1    Title: Assigning occurrence data to cryptic taxa improves climatic niche assessments:

2    Biodecrypt, a new tool tested on European butterflies

3

4    Leonardo Platania[1]†, Mattia Menchetti[2]†, Vlad Dincă[3], Cecília Corbella[1], Isaac Kay-lavelle[1]

5    Roger Vila[1], Martin Wiemers[4,5], Oliver Schweiger[5], Leonardo Dapporto[2]*

6

7    [1] Institut de Biologia Evolutiva (CSIC - Universitat Pompeu Fabra), Barcelona, Spain

8    [2] ZEN Lab, Dipartimento di Biologia dell'Università di Firenze, via Madonna del Piano 6

9    50019 Sesto Fiorentino, Italy

10   [3] Ecology and Genetics Research Unit, PO Box 3000, University of Oulu, 90014, Finland

11   [4] Senckenberg Deutsches Entomologisches Institut, Eberswalder Str. 90, 15374 Müncheberg,

12   Germany

13   [5] Helmholtz Centre for Environmental Research - UFZ, Department of Community Ecology,

14   06120 Halle, Germany

15   Corresponding author: leonardo.dapporto@unifi.it

16

17

18   Running title: *biodecrypt*, assigning occurrence data

19

20

21

22

23

24

25

26 **Abstract**

27 **Aim**

28 Occurrence data are fundamental to macroecology, but accuracy is often compromised when

29 multiple units are lumped together (e.g. in recently separated cryptic species or citizen

30 science records). Using amalgamated data leads to inaccuracy in species mapping, to biased

31 beta-diversity assessments and to potentially erroneously predicted responses to climate

32 change. We provide a set of R functions (biodecrypt) to objectively attribute undetermined

33 occurrences to the most probable taxon based on a subset of identified records.

34

35 **Innovation**

36 Biodecrypt assumes that unknown occurrences can only be attributed at certain distances

37 from areas of sympatry. The function draws concave hulls based on the subset of identified

38 records; subsequently, based on hull geometry, it attributes (or not) unknown records to a

39 given taxon. Concavity can be imposed with an alpha value and sea or land areas can be

40 excluded. A cross-validation function tests attribution reliability and another function

41 optimizes the parameters (alpha, buffer, distance ratio between hulls). We applied the

42 procedure to 16 European butterfly complexes recently separated into 33 cryptic species for

43 which most records were amalgamated. We compared niche similarity and divergence

44 between cryptic taxa, and we re-calculated and contributed updated CLIMBER variables for

45 climatic preferences.

46

47 **Main conclusions**

48 Biodecrypt showed a cross-validated correct attribution of known records always ≥98% and

49 attributed more than 80% of unknown records to the most likely taxon in parapatric species.

50 The functions determined where records can be assigned even for largely sympatric species,

51    and highlighted areas where further sampling is required. All the cryptic taxa showed

52    significantly diverging climatic niches, reflected in different values of mean temperature and

53    precipitation compared to the values originally provided in the CLIMBER database. The

54    substantial fraction of cryptic taxa existing across different taxonomic groups and their

55    divergence in climatic niches highlights the importance of using reliably assigned occurrence

56    data in macroecology.

57

60

61

62

## Introduction

A solid record of species occurrence data is key to understand the multiple factors defining their large-scale geographic distributions and, by means of ecological niche modelling, to assess and project their responses to changing environmental conditions in terms of range expansion or contraction (Franklin, 2010; Schweiger et al., 2012; Hortal et al., 2015; Thuiller et al., 2016). In addition, resulting species-specific niche characteristics have provided conservation biogeography with a powerful set of species features, such as measures of mean and variation in multiple climatic characteristics (e.g. CLIMBER variables; Schweiger, Harpke, Wiemers, & Settele 2014), useful for assessing community-wide responses to global change (Devictor et al., 2012; Zografou et al., 2014; Herrando et al., 2019).

Generating reliable species occurrence data at continental scale requires an enormous effort and such databases have been assembled over decades of field research, often based on, or improved by, citizen science projects (Dennis, Morgan, Brereton, Roy, & Fox, 2017; Titeux et al., 2017). In addition, proper definition and discrimination of species are necessary for reliable niche modelling, as well as for the identification of environmental preferences of species and the derived specific indices. However, the existence of a considerable fraction of cryptic species in almost all groups of living organisms (Bickford et al., 2007) poses a serious challenge to our understanding of diversity in general, and to this line of research in particular. When a taxon, believed to represent a single species, is recognized as two or multiple cryptic species, the occurrence data accumulated for decades suddenly become obsolete. Researchers often amalgamate the occurrence data for cryptic taxa, but this approach ignores a substantial fraction of diversity in terms of species identity, distribution, evolution, and its potential dynamics in changing environments.

Most complexes of cryptic taxa are parapatric with minimal areas of sympatry, frequently because they evolved in allopatry and the achievement of secondary sympatry is delayed by

88    (i) limited dispersal and competition due to a still incomplete separation of ecological niches

89    or by (ii) reproductive interference due to the lack of a pre-mating barrier (Pigot & Tobias,

90    2013, 2015; Vodă, Dapporto, Dincă, & Vila, 2015). Only a minor fraction of cryptic taxa are

91    largely sympatric and these typically show strong reproductive barriers (Dincă, Lukhtanov,

92    Talavera, & Vila, 2011; Dincă et al., 2013). Because cryptic species tend to show parapatric

93    distributions (Waters, 2011; Vodă et al., 2015; Dapporto et al., 2017; Scalercio et al., 2020),

94    they encompass a conspicuous fraction of beta-diversity (Vodă et al., 2015) and, since they

95    inhabit different areas, are expected to be adapted to different climates (Toews, Mandic,

96    Richards, & Irwin, 2014). For this reason, cryptic species are supposed to react differently to

97    climatic changes and any modelling based on amalgamated records is likely inaccurate

98    (Lecocq, Harpke, Rasmont, & Schweiger, 2019).

99    We provide a methodology to objectively attribute occurrence data previously referring to a

100   single taxon to the most likely taxon among two or more newly recognized entities, based on

101   a subset of ascertained records and on justifiable geographic rules. We applied this procedure

102   to the cryptic butterfly species of Europe recently separated into different taxa (Wiemers et

103   al., 2018) that were amalgamated in the Distribution Atlas of Butterflies in Europe (Kudrna et

104   al., 2011; Kudrna, Pennerstorfer, & Lux, 2015). This atlas, of which several editions have

105   been published, represents the most comprehensive source of occurrence data for European

106   butterflies, and the data from the 2011 edition were used to calculate the widely used

107   CLIMBER variables describing species ranges and their climatic preferences over Europe

108   (Schweiger et al., 2014).

109   We provide six new R functions, added to the recluster R package (Dapporto et al., 2013), for

110   parameter optimisation, record attribution to potential cryptic taxa and for testing the

111   reliability of the procedure. Finally, we provide new climatic variables for the species

112   included in this study and show that cryptic taxa differ substantially in their climatic niches.

113

**Methods**

115

*The algorithm*

The objective of the algorithm is to reliably attribute species membership to a set of ambiguous records belonging to two (or more) cryptic species based on the distribution of a subset of accurately determined records. The main idea is that records from an area where only one taxon occurs can be attributed with confidence, while records from the areas of sympatry or too far from any ascertained record cannot be reliably attributed. For this purpose, we developed a series of R (R Core Team, 2019) functions added to the recluster R package (biodecrypt, biodecrypt.view, biodecrypt.cross, biodecrypt.wrap, optimize.biodecrypt, plot.biodecrypt). The main inputs for the functions are a matrix with longitude and latitude for all the occurrence data and a vector (in the same order) providing their identification. The records identified to species-level (identified records) must be indicated in the vector with a sequential numeric value (1, 2... n), which represents the verified membership to the $n^{th}$ species. The occurrence data with unknown identification (unidentified records) are marked with a 0 (Fig. 1). Based on this vector and on the geographical coordinates of identified records, biodecrypt builds hulls of distribution for each species. In a highly simplified hypothesis, the distribution of a species can be approximated by a convex hull among the geographic coordinates of identified records. Nevertheless, areas of distribution can be largely concave, mostly in geomorphologically highly heterogeneous regions. This is the case for Southern Europe, characterized by the presence of three major peninsulas and several insular systems with contrasting species assemblages (Vodă et al., 2015), as well as by a complex quaternary biogeography (Schmitt, 2007; Dapporto et al., 2019). For this reason, biodecrypt and biodecrypt.view use the function

138    getDynamicAlphaHull of the rangeBuilder R package (https://github.com/ptitle/rangeBuilder,

139    accessed 2020/02/05) to construct concave alpha-hulls.

140    An alpha-hull is a piecewise series of linear simple curves in the Euclidean plane associated

141    with the shape of a finite set of points (Edelsbrunner, Kirkpatrick, & Seidel, 1983). Alpha-

142    hulls generalize the concept of the convex hull since every convex hull is an alpha-hull,

143    whereas not every alpha-hull is a convex hull. Alpha-hulls are not necessarily convex and

144    two points inside an alpha-hull can be connected by a segment not completely lying inside

145    the hull itself. The boundary of the alpha-hull is formed by arcs with radius $\alpha$ (see Fig. 2 for a

146    polygon with $\alpha = 3$). For $\alpha=0$ the hull is reduced to the set of points. For increasing values,

147    the area encompassed by the alpha-hull increases in the form of separate concave polygons

148    connecting an increasing number of points, which in some cases remain isolated. For a very

149    high value of $\alpha$, an alpha-hull corresponds to the convex hull connecting the set of points. For

150    these reasons, alpha-hulls are particularly suitable to model disjunct organism distributions

151    including dot-like populations, which – using convex hulls – appear as continuous areas.

