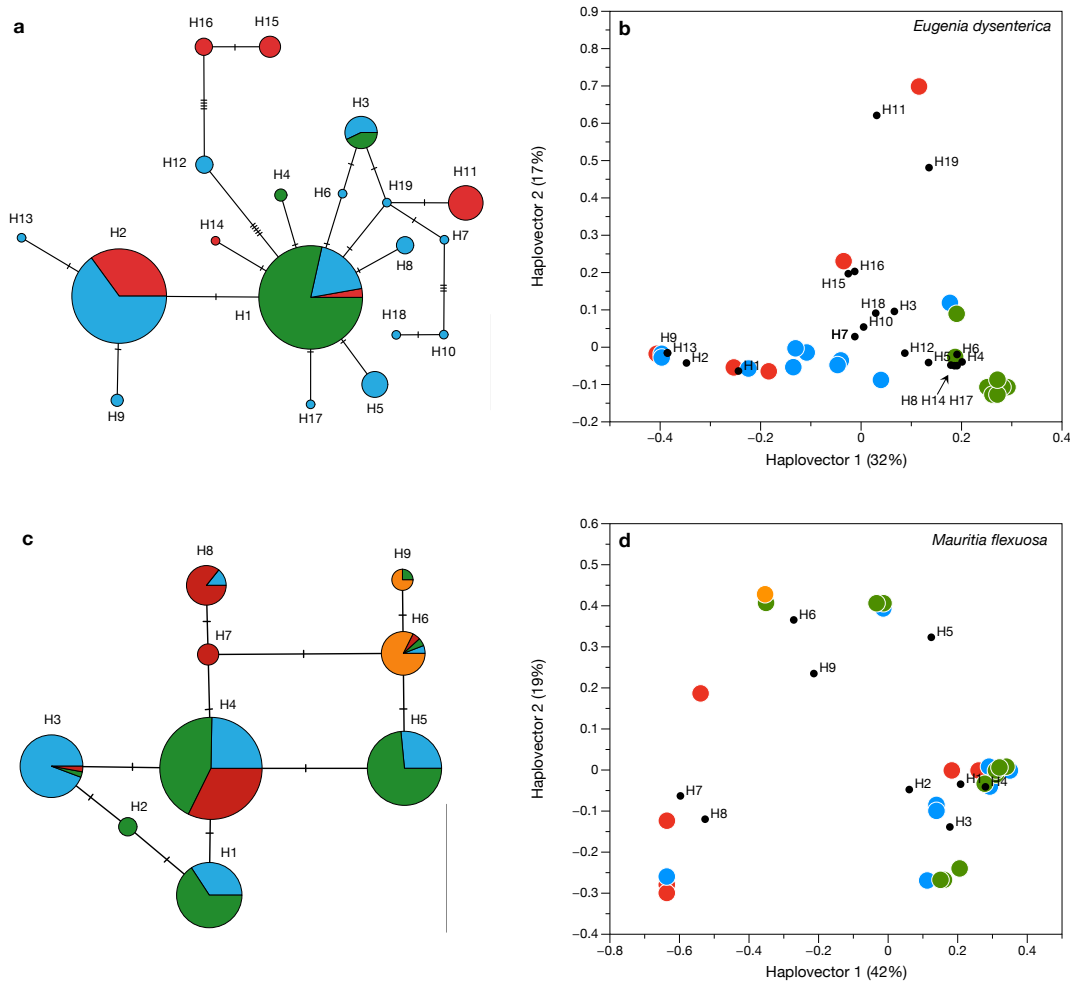| Source | Case species | Question | Method | Results |
|---|---|---|---|---|
| Lima et al. (2017) | *Eugenia dysenterica* | Differences among Cerrado regions? | AMOVA | $F_{CT}= 0.164$; $P < 0.001$ |
| | | | ADONIS on $D_P$ | $R^2 = 0.20$; $F_{Obs} = 5.38$; $P_{\text{site shuffle}} = 0.001$, $P_{\text{network shuffle}} = 0.973$ |
| | | | OLS on haplovectors | *Haplovector 1 (32%)*: $R^2 = 0.58$; $F_{Obs} = 11.44$; $P_{\text{site shuffle}} = 0.001$, $P_{\text{network shuffle}} = 0.936$ |
| | | | | *Haplovector 2 (17%)*: $R^2 = 0.15$; $F_{Obs} = 2.36$; $P_{\text{site shuffle}} = 0.016$, $P_{\text{network shuffle}} = 0.024$ |
| | | | | *Haplovector 3 (11%)*: $R^2 < 0.01$; $F_{Obs} = 0.34$; $P_{\text{site shuffle}} = 0.608$, $P_{\text{network shuffle}} = 0.854$ |
| | | | | *Haplovector 4 (9%)*: $R^2 < 0.01$; $F_{Obs} = 0.67$; $P_{\text{site shuffle}} = 0.432$, $P_{\text{network shuffle}} = 0.111$ |
| Lima et al. (2014) | *Mauritia flexuosa* | Differences among basins? | AMOVA | $F_{CT} = 0.387$; $P < 0.050$ |
| | | | ADONIS on $D_P$ | $R^2 = 0.23$; $F_{Obs} = 2.24$; $P_{\text{site shuffle}} = 0.010$, $P_{\text{network shuffle}} = 0.005$ |
| | | | OLS on haplovectors | *Haplovector 1 (42%)*: $R^2 = 0.24$; $F_{Obs} = 2.78$; $P_{\text{site shuffle}} =0.021$, $P_{\text{network shuffle}} = 0.007$ |
| | | | | *Haplovector 2 (19%)*: $R^2 = 0.07$; $F_{Obs} = 1.23$; $P_{\text{site shuffle}} = 0.211$, $P_{\text{network shuffle}} = 0.460$ |
| | | | | *Haplovector 3 (13%)*: $R^2 = 0.03$; $F_{Obs} = 0.93$; $P_{\text{site shuffle}} = 0.304$, $P_{\text{network shuffle}} = 0.695$ |
| | | | | *Haplovector 4 (6%)*: $R^2 = 0.03$; $F_{Obs} = 1.28$; $P_{\text{site shuffle}} = 0.204$, $P_{\text{network shuffle}} = 0.680$ |

