

1 **BAGS: an automated Barcode, Audit & Grade System for DNA barcode reference**
2 **libraries.**

3

4 João T Fontes^{1,2+}, Pedro E Vieira^{1,2+}, Torbjørn Ekrem³, Pedro Soares^{1,2}, Filipe O Costa^{1,2}

5 ¹ CBMA - Centre of Molecular and Environmental Biology, Department of Biology, University
6 of Minho, Campus de Gualtar, 4710-057 Braga, Portugal

7 ² Institute of Science and Innovation for Bio-Sustainability (IB-S), University of Minho, Portugal

8 ³ Department of Natural History, NTNU University Museum, Trondheim, Norway

9 + João T Fontes and Pedro E Vieira are joint first authors

10 **Correspondence**

11 Filipe O Costa

12 Email: fcosta@bio.uminho.pt

13 <https://orcid.org/0000-0001-5398-3942>

14

15 **Running head**

16 Auditing and annotation of DNA barcodes.

17 **Abstract**

18 Biodiversity studies greatly benefit from molecular tools, such as DNA metabarcoding, which
19 provides an effective identification tool in biomonitoring and conservation programmes. The
20 accuracy of species-level assignment, and consequent taxonomic coverage, relies on
21 comprehensive DNA barcode reference libraries. The role of these libraries is to support
22 species identification, but accidental errors in the generation of the barcodes may compromise
23 their accuracy. Here we present an R-based application, BAGS (Barcode, Audit & Grade
24 System; <https://github.com/tadeu95/BAGS>), that performs automated auditing and annotation
25 of cytochrome c oxidase subunit I (COI) sequences libraries, for a given taxonomic group of
26 animals, available in the Barcode of Life Data System (BOLD). This is followed by
27 implementing a qualitative ranking system that assigns one of five grades (A to E) to each
28 species in the reference library, according to the attributes of the data and congruency of
29 species names with sequences clustered in Barcode Index Numbers (BINs). Our goal is to
30 allow researchers to obtain the most useful and reliable data, highlighting and segregating
31 records according to their congruency. Different tests were performed to perceive its
32 usefulness and limitations. BAGS fulfils a significant gap in the current landscape of DNA
33 barcoding research tools by quickly screening reference libraries to gauge the congruence
34 status of data and facilitate the triage of ambiguous data for posterior review. Thereby, BAGS
35 has the potential to become a valuable addition in forthcoming DNA metabarcoding studies, in
36 the long term contributing to globally improve the quality and reliability of the public reference
37 libraries.

38

39 **Keywords**

40 reference libraries, quality control, annotation, BOLD systems, DNA metabarcoding, R.

41

42 1. INTRODUCTION

43 The availability of well-curated comprehensive reference libraries is fundamental for accurate
44 DNA barcode-based species identification (Cariani et al., 2017; Ekrem, Willassen, & Stur,
45 2007; Leese et al., 2016; Oliveira et al., 2016). The demand for high quality reference libraries
46 has increased considerably since the introduction and extended use of DNA metabarcoding
47 for biodiversity assessments and biomonitoring (Leese et al., 2018; Weigand et al., 2019). Due
48 to the large number of reads from high-throughput sequencing (HTS) instruments, the required
49 bioinformatics often include automated systems to match query sequences to reference
50 sequences in DNA sequence repositories (e.g. Bengtsson-Palme et al., 2018), such as the
51 Barcode of Life Data Systems (BOLD; Ratnasingham & Hebert, 2007) or NCBI GenBank
52 (Sayers et al., 2019). With a few exceptions, such as R-Syst::diatom (Rimet et al., 2016), the
53 UNITE database (Nilsson et al., 2018) or MIDORI (Machida, Leray, Ho, & Knowlton, 2017),
54 which are reference libraries compiled and curated for specific taxa, typically, there is no
55 supervision or quality control of the reference dataset. Therefore, inaccurate records in
56 reference libraries may result in recurrent identification errors which can be perpetuated over
57 time and across studies without being detected (Keller et al., 2020; Leese et al., 2016; Weigand
58 et al., 2019).

59 Errors or discordances can have operational or biological explanations. Operational errors
60 include morphology-based misidentifications, cross-contamination of samples, mislabelling,
61 accidental mistakes when recording data, among others (Packer, Gibbs, Sheffield, & Hanner,
62 2009; Pentinsaari, Ratnasingham, Miller, & Hebert, 2019; Rulik et al., 2017). Possible
63 biological reasons for discordances include recently diverged species and incomplete lineage
64 sorting, introgression, insufficient discrimination capacity of the barcode marker, phenotypic
65 plasticity, among others (Costa & Antunes, 2012; Lin, Stur, & Ekrem, 2018; Weber, Stöhr, &
66 Chenuil, 2019; Weigand, Jochum, Pfenninger, Steinke, & Klussmann-Kolb, 2011). Although
67 some data quality assurance and quality control (QA/QC) criteria have been implemented

68 upstream and along the DNA barcode production workflow (e.g. Hanner, 2005), no
69 comprehensive tool for downstream quality control of the taxonomic accuracy in DNA barcode
70 reference libraries is available to check QA/QC in a standardized way. Some QA/QC measures
71 are implemented in BOLD: Labelling of barcode compliant records, flagging of sequences that
72 are likely contaminations or based on misidentified specimens, flagging of sequences with stop
73 codons (Ratnasingham & Hebert, 2007), and the possibility to run BIN-discordance reports
74 (Ratnasingham & Hebert, 2013). However, there are several sources of potential discordance
75 or errors that remain unscreened or unexplored through existing systems (Meiklejohn,
76 Damaso, & Robertson, 2019; Mioduchowska, Czyz, Gołdyn, Kur, & Sell, 2018; Siddall,
77 Fontanella, Watson, Kvist, & Erséus, 2009; Weigand et al., 2019).