152    The getDynamicAlphaHull function supports an initial alpha value that determines a starting

153    custom degree of concavity, which is increased until a given fraction of identified records are

154    included in the resulting hulls (default 95%) and the number of separate polygons is lower

155    than a custom number (default 10). This function can also remove sea or ground areas from

156    the hulls when terrestrial or aquatic organisms, respectively, are under study, thus improving

157    the precision of the hull geometry. After the construction of the alpha-hulls, biodecrypt

158    attempts the attribution of unidentified records to the most likely species (Fig. 1). For this

159    aim, biodecrypt also requires a buffer and a ratio value (explained below). Based on hull

160    geometry and their relative position, each unidentified record could be either: i) inside more

161    than one hull, ii) inside a single hull, or iii) outside all hulls. The three cases are treated

162    separately.

163

*Cases inside more than one hull*

In this case, the function cannot attribute the unidentified records to a species (case 1 in Fig. 2) and only the a priori identified records belonging to intersection areas are passed to the final vector as identified.

168

*Cases inside a single hull*

The unidentified records falling inside a single hull are attributed to that species if their distance to any other hull is higher than the buffer value (in km) provided by the user (case 2 in Fig. 2). Unidentified records inside the buffer of another hull are not attributed (case 3 in Fig. 2).

174

*Cases outside all hulls*

The unidentified records which do not fall inside any hull are attributed to the closest hull if: i) the distance from the second nearest hull is higher than the buffer and if ii) the ratio between the minimum distance to the second closest hull and to the closest hull is more than the ratio value indicated by the user. For example, in Fig. 2 point 4 is not attributed while point 5 is attributed to *Polyommatus celina*.

181

*Check for distances from the nearest identified record*

As described above, the attribution of unknown records is strictly determined by the distance from the hulls. The biodecrypt function also contains an option ("checkdist") to check if cases attributed to a given species based on relative distance from hulls are closer to an identified record of another species, which may occasionally occur. If this option is selected (default) these cases are not attributed to any species (Fig. 2, case 6).

188

189 Different alpha values can better fit the distribution of a given cryptic species and the optimal

190 alpha value can be evaluated by series of cross-validation analyses using biodecrypt.wrap

191 (see below), or according to the researcher's perception. For this reason, we implemented the

192 biodecrypt.view function providing a visual representation of the alpha-hulls for the different

193 cryptic taxa and given alpha values. The alpha values can be modified until an optimal

194 representation is obtained (Fig. 1). The biodecrypt and biodecrypt.view functions also

195 provide three measures: i) the area occupied by each hull (in $km^2$), ii) the area of all the

196 possible intersections between pairs of hulls (in $km^2$) and iii) the fraction of area of

197 intersection between pairs of hulls.

198

199 *Cross-validation*

200 A third function (biodecrypt.cross) wraps the biodecrypt function to carry out cross-

201 validation of identified cases and to verify the robustness of the attribution of unknown cases

202 (Fig. 2). This function requires the same input of biodecrypt (coordinates and vector with

203 attribution and distance ratio, buffer and alpha values) and a "runs" value defining the

204 number of different runs, thus the fraction of test records in each run. The analysis is repeated

205 as many times as defined in "runs" (a "runs" value of 10 will perform a ten-fold cross-

206 validation). In each run, a randomly selected fraction of 1/"runs" identified records are

207 regarded as unidentified (0 value) and the biodecrypt function is carried out to attribute them.

208 The blind attribution of identified records is compared with their membership and two values

209 are provided: the percentages of cases attributed to a wrong species (misidentified records,

210 MIR) and the percentage of cases not attributed to any species (non-attributed identified

211 records, NIR). MIR and NIR represent measures for the power of the function to correctly

212 attribute unknown records to a given species (NIR) and to avoid mis-identification (MIR).

213    The function also has an option to calculate the percentage of non-attributed unidentified

214    records (NUR) representing the fraction of unknown records that could not be attributed to a

215    species after the biodecrypt function was completed using the parameters provided by the

216    user and the complete set of identified and non-identified records.

217    We also provide a function (biodecrypt.wrap) that replicates the cross-validation analysis by

218    using all possible combinations of a series of distance ratio, alpha and buffer values to

219    compare their resulting MIR, NIR and NUR. To optimise the three parameters for each

220    species, we introduced a combination of $MIR^2+NIR+NUR$ as a penalty value for the different

221    combinations of the parameters. Since the method showing the lowest penalty in cross-

222    validation might not necessarily be the optimal value for the final analysis, all the

223    combinations showing a penalty value not higher than a certain threshold compared with the

224    analysis showing the lowest penalty should be considered as similarly good. We provided a

225    value of 10 as a default, representing a variation of about 3 for each addendum of the penalty.

226    The optimal parameters can then be calculated as mean values of distance ratio, alpha and

227    buffer among those used in these cross-validation analyses, weighted by 1/penalty to provide

228    an increasing contribution to the solutions with low penalty values. This is done by

229    optimise.biodecrypt, calculating the optimal values of alpha, buffer and distance ratio based

230    on biodecrypt.wrap results.

231    A plot.biodecrypt function can be applied to the results of biodecrypt to inspect the solution

232    of the analysis and to locate the NUR records. The same function can be applied to the cross-

233    validation results to locate NIR and MIR records.

234

235    *Occurrence data used in this study*

236    As a main source for occurrence data of amalgamated data we used the Distribution Atlas of

237    Butterflies in Europe (Kudrna et al., 2011) which also served as the basis for the calculation

238    of the CLIMBER variables describing climatic preferences of European butterflies

239    (Schweiger et al., 2014). An earlier edition of this atlas (Kudrna, 2002) was also used to

240    generate the Climatic Risk Atlas of European Butterflies (Settele et al., 2008). As a

241    supplementary source of occurrence data for both amalgamated and split species, we used

242    specimens belonging to Roger Vila's collection (Institut de Biologia Evolutiva, Barcelona).

243    The main source for occurrence of cryptic species with known attribution was published data,

244    in most cases represented by genitalia assessment and/or by mitochondrial DNA sequences

245    (Appendix S1 for details). For ten species, a series of 52 specimens from the contact zones

246    have also been specifically sequenced (DNA barcoded) for this study (sequencing methods in

247    Appendix S1) that are included in the "DS-WEUP" BOLD project.

248

249    *Case study species*

250    Following the latest taxonomic assessment for European butterflies (Wiemers et al., 2018),

251    20 cases amalgamated in the Distribution Atlas should be divided into 41 distinct taxa with

252    parapatric or sympatric ranges. Of these, we compiled sufficient identified records for 16

253    cases amalgamated in the above-mentioned atlas, which represent 33 cryptic species. These

254    species were used in this study : 1) *Carcharodus alceae* and *C. tripolinus*, which co-exist in

255    southern Iberia (Dincă et al., 2015); 2) *Spialia sertorius* and *S. rosae*, largely sympatric in

256    Iberia (Hernández-Roldán et al., 2016), 3) *Pyrgus malvae* and *P. malvoides*, parapatric with a

257    contact along central France and the Alps (Koren, Beretta, Črne, & Verovnik, 2013; Litman

258    et al., 2018), 4) *Iphiclides podalirius* and *I. feisthamelii*, parapatric with a contact zone in the

259    Pyrenees and southern France (Wiemers & Gottsberger, 2010; Gaunet et al., 2019); 5)

260    *Zerynthia polyxena* and *Z. cassandra*, parapatric with contact zone in northern Italy (Zinetti

261    et al., 2013); 6) *Pontia daplidice* and *P. edusa,* parapatric with a contact zone in northwestern

262    Italy (Porter, Wenger, Geiger, Scholl, & Shapiro, 1997); 7) *Leptidea sinapis/reali/juvernica*,

263   with *L. reali* and *L. juvernica* being allopatric, but each sympatric with respect to *L. sinapis*

264   (Dincă, Lukhtanov et al., 2011, Dinca et al., 2013); 8) *Lycaena tityrus* and *L. bleusei*

265   parapatric in Iberia (Dincă et al., 2015); 9) *Polyommatus icarus* and *P. celina*, parapatric with

266   contact in southern Iberia (Dincă, Dapporto, & Vila, 2011); 10) *Lysandra coridon* and *L.*

267   *caelestissima*, parapatric in central Iberia (Talavera, Lukhtanov, Rieppel, Pierce, & Vila,

268   2013); 11) *Melitaea athalia* and *M. celadussa* parapatric with a contact zone from southern

269   France through the Alps. We also applied the assignment procedure to a series of species

270   showing almost complete allopatric distribution that were split in Wiemers et al. (2018) but

271   not yet considered in the CLIMBER dataset: 12) *Aglais urticae* and *A. ichnusa*, 13) *Iolana*

272   *iolas* and *I. debilitata,* 14) *Pseudochazara anthelea* and *P. amalthea.* For two largely

273   allopatric species: 15) *Erebia hispania* and *E. rondoui* and 16) *Zizeeria knysna* and *Z.*

274   *karsandra*, applying the procedure was pointless because of their clear allopatry (Appendix

275   S1), but we separated their records and re-calculated the CLIMBER variables (see below). In-

276   depth descriptions of the markers used to generate the subset of identified records of each

277   species are provided in the Appendix S1. We excluded four species groups from the study

278   because knowledge regarding their distribution was still incomplete. This is caused by

279   uncertainty in the identification of records due to the absence of unequivocal morphological

280   markers, sharing of DNA barcodes, and/or to their insufficiently assessed distribution

281   (*Hipparchia semele/neapolitana/blachieri*, *Pieris napi/balcana*, *Lycaena hippothoe/candens,*

282   *Melitaea phoebe/ornata*).