441

442

19

443
444
445 **Fig. 1.** Expected distribution of haplotypes (colors denote different haplotypes) across localities (squares) and
446 the respective probabilities of homogeneity of haplotype composition among localities, under site and network
447 shuffle null models. Scenario 1): Localities contain different haplotypes ($P_{\text{site shuffle}} \leq 0.05$), but variation of
448 haplotype composition among sites is not associated with evolutionary structure depicted by the haplotype
449 network ($P_{\text{network shuffle}} > 0.05$). Scenario 2): Localities contain different haplotypes ($P_{\text{site shuffle}} \leq 0.05$), and
450 haplotype distribution across sites is mediated by evolutionary relatedness among them ($P_{\text{network shuffle}} \leq 0.05$).
451

452

**Fig. 2.** Haplotype networks and their respective scatter plots of haplovectors computed for two datasets. 1a-b) *Eugenia dysenterica* (Lima et al. 2017). Red, blue and green circles indicate Northeast, Central and Southeastern biogeographic regions of the Brazilian Cerrado biome, respectively. a) MJN haplotypic network; b) Scatter plot for the two first haplovectors. Black circles indicate haplotypes (H1-H19). 1c-d) *Mauritia flexuosa* (Lima et al. 2014). Red, blue, orange and green circles indicate Amazon, Paraná, São Francisco and Araguaia/Tocantins river basins in South America, respectively. c) MJN haplotypic network; d) Scatter plot for the two first haplovectors. Black circles indicate haplotypes (H1-H9).

460