283   To identify the best combination of alpha, buffer and distance ratio parameters, we ran the

284   biodecrypt.wrap function for each species in 80 possible combinations of four alpha values

285   (1, 5, 10, 15), four distance ratio values (2, 3, 4, 5) and five buffer values (0, 40, 80, 120,

286   160).

287

288    *Dependency of attribution on range overlap and on parameters*

289    The most problematic cases are represented by records in or close to areas of sympatry

290    between species, which cannot be attributed with confidence. To verify the effect of the

291    degree of sympatry, we correlated the percentage of MIR, NIR and NUR and of optimised

292    distance ratio, alpha and buffer values with the percentage of the overlapping area between

293    species. The significance of the correlations was tested with Spearman tests. We also

294    evaluated the relationship of MIR, NIR and NUR from the three different parameters by

295    using Generalized Additive Mixed Models. We collated the output of the wrap analyses of

296    the 12 species to which biodecrypt.wrap was applied (960 biodecrypt.cross analyses) and

297    modelled MIR, NUR and NIR in three separate analyses against smoothed (k=2) alpha, ratio

298    and buffer using species as a random factor.

299

300    *Evaluation of niche overlap among cryptic taxa and calculation of CLIMBER variables*

301    To evaluate the potential impact of separation of cryptic taxa on macroecological studies, we

302    evaluated climatic niche overlap among the taxa separated in this study. We used an approach

303    based on a PCA of the climate space in Europe and a density smoothing of the occurrence

304    points for each target species within this space (Broennimann et al., 2012). This is followed

305    by the calculation of niche overlap based on Schoener's D (Schoener, 1968) and the modified

306    Hellinger metric I (Appendix S1 for details). Both indices range from 0 to 1, indicating no

307    niche overlap (0) to full overlap (1). We verified whether the observed overlap is

308    significantly different among separated taxa as done by Warren, Glor, & Turelli (2008). We

309    performed two one-sided tests based on a randomised null model approach, one for niche

310    conservatism and one for niche divergence by testing (i) niche equivalency, i.e. without

311    considering overall available niche space, and (ii) niche similarity, i.e. accounting for the

312    available niche space in Europe (Warren et al., 2008). Analyses were performed with the R

313    package ecospat (Broennimann, Di Cola, & Guisan, 2016) and ade4 (Dray & Dufour, 2007).

314    We also recalculated the CLIMBER variables for the improved distribution data of 33 cryptic

315    taxa following the same procedure applied by Schweiger et al. (2014). For each cryptic

316    group, the former values of mean temperature and precipitation in the distribution area for the

317    amalgamated species complex and for the separated cryptic species were plotted in bivariate

318    plots, together with the values of all the other European species, to illustrate divergence in

319    climatic preferences among cryptic taxa. In order to assess if cryptic pairs have smaller or

320    larger differences in mean temperature and precipitation compared to their congeneric

321    species, we proceeded as follows: We scaled and centered the complete dataset of CLIMBER

322    for mean temperature and precipitation (mean=0, sd=1). Within each genus of the 16

323    complexes examined we calculated the Euclidean distances in scaled mean temperature and

324    precipitation between all pairs of congeneric species included in CLIMBER. Then, we

325    compared the distances between pairs of cryptic taxa separated in this study and between all

326    the other congeneric taxa by using a generalized linear mixed model with a gamma family

327    distribution (glmer function of the lme4 R package) (Bates et al., 2015), including genus as a

328    random     factor.     Script     and     data     are     uploaded     in     Dryad

329    (https://doi.org/10.5061/dryad.hmgqnk9dh).

331    **Results**

333    Cross-validation analyses obtained by biodecrypt.wrap identified the models under a penalty

334    threshold of 10 compared to the model with the lowest penalty, and allowed setting of

335    appropriate alpha, buffer and distance ratio for each species (Fig. 1, Table 1). Compared to

336    the relatively large range of the three parameters tested, the optimization by penalty produced

337    similar values for all the species, with optimal ratio ranging from 2.1 to 3.2, alpha from 5.6 to

338    11.0 and buffer from 31.7 to 96.1 (Table 1). When using the optimised alpha values, the

parapatric taxa revealed very limited overlap, with intersections among the hulls usually

lower than 5% (Table 1). All species showed very low values of misidentified records (MIR),

which were at most 2.0% in the case of *Z. polyxena/cassandra* and *I. podalirius/feisthamelii*.

The percentage of non-attributed identified records (NIR) was much higher and ranged

between 1.2 and 84.7%, with very high values in sympatric species (*Spialia* and *Leptidea*

groups, Fig. 3; Table 1) because, based on the assignment algorithm, all the test records

belonging to the overlap areas cannot be attributed in a cross-validation analysis. The

percentage of non-attributed unknown records (NUR) also varied considerably but, except for

the sympatric *Leptidea* (71.1% NUR) it did not exceed 20.0% of the unknown records (Fig.

3; Table 1). The percentage NIR and NUR correlated with the percentage of area of overlap

(Spearman test: NIR, Rho=0.932, P<0.001; NUR, Rho=0.587, P=0.027; Supplementary

Figure S1). Conversely, the percentage MIR, and the optimised parameters revealed no

correlations with the area of overlap (Spearman test: MIR, Rho=0.432, P=0.160; alpha,

Rho=0.030, P=0.919; buffer, Rho=0.048, P=0.870; distance ratio, Rho=-0.073, P=0.804;

Supplementary Figure S1).

When the different solutions obtained by biodecrypt.wrap for all species were compared in

GAMM analyses it emerged that ratio had no significant effect on MIR, while it strongly

increased NIR and NUR with an almost linear trend (Fig. 4 and effective degrees of freedom

close to 1 in Table 2). Increasing the buffer had a strong effect in reducing MIR, but it also

increased NUR and NIR (with a strongly curvilinear effect for MIR flattening around 100

km), while high values of alpha reduced the number of MIR and NUR and slightly increased

NIR (Fig. 4; Table 2).

The analysis of niche overlap (Table 3) showed that climatic niches are more similar than

expected by chance for all species pairs, as indicated by non-significant divergence for the

niche similarity test (SDD and SDI in Table 3), but only *S. sertorius/rosa*e and *L.*

364  *sinapis/reali* also showed a significant conservatism in terms of niche similarity (SCD and

365  SCI in Table 3). However, significant divergence for the niche equivalency tests indicated

366  that the climatic niches differ considerably for all species pairs (EDD and EDI in Table 3).

367  Taken together, these results show that the climatic niches of all species pairs are

368  significantly different, but still more similar than expected by a random distribution across

369  Europe. These results were consistent for both measures of niche overlap, D and I.

370  CLIMBER variables were calculated for the 33 cryptic species and are available in Appendix

371  S2. When the mean temperature and precipitation for the formerly amalgamated species were

372  plotted together with the data of the newly separated cryptic taxa, it became obvious that the

373  values for the separated species diverge considerably, in some cases spreading all across the

374  main sector of the space occupied by most European species (Fig. 5a-f). A Generalised

375  Linear Mixed Model revealed that differences in mean temperature and precipitation among

376  the cryptic taxa separated here are not smaller than the differences among all congeneric

377  species included in CLIMBER (cryptic taxa, mean difference=1.29=±0.63; congeneric taxa

378  mean difference=1.68±1.21, Estimate=-0.020, Standard Error=0.121, df=1825, t=-0.161,

379  P=0.872).

380

381  **Discussion**

382  Recently diverged taxa (Pigot & Tobias, 2013), as well as cryptic species of butterflies (Vodă

383  et al., 2015; Dapporto et al., 2017; Scalercio et al., 2020), tend to have parapatric

384  distributions with narrow contact zones. This phenomenon facilitates the reliable attribution

385  of ambiguous records to newly discovered species following a hull-based procedure based on

386  the distribution of a subset of identified records. After the parameters used to build the hulls

387  and to attribute records to species had been optimised in a series of cross-validation analyses,

388  the procedure showed a very high attribution of identified records to their correct species

389   (MIR always lower or equal to 2%). Moreover, following this procedure, the fraction of

390   unknown records which remained non-attributed (NUR) was at most 20% for parapatric taxa.

391   The number of incorrectly attributed specimens does not considerably increase with the

392   degree of sympatry. However, the values of optimised parameters did not depend on the

393   degree of overlap, and they are more likely imposed by the geometry of the distribution areas.

394   The comparison of a large set of cross validation analyses (GAMM approach) revealed how

395   the parameters can be varied to impose stricter or wider inclusion possibilities for a given

396   taxon. Theoretically, higher distance ratio and buffer, and lower alpha values are expected to

397   decrease the number of incorrectly attributed records, but to increase the number of non-

398   attributed records. Accordingly, increasing buffer strongly reduced MIR but increased NIR

399   and NUR (Fig. 4), while the ratio and the alpha values had much lighter effects on MIR. This

400   is probably due to the relatively simple (not interdigitated) shape of the distribution

401   boundaries among cryptic species of butterflies for which a high alpha value (producing

402   slightly concave hulls) and an optimized distance from other hulls (buffer) can avoid most

403   MIR. These observations can help in refining the solutions. In the genus *Iphiclides*,

404   unavailability of identified specimens from northern areas did not allow to attribute

405   northernmost unidentified records (most probably belonging to *I. podalirius*) and this could

406   have slightly affected the assessment of climatic niches. When including a few misclassified

407   specimens does cause more problems than generating several non-attributed unrecognized

408   record (NUR), a computational strategy could be excluding the squared exponent of MIR

409   from the penalty calculation in biodecrypt.optimise (e.g. MIR+NIR+NUR). This can reduce

410   the maximum ratio for attribution, thus lowering the number of NUR with a very small

411   increase in MIR (Fig. 4). Finally, the values of the procedure parameters can also be adjusted

412   according to researcher's perception.

413    The framework presented here provides a standardised method that allows researchers to take

414    advantage of previous efforts of data collection even when modern investigation methods

415    identify new layers of biodiversity. This is not necessarily restricted to cryptic species but can

416    also be relevant for analyses of evolutionary significant units such as those highlighted by

417    species delimitation methods (Lecocq et al., 2019). Moreover, biodecrypt can also be used in

418    cases where very similar species are treated as a single unit in citizen science projects.

419    Clearly, in the presence of largely sympatric species the approach described here is less

420    efficient because all the occurrence data for the overlapping zones must be directly verified.

421    In general, this can happen when climatic distance governs limits between species

422    distributions more than geographic distance, as revealed by the existence of several species

423    with small ranges being restricted to climatically rare areas (Ohlemüller et al., 2008). In our

424    dataset this could be the case of two pairs, *Leptidea sinapis*/*L. reali* and *Spialia sertorius*/*S.*

425    *rosae*. In fact, *L. reali* and *S. rosae* are specialized to mountain areas, while *L. sinapis* and *S.*

426    *sertorius* are climate generalists. Although the distance-based biodecrypt algorithm is less

427    efficient in attributing unknown records in areas where the pairs of species are mostly

428    separated by climatic distances, it remains conservative in avoiding mis-identifications.

429    Moreover, even predominantly sympatric taxa often do not co-occur in large areas of their

430    distributions (Fig. 3) and in these cases, biodecrypt can highlight areas where the attribution

431    to a given taxon is reliable based on distances as opposed to areas where climatic distances

432    might exert an important complementary effect. The detection of poorly studied areas is

433    another notable result stemming from our analysis. In fact, the distribution of non-attributed

434    unknown records can highlight where further research efforts are needed to confirm the

435    distribution of recently recognized species by means of dedicated field research or by

436    morphological examination and DNA sequencing from existing natural history collections

437    (e.g. Kharouba, Lewthwaite, Guralnick, Kerr, & Vellend, 2019). An example of such regions

438    is north-western France for *Pontia daplidice/edusa*, *Melitaea athalia/celadussa*, and *Pyrgus*

439    *malvae/malvoides* (Fig. 3).

440    It has already been demonstrated that considering cryptic taxa as amalgamated entities

441    overlooks a large fraction of beta-diversity (Vodă et al., 2015). Our results also indicate that

442    amalgamating cryptic taxa may affect macroecological studies investigating the response of

443    species and communities to climatic changes (Settele et al., 2008; Devictor et al., 2012).We

444    found significant niche divergence for all the pairs of cryptic species and a particularly low

445    level of niche conservatism (two species pairs out of sixteen). This is reflected in the strong

446    differences in the values of CLIMBER variables, which are widely used in macroecology to

447    assess butterfly responses to current and future climatic conditions. The differences in mean

448    temperature and precipitation experienced by cryptic species pairs in Europe are considerable

449    and in the same order of magnitude as the differences among other congeneric species. In

450    some cases, the variables describing climatic preferences of cryptic taxa were so divergent

451    that they spread all over the space defined by the mean temperature and precipitation of most

452    European species.

453    Failing to consider the divergent climatic niches of cryptic taxa can have considerable

454    impacts on climatic risk assessments based on species distribution modelling. It has been

455    shown, in the case of locally adapted subspecies, that individual responses of taxa to climate

456    change do not necessarily resemble the modelled response of the amalgamated group, which

457    might lead to severe over- or underestimation of the risks (Lecocq et al., 2019). Also,

458    assessments of the response of communities, such as based on the community-weighted mean

459    temperature index, might be improved by considering cryptic species. The replacement of a

460    cool-adapted species by a warm-adapted species of an initially amalgamated group will

461    clearly contribute to an increase in the community-weighted mean temperature index, while

462    ignoring climatic niche divergence within the amalgamated group will be interpreted as no

463    change. Such an effect might, partly, contribute to an observed time lag in the response of

464    butterfly communities to climate change (Devictor et al., 2012). Our dataset comprised 33

465    species, representing 6.7% of the 496 species recorded in Europe (Wiemers et al., 2018).

466    Notably, some of these complexes represent widely distributed and common butterflies in

467    Europe (*Polyommatus icarus/celina*, *Aglais urticae/ichnusa, Carcharodus alceae/tripolinus,*

468    *Melitaea athalia/celadussa, Pontia daplidice/edusa* and *Leptidea sinapis/reali/juvernica*).

469    Given the range, density and the large differences in climatic preferences displayed by

470    cryptic taxa, their separation can be an important improvement for monitoring and analysing

471    community-level responses to climate change. This information can also help in connecting

472    large-scale and long-term evolutionary processes with ecological processes through the

473    analysis of species interactions with their abiotic and biotic environments, or shed light on the

474    phylogenetic conservation of niche characteristics relevant to climate-change (e.g. Devictor

475    et al., 2012; Kharouba et al., 2019; Ohlemüller et al., 2008).

476    In conclusion, it is our hope that the framework here presented, which allows the objective

477    attribution of undetermined occurrence data to recently assessed taxonomic units, will benefit

478    biodiversity mapping, highlight gaps in current knowledge and improve macroecological

479    analyses.

480 **References**

481

482 Bates, D., Maechler, M., Bolker, B., Walker, S., Christensen, R.H.B., Singmann, H., …
483      Green, P. (2015) Package 'lme4.' *Convergence*, **12**, 2.

484 Bickford, D., Lohman, D.J., Sodhi, N.S., Ng, P.K.L., Meier, R., Winker, K., … Das, I. (2007)
485      Cryptic species as a window on diversity and conservation. *Trends in ecology &*
486      *evolution*, **22**, 148–155.

487 Broennimann, O., Di Cola, V. & Guisan, A. (2016) ecospat: Spatial ecology miscellaneous
488      methods. R package version 2.1. 1.

489 Broennimann, O., Fitzpatrick, M.C., Pearman, P.B., Petitpierre, B., Pellissier, L., Yoccoz,
490      N.G., … Zimmermann, N.E. (2012) Measuring ecological niche overlap from
491      occurrence and spatial environmental data. *Global ecology and biogeography*, **21**, 481–
492      497.

493 Dapporto, L., Cini, A., Menchetti, M., Vodă, R., Bonelli, S., Casacci, L.P., … Biermann, H.
494      (2017) Rise and fall of island butterfly diversity: Understanding genetic differentiation
495      and extinction in a highly diverse archipelago. *Diversity and Distributions*, **23**, 1169–
496      1181.

497 Dapporto, L., Cini, A., Vodă, R., Dincă, V., Wiemers, M., Menchetti, M., … Vila, R. (2019)
498      Integrating three comprehensive data sets shows that mitochondrial DNA variation is
499      linked to species traits and paleogeographic events in European butterflies. *Molecular*
500      *Ecology Resources*, **19**, 1623–1636.

501 Dapporto, L., Ramazzotti, M., Fattorini, S., Talavera, G., Vila, R. & Dennis, R.L.H. (2013)
502      recluster: an unbiased clustering procedure for beta-diversity turnover. *Ecography*, **36**,
503      1070–1075.

504 Dennis, E.B., Morgan, B.J.T., Brereton, T.M., Roy, D.B. & Fox, R. (2017) Using citizen
505      science butterfly counts to predict species population trends. *Conservation biology*, **31**,
506      1350–1361.

507 Devictor, V., Van Swaay, C., Brereton, T., Brotons, L., Chamberlain, D., Heliölä, J., …

508 Lindström, Å. (2012) Differences in the climatic debts of birds and butterflies at a
509  continental scale. *Nature climate change*, **2**, 121–124.

510 Dincă, V., Dapporto, L. & Vila, R. (2011a) A combined genetic-morphometric analysis
511  unravels the complex biogeographical history of *Polyommatus icarus* and *Polyommatus*
512  *celina* Common Blue butterflies. *Molecular Ecology*, **20**, 3921–3935.

513 Dincă, V., Lukhtanov, V.A., Talavera, G. & Vila, R. (2011b) Unexpected layers of cryptic
514  diversity in wood white *Leptidea* butterflies. *Nature communications*, **2**, 1–8.

515 Dincă, V., Montagud, S., Talavera, G., Hernández-Roldán, J., Munguira, M.L., García-
516  Barros, E., … Vila, R. (2015) DNA barcode reference library for Iberian butterflies
517  enables a continental-scale preview of potential cryptic diversity. *Scientific Reports*, **5**,
518  12395.

519 Dincă, V., Wiklund, C., Lukhtanov, V.A., Kodandaramaiah, U., Norén, K., Dapporto, L., …
520  Friberg, M. (2013) Reproductive isolation and patterns of genetic differentiation in a
521  cryptic butterfly species complex. *Journal of Evolutionary Biology*, **26**, 2095–2106.

522 Dray, S. & Dufour, A.-B. (2007) The ade4 package: implementing the duality diagram for
523  ecologists. *Journal of statistical software*, **22**, 1–20.

524 Edelsbrunner, H., Kirkpatrick, D. & Seidel, R. (1983) On the shape of a set of points in the
525  plane. *IEEE Transactions on information theory*, **29**, 551–559.

526 Franklin, J. (2010) *Mapping species distributions: spatial inference and prediction*,
527  Cambridge University Press.

528 Gaunet, A., Dincă, V., Dapporto, L., Montagud, S., Vodă, R., Schär, S., … Vila, R. (2019)
529  Two consecutive *Wolbachia*-mediated mitochondrial introgressions obscure taxonomy
530  in Palearctic swallowtail butterflies (Lepidoptera, Papilionidae). *Zoologica Scripta*, **48**,
531  507–519.

532 Hernández-Roldán, J.L., Dapporto, L., Dincă, V., Vicente, J.C., Hornett, E.A., Šíchová, J., …
533  Vila, R. (2016) Integrative analyses unveil speciation linked to host plant shift in *Spialia*
534  butterflies. *Molecular Ecology*, **25**, 4267–4284.

535 Herrando, S., Titeux, N., Brotons, L., Anton, M., Ubach, A., Villero, D., … Stefanescu, C.

(2019) Contrasting impacts of precipitation on Mediterranean birds and butterflies. *Scientific reports*, **9**, 1–7.

Hortal, J., de Bello, F., Diniz-Filho, J.A.F., Lewinsohn, T.M., Lobo, J.M. & Ladle, R.J. (2015) Seven shortfalls that beset large-scale knowledge of biodiversity. *Annual Review of Ecology, Evolution, and Systematics*, **46**, 523–549.

Kharouba, H.M., Lewthwaite, J.M.M., Guralnick, R., Kerr, J.T. & Vellend, M. (2019) Using insect natural history collections to study global change impacts: challenges and opportunities. *Philosophical Transactions of the Royal Society B*, **374**, 20170405.

Koren, T., Beretta, S., Črne, M. & Verovnik, R. (2013) On the distribution of *Pyrgus malvoides* (Elwes & Edwards, 1897) (Lepidoptera: Hesperiidae) at the eastern part of its range. *Entomologist's Gazette*, **64**, 225–234.

Kudrna, O. (2002) The distribution atlas of European butterflies. *Oedippus*, **20**, 1–343.

Kudrna, O., Harpke, A., Lux, K., Pennerstorfer, J., Schweiger, O., Settele, J. & Wiemers, M. (2011) *Distribution atlas of butterflies in Europe*, Gesellschaft für Schmetterlingsschutz eV Halle.

Kudrna, O., Pennerstorfer, J. & Lux, K. (2015) *Distribution atlas of European butterflies and skippers*, Peks.

Lecocq, T., Harpke, A., Rasmont, P. & Schweiger, O. (2019) Integrating intraspecific differentiation in species distribution models: Consequences on projections of current and future climatically suitable areas of species. *Diversity and Distributions*, **25**, 1088–1100.

Litman, J., Chittaro, Y., Birrer, S., Praz, C., Wermeille, E., Fluri, M., … Gonseth, Y. (2018) A DNA barcode reference library for Swiss butterflies and forester moths as a tool for species identification, systematics and conservation. *PLoS ONE*, **13**.

Ohlemüller, R., Anderson, B.J., Araújo, M.B., Butchart, S.H.M., Kudrna, O., Ridgely, R.S. & Thomas, C.D. (2008) The coincidence of climatic and species rarity: high risk to small-range species from climate change. *Biology letters*, **4**, 568–572.

Pigot, A.L. & Tobias, J.A. (2015) Dispersal and the transition to sympatry in vertebrates.

564        *Proceedings of the Royal Society B: Biological Sciences*, **282**, 20141929.

565    Pigot, A.L. & Tobias, J.A. (2013) Species interactions constrain geographic range expansion
566        over evolutionary time. *Ecology letters*, **16**, 330–338.

567    Porter, A.H., Wenger, R., Geiger, H., Scholl, A. & Shapiro, A.M. (1997) The *Pontia*
568        *daplidice-edusa* hybrid zone in northwestern Italy. *Evolution*, **51**, 1561–1573.

569    R Core Team (2019) R: a language and environment for statistical computing, version 3.0. 2.
570        Vienna, Austria: R Foundation for Statistical Computing; 2013.

571    Scalercio, S., Cini, A., Menchetti, M., Vodă, R., Bonelli, S., Bordoni, A., … Vila, R. (2020)
572        How long is 3 km for a butterfly? Ecological constraints and functional traits explain
573        high mitochondrial genetic diversity between Sicily and the Italian Peninsula. *Journal of*
574        *Animal Ecology*.

575    Schmitt, T. (2007) Molecular biogeography of Europe: Pleistocene cycles and postglacial
576        trends. *Frontiers in zoology*, **4**, 11.

577    Schoener, T.W. (1968) The *Anolis* lizards of Bimini: resource partitioning in a complex
578        fauna. *Ecology*, **49**, 704–726.

579    Schweiger, O., Harpke, A., Wiemers, M. & Settele, J. (2014) CLIMBER: Climatic niche
580        characteristics of the butterflies in Europe. *ZooKeys*, 65.

581    Schweiger, O., Heikkinen, R.K., Harpke, A., Hickler, T., Klotz, S., Kudrna, O., … Settele, J.
582        (2012) Increasing range mismatching of interacting species under global change is
583        related to their ecological characteristics. *Global Ecology and Biogeography*, **21**, 88–99.

584    Settele, J., Kudrna, O., Harpke, A., Kühn, I., Van Swaay, C., Verovnik, R., … Hickler, T.
585        (2008) *Climatic risk atlas of European butterflies*, Pensoft Sofia.

586    Talavera, G., Lukhtanov, V.A., Rieppel, L., Pierce, N.E. & Vila, R. (2013) In the shadow of
587        phylogenetic uncertainty: the recent diversification of *Lysandra* butterflies through
588        chromosomal change. *Molecular Phylogenetics and Evolution*, **69**, 469–478.

589    Thuiller, W., Georges, D., Engler, R., Breiner, F., Georges, M.D. & Thuiller, C.W. (2016)
590        Package 'biomod2'. Species distribution modeling within an ensemble forecasting

591 framework. *Ecography*, **32**, 369–373.

592 Titeux, N., Maes, D., Van Daele, T., Onkelinx, T., Heikkinen, R.K., Romo, H., … van
593 Swaay, C.A.M. (2017) The need for large-scale distribution data to estimate regional
594 changes in species richness under future climate change. *Diversity and Distributions*, **23**,
595 1393–1407.

596 Toews, D.P.L., Mandic, M., Richards, J.G. & Irwin, D.E. (2014) Migration, mitochondria,
597 and the yellow-rumped warbler. *Evolution*, **68**, 241–255.

598 Vodă, R., Dapporto, L., Dincă, V. & Vila, R. (2015) Cryptic matters: overlooked species
599 generate most butterfly beta-diversity. *Ecography*, **38**, 405–409.

600 Warren, D.L., Glor, R.E. & Turelli, M. (2008) Environmental niche equivalency versus
601 conservatism: quantitative approaches to niche evolution. *Evolution: International
602 Journal of Organic Evolution*, **62**, 2868–2883.

603 Waters, J.M. (2011) Competitive exclusion: phylogeography's 'elephant in the room'?
604 *Molecular Ecology*, **20**, 4388–4394.

605 Wiemers, M., Balletto, E., Dincă, V., Fric, Z.F., Lamas, G., Lukhtanov, V., … Verovnik, R.
606 (2018) An updated checklist of the European Butterflies (Lepidoptera, Papilionoidea).
607 *ZooKeys*, **811**, 9–45.

608 Wiemers, M. & Gottsberger, B. (2010) Discordant patterns of mitochondrial and nuclear
609 differentiation in the Scarce Swallowtail *Iphiclides podalirius feisthamelii* (Duponchel,
610 1832)(Lepidoptera: Papilionidae). *Entomologische Zeitschrift*, **120**, 111–115.

611 Zinetti, F., Dapporto, L., Vovlas, A., Chelazzi, G., Bonelli, S., Balletto, E. & Ciofi, C. (2013)
612 When the rule becomes the exception. No evidence of gene flow between two *Zerynthia*
613 cryptic butterflies suggests the emergence of a new model group. *PLoS ONE*, **8**.

614 Zografou, K., Kati, V., Grill, A., Wilson, R.J., Tzirkalli, E., Pamperis, L.N. & Halley, J.M.
615 (2014) Signals of climate change in butterfly communities in a Mediterranean protected
616 area. *PLoS ONE*, **9**.

617

618    Data Accessibility Statement:

619

620    Occurrence data and R scripts to replicate the analysis are available in Dryad

621    https://doi.org/10.5061/dryad.hmgqnk9dh.

622    Newly generated COI data are available in the "DS-WEUP" BOLD project.

623    The biodecrypt functions are available in the recluster R package and in https://cran.r-

624    project.org/web/packages/recluster/index.html                and             in             Dryad

625    https://doi.org/10.5061/dryad.hmgqnk9dh.

626

627    Table 1. Results of cross-validation, parameters estimated by optimization and range overlap.

628

| Species | MIR | NIR | NUR | Alpha | Buffer | Ratio | Overlap% |
|---|---|---|---|---|---|---|---|
| *S. sertorius/rosae* | 1.9 | 41.2 | 18.2 | 11.0 | 96.1 | 2.3 | 15.6 |
| *P. malvae/malvoides* | 0.0 | 33.6 | 7.6 | 9.7 | 68.7 | 2.4 | 2.3 |
| *C. alceae/tripolinus* | 0.5 | 35.2 | 5.0 | 7.3 | 92.4 | 2.4 | 4.4 |
| *Z. polyxena/cassandra* | 2.0 | 12.0 | 7.6 | 5.6 | 35.3 | 2.3 | 0.1 |
| *I. podalirius/feisthamelii* | 2.0 | 9.1 | 6.2 | 7.7 | 39.0 | 2.1 | 0.2 |
|  |  |  |  |  |  |  |  |
| *P. daplidice/edusa* | 0.6 | 4.3 | 3.9 | 9.0 | 50.5 | 2.6 | 0.0 |
| *L. sinapis/reali/juvernica* | 1.9 | 84.7 | 71.1 | 8.5 | 52.5 | 3.2 | Lr-Lj 0.0 |
|  |  |  |  |  |  |  | Lr-Ls 4.1 |
|  |  |  |  |  |  |  | Ls-Lj 50.1 |
| *L. tityrus/bleusei* | 0.0 | 2.7 | 2.0 | 9.2 | 74.0 | 2.9 | 0.0 |
| *I. iolas/debilitata* | 0.0 | 4.5 | 1.5 | 7.6 | 80.5 | 2.4 | 0.0 |
| *P. icarus/celina* | 0.7 | 11.4 | 3.1 | 8.5 | 52.5 | 3.2 | 1.9 |
| *L. coridon/caelestissima* | 0.0 | 1.2 | 0.3 | 8.9 | 76.4 | 2.7 | 0.0 |
|  |  |  |  |  |  |  |  |
| *M. athalia/celadussa* | 1.9 | 10.2 | 7.2 | 10.4 | 31.7 | 2.7 | 1.0 |
| *A. urticae/ichnusa* | NA | NA | 0.2 | 8.0 | 60.0 | 2.5 | 0.0 |
| *P. anthelea/amalthea* | NA | NA | 20.0 | 8.0 | 60.0 | 2.5 | 0.0 |

629

630    Table 1. Misidentified identified records (MIR), non-attributed identified records (NIR) and

631    non-attributed non-identified records (NUR) obtained for each species after optimization of

632    alpha, buffer and distance ratio parameters. The parameters have been estimated with series

633    of cross-validation analysis. The alpha, buffer and ratio of the solutions showing a penalty

634    (MIR$^2$+NIR+NUR) not higher than 10 compared to the model with lowest penalty have been

635    averaged weighted by their inverse penalty by biodecrypt.optimise. The percentage of range

636    overlap among hulls in the solution obtained by using these parameters is also reported. For

637    *Leptidea*, the pairwise overlap values between *L. sinapis* (Ls), *L. reali* (Lr) and *L. juvernica*

638    (Lj) are provided.

639

640    Table 2. Generalized Additive Mixed Model (GAMM) results for the influence of the three

641    parameters on record attribution.

642

| | MIR | | | NIR | | | NUR | | |
|---|---|---|---|---|---|---|---|---|---|
| | edf | F | P | edf | F | P | edf | F | P |
| Ratio | 1 | 0.486 | 0.486 | 1 | 219.2 | <0.001 | 1.168 | 933.36 | <0.001 |
| Buffer | 1.992 | 430.06 | <0.001 | 1.546 | 832.8 | <0.001 | 1.825 | 186.48 | <0.001 |
| Alpha | 1.611 | 16.972 | <0.001 | 1.54 | 11.7 | <0.001 | 1.967 | 95.09 | <0.001 |

643

644    Table 2. The effect of ratio, buffer and alpha of misidentified records (MIR), non-attributed

645    identified records (NIR) and non-attributed unidentified records (NUR) as verified by

646    GAMMs (edf, effective degrees of freedom, F and P values are provided).

647

648    Table 3. Tests for niche equivalency, divergence and conservatism.

649

| Group | D | I | ECD P | ECI P | EDD P | EDI P | SCD P | SCI P | SDD P | SDI P |
|---|---|---|---|---|---|---|---|---|---|---|
| *C. alceae/tripolinus* | 0.179 | 0.367 | 1.000 | 1.000 | **0.010** | **0.010** | 0.069 | 0.089 | 0.881 | 0.842 |
| *C. sertorius/rosae* | 0.155 | 0.376 | 1.000 | 1.000 | **0.010** | **0.010** | **0.030** | **0.030** | 0.950 | 0.960 |
| *P. malvae/malvoides* | 0.252 | 0.429 | 1.000 | 1.000 | **0.010** | **0.010** | 0.178 | 0.208 | 0.772 | 0.762 |
| *I. podalirius/feisthamelii* | 0.351 | 0.485 | 1.000 | 1.000 | **0.010** | **0.010** | 0.248 | 0.267 | 0.752 | 0.743 |
| *Z. polyxena/cassandra* | 0.341 | 0.538 | 1.000 | 1.000 | **0.010** | **0.010** | 0.059 | 0.069 | 0.931 | 0.931 |
| *P. daplidice/edusa* | 0.481 | 0.608 | 1.000 | 1.000 | **0.010** | **0.010** | 0.079 | 0.119 | 0.960 | 0.931 |
| *L. sinapis/reali* | 0.274 | 0.507 | 1.000 | 1.000 | **0.020** | **0.030** | 0.069 | **0.050** | 0.970 | 0.970 |
| *L. sinapis/juvernica* | 0.266 | 0.491 | 1.000 | 1.000 | **0.010** | **0.010** | 0.129 | 0.079 | 0.782 | 0.861 |
| *L. reali/juvernica* | 0.179 | 0.303 | 1.000 | 1.000 | **0.010** | **0.010** | 0.208 | 0.178 | 0.832 | 0.861 |
| *L. tityrus/bleusei* | 0.131 | 0.326 | 0.990 | 1.000 | **0.010** | **0.010** | 0.168 | 0.158 | 0.812 | 0.832 |
| *I. iolas/debilitata* | 0.255 | 0.429 | 1.000 | 1.000 | **0.010** | **0.010** | 0.149 | 0.168 | 0.901 | 0.871 |
| *P. icarus/celina* | 0.115 | 0.294 | 1.000 | 1.000 | **0.010** | **0.010** | 0.347 | 0.347 | 0.634 | 0.653 |
| *L. coridon/caelestissima* | 0.012 | 0.104 | 1.000 | 1.000 | **0.010** | **0.010** | 0.238 | 0.248 | 0.703 | 0.693 |
| *M. athalia/celadussa* | 0.271 | 0.429 | 1.000 | 1.000 | **0.010** | **0.010** | 0.257 | 0.317 | 0.723 | 0.703 |
| *A. urticae/ichnusa* | 0.054 | 0.226 | 1.000 | 1.000 | **0.010** | **0.010** | 0.267 | 0.238 | 0.723 | 0.812 |
| *P. anthelea/amalthea* | 0.060 | 0.120 | 1.000 | 1.000 | **0.010** | **0.010** | 0.059 | 0.277 | 0.941 | 0.752 |

650

651    Table 3. Niche overlap metrics Shoener's D (D) and the modified Hellinger metric I (I)

652    calculated in multivariate climatic space for pairs of the studied taxa. Four tests have been

653    performed per metric and respective P values are provided for conservatism based on niche

654    equivalency for D and I (ECD P; ECI P), divergence based on niche equivalency (EDD P;

655    EDI P), conservatism based on niche similarity (SCD P, SCI P) and divergence based on

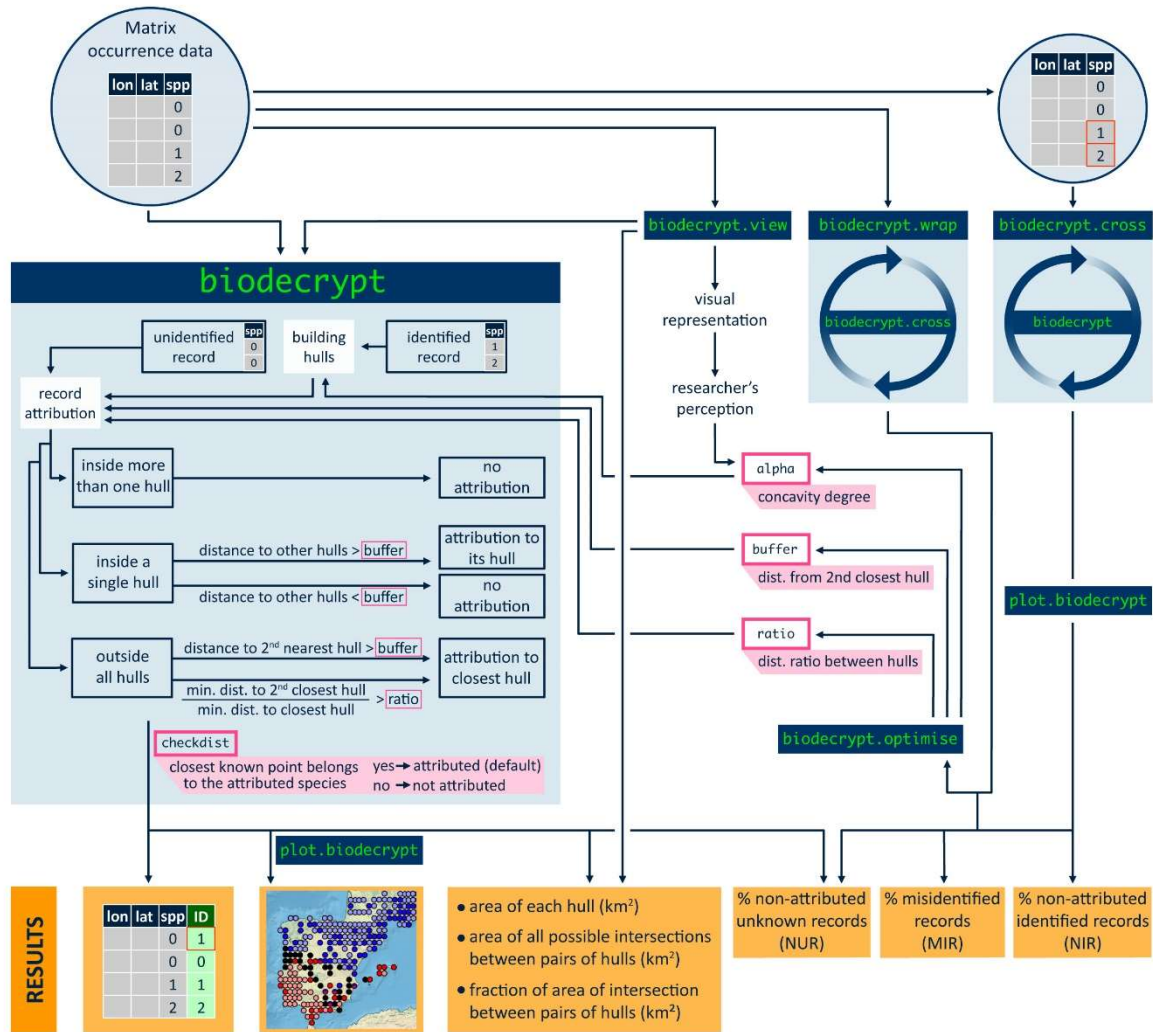656    niche similarity (SDD P, SDI P). Significant values are indicated in bold.

657

Figure 1. A diagram showing the workflow of the six functions of the biodecrypt family. Function names are indicated as blue boxes. Circles with arrows indicate wrapping of a function inside another one. Options/Parameters are indicated as pink boxes and a short description is given. Results are shown in orange boxes.
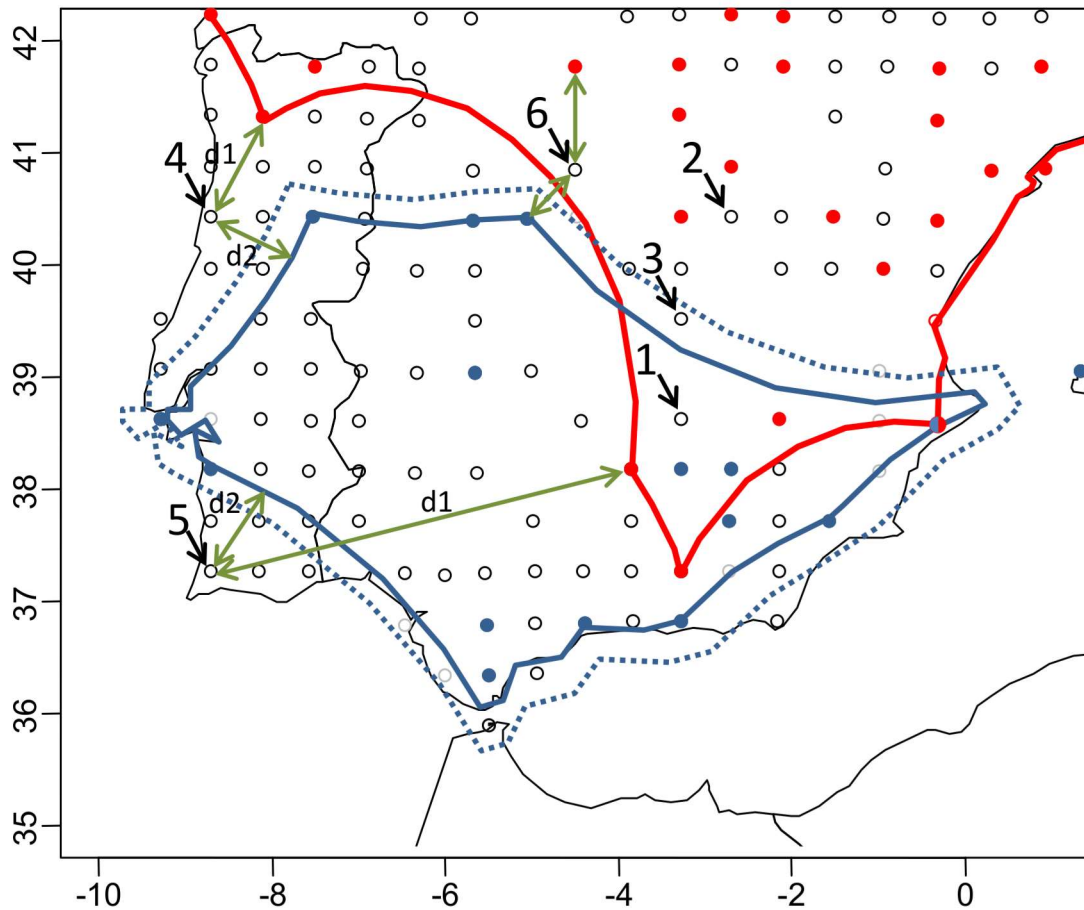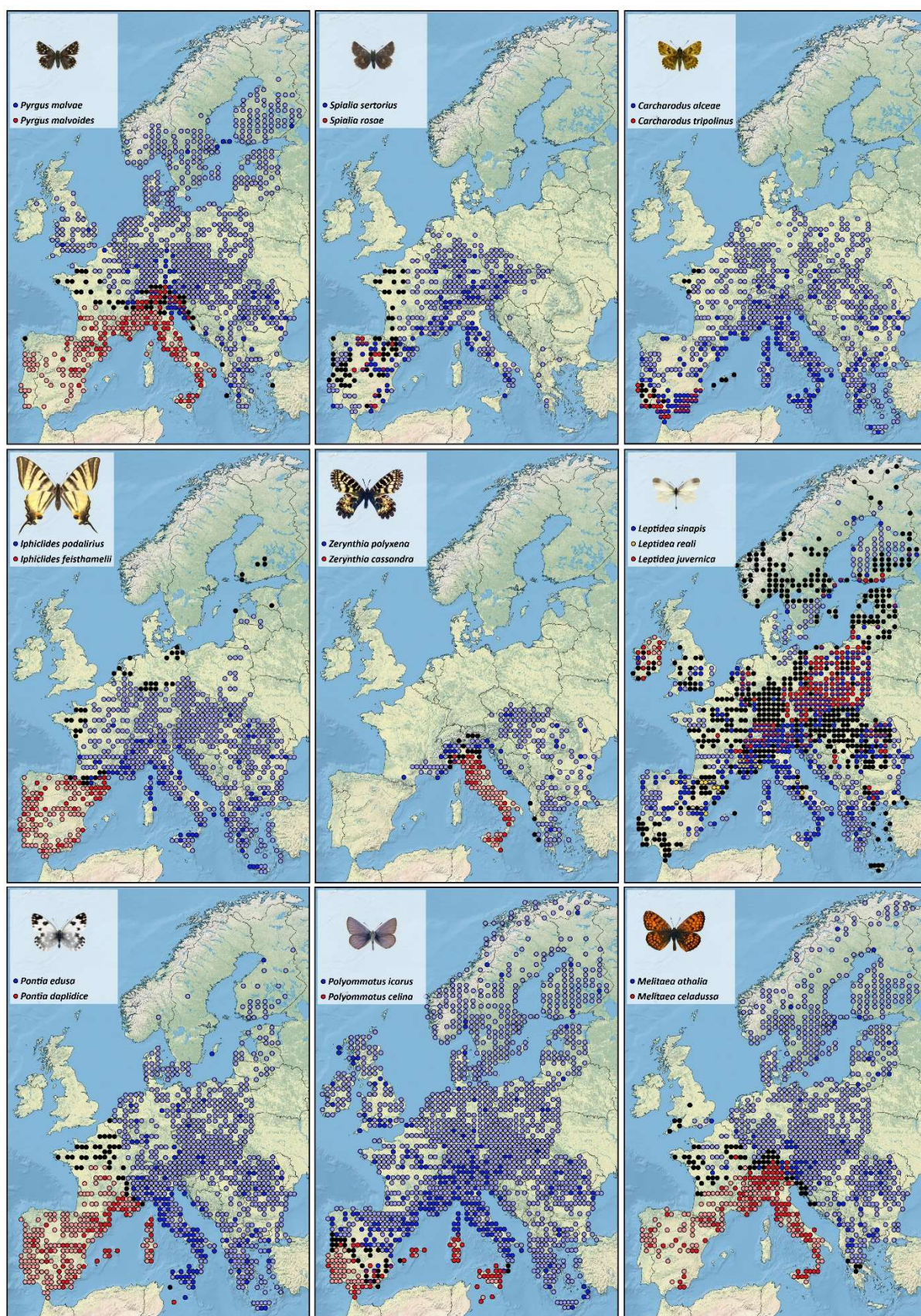
664

665

666 Figure 2. An example of the assignment procedure based on the overlapping distributions of

667 *Polyommatus icarus* (red) and *P. celina* (blue) in Iberia. Red and blue dots represent the sites

668 from where specimens of the respective species have been sequenced, empty circles represent

669 sites with unidentified records. The continuous red and blue lines represent the hulls obtained

670 for the two species based on occurrence of sequenced specimens (alpha = 3 in this example

671 compared to the optimized value of 8.5 to show the effect of a low alpha value on polygon

672 shape), the dotted blue line represents the buffer of the *P. celina* hull (for clarity, the buffer of

673 the *P. icarus* hull is not represented). Unidentified record 1 is not attributed since it falls

674 inside both hulls; record 2 is attributed to *P. icarus*, record 3 is not attributed since it falls in

675 the buffer of the *P. celina* hull; record 4 (external to both hulls) is not attributed because the

676 distance ratio (d1/d2; where d1 represents the larger distance) is smaller than the default ratio

677    of 2.5; while record 5 is attributed to *P. celina* because d1/d2 is larger than 2.5 (default in

678    biodecrypt). If checkdist is selected, point 6 is not attributed to *P. icarus* because the closest

679    known record is of *P. celina*.

680

681

682

683　Figure 3. Distribution maps obtained for nine groups of cryptic species by using the

684　optimised parameters reported in Table 1. Dots with darker colour represent records

685     identified based on DNA sequences (mostly mitochondrial), genitalic morphology or other

686     markers. Circles with paler colours represent records attributed by the algorithm. Black

687     circles are non-attributed records. The physical map used is freely available from Natural
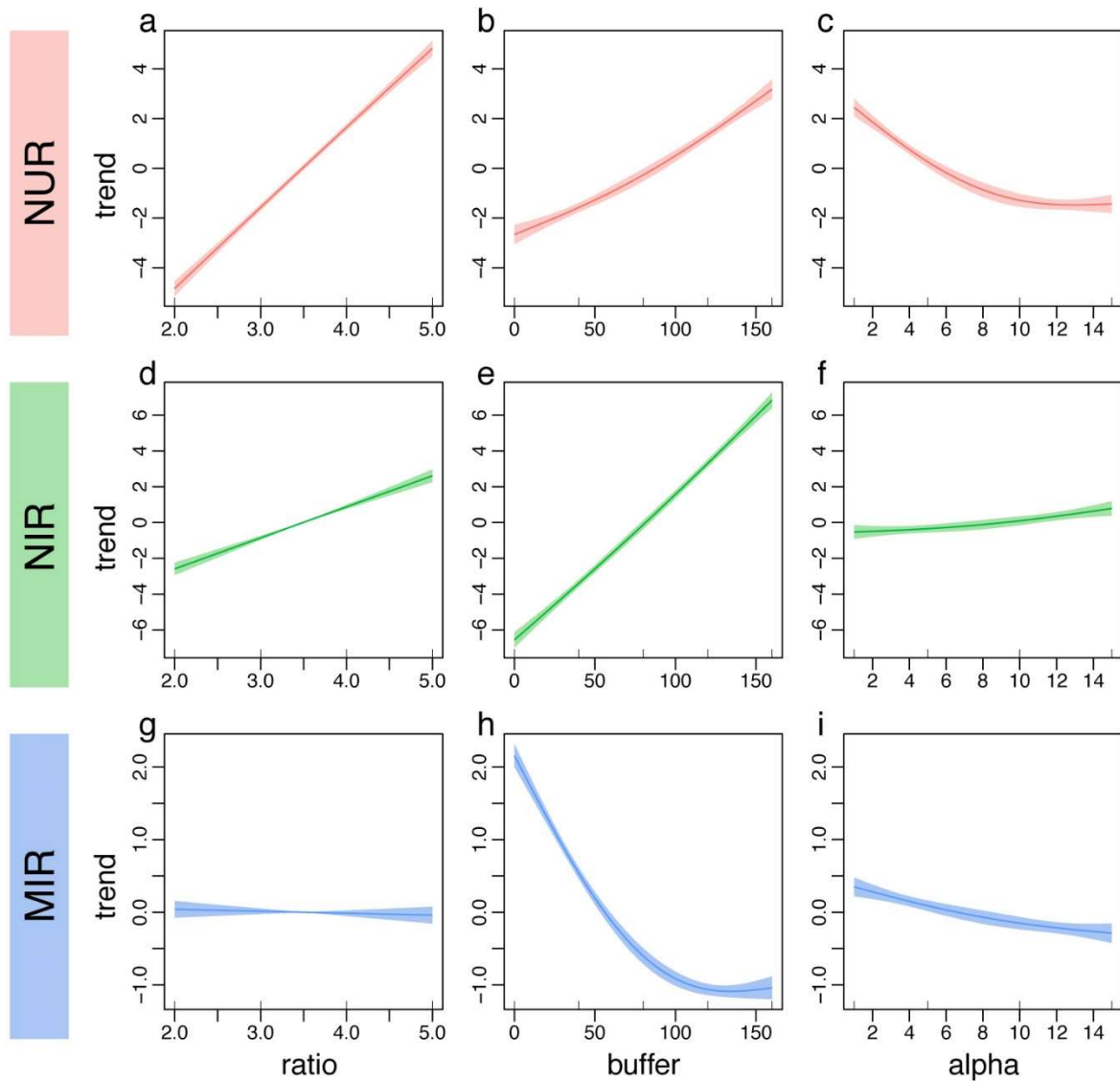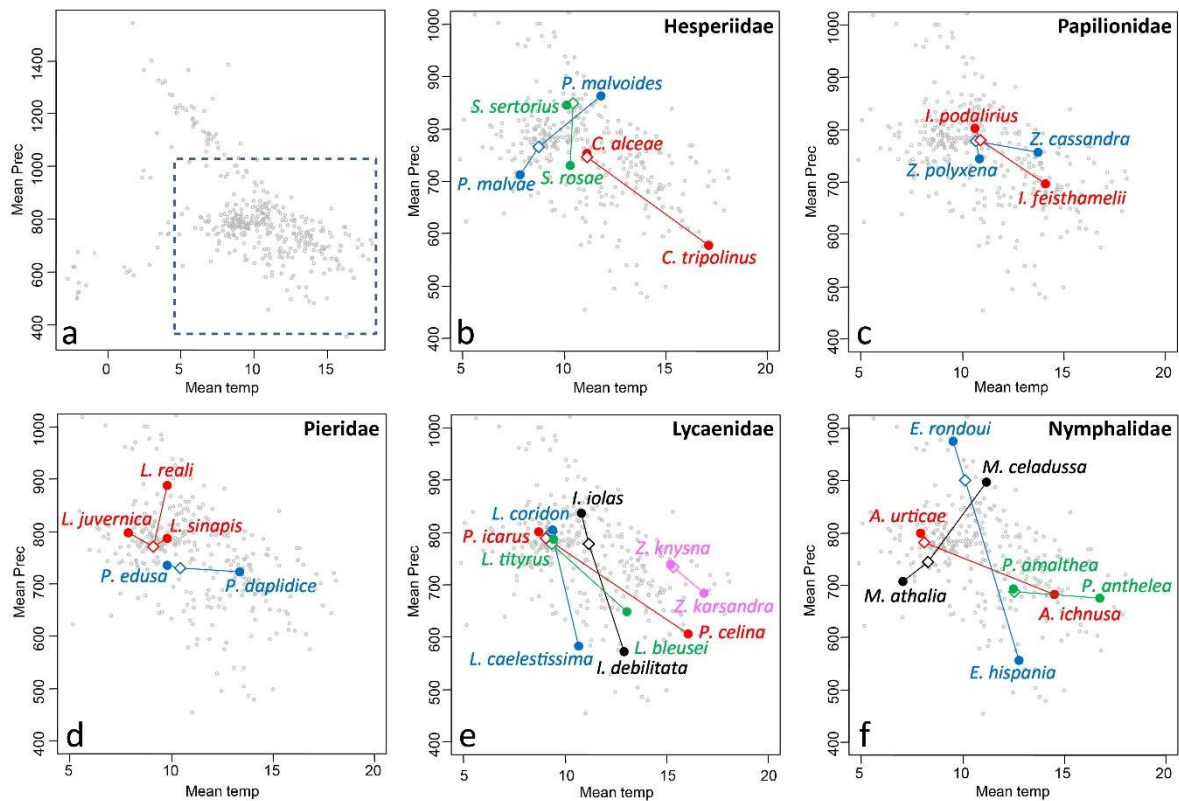
688     Earth (www.naturalearthdata.com).

689

Figure 4. The standardized effect (trends) of ratio, buffer and alpha obtained after Generalized Additive Mixed Model on non-attributed unknown records (NUR), non-attributed known records (NIR) and misidentified known records (MIR). Data belongs to series of biodecrypt.wrap analyses.

698

Figure 5. Climatic space defined by mean annual temperature and annual precipitation for all European species included in the CLIMBER dataset (a). The blue rectangle represents the climate space where the species under study are located and which is used for figures b-f. For the five families, the resulting values for the cryptic taxa separated in this study are represented by filled dots connected with a line to the former values of the amalgamated taxon (diamonds). Each cryptic complex within a family is represented by a different colour (b-f) and separated cryptic species names are reported.