78 The origin of discordances and inaccuracies in DNA barcode data and DNA databases in
79 general are well known (Harris et al., 2003; Meiklejohn et al., 2019; Mioduchowska et al., 2018;
80 Pentinsaari et al., 2019; Siddall et al., 2009; Vilgalys, 2003), however, relatively few studies
81 have addressed the problem of compilation, and quality control of reference libraries,
82 particularly concerning taxonomic reliability (Leese et al., 2018; Weigand et al., 2019). For
83 instance, CO-ARBitrator (Heller, Casaletto, Ruiz, & Geller, 2018) detects sequences
84 mislabelled as cytochrome c oxidase subunit I (COI), but which are originating from non-
85 homologous loci. The “coil” R package (Nugent, Elliot, Ratnasingham, & Adamowicz, 2020) is
86 also useful in detecting errors in animal barcoding and metabarcoding data by placing
87 sequences in a reading frame and translating them to amino acids. While both packages
88 successfully detect cases of non-homologous barcode sequences, they do not address the
89 issue of taxonomic congruency.

90 Recently, Rulik et al. (2017) proposed a pre-processing system for large datasets aiming to
91 generate high quality DNA barcodes by verifying taxonomic consistency. However, this system
92 requires a phylogenetic backbone for implementation, and it is meant to be used before

93 uploading data to reference libraries. It therefore does not consider global congruence with
94 other data already available in either BOLD or GenBank.

95 A large number of COI sequences are currently available in GenBank (Porter & Hajibabaei,
96 2018) and although a fair portion of the records may not abide to the formal barcode data
97 standards (Ratnasingham & Hebert, 2013), they still constitute a useful resource that should
98 not be overlooked. In fact, many metabarcoding-based studies report taxonomic assignments
99 based on all available COI data, thereby including non-barcode compliant records. This
100 reinforces the need for a barcode compilation, auditing and annotation system that provides
101 an indication of the taxonomic reliability of the records for end-users of reference libraries.

102 Costa et al. (2012) proposed a ranking system to be implemented at the post-barcoding end
103 of the barcode production pipeline, which considered all available sequence data for a given
104 species (thus both barcode compliant and non-compliant). The ranking system attributes five
105 different grades to species records (A to E), depending essentially on the level of congruency
106 between morphospecies and the respective COI barcode clusters. Later, the system was
107 updated to use Barcode Index Numbers (BINs; Ratnasingham & Hebert, 2013) as the
108 reference DNA barcode clustering method (e.g. Knebelsberger et al., 2014; Oliveira et al.,
109 2016). In global terms, the goal was to provide end-users of reference libraries with a system
110 to sort out and annotate species that can be confidently identified with current data, from
111 ambiguous or inaccurate records that need revision, or to flag cases of suspected hidden
112 diversity. The implementation of this ranking system to a compilation of COI barcodes from
113 European fish revealed that the majority of species could be confidently identified with DNA
114 barcodes (Oliveira et al., 2016), and a number of ambiguous records could be clarified upon
115 careful revision. However, in these implementations of the ranking system, the attribution of
116 the grades was dependent on individual analyses of each species' data, a strategy which would
117 be impractical for the large DNA metabarcoding reference libraries involving hundreds or
118 thousands of species.

119 To address this problem, we here introduce BAGS, an R-based application for automated
120 auditing and annotation of DNA barcode reference libraries. We adapt the proposed Oliveira
121 et al. (2016) ranking system, essentially based on match/mismatch between BINs and
122 morphospecies identifications. BAGS can be applied to user-provided species lists or large
123 taxon-specific datasets composed of all available COI barcode sequences in BOLD, including
124 those mined from GenBank. BAGS also aims to facilitate revision and curation of barcode
125 reference libraries, thereby contributing to improve their quality.

126

127 **2. METHODS**

128 **2.1. Overview of BAGS**

129 BAGS features automated compilation of quality-filtered COI sequence datasets from BOLD,
130 allowing for selection or exclusion of marine taxa through matching with the World Register of
131 Marine Species (WoRMS) checklists (WoRMS Editorial Board, 2020). It delivers taxon-specific
132 libraries annotated with qualitative grades based on BIN/morphospecies congruence and on
133 the amount of available data for each species (A to E, see below for details), which can be
134 downloaded whole or sorted by grade. A user-friendly interface allows for minimal operation
135 for users non-familiar with R (R Development Core Team, 2019), while providing a grasp of
136 the overall quality of the reference library through a graphical output of the proportion of records
137 and species assigned to each of the five grades. However, since BAGS can also be run locally,
138 the more experienced R users have the option to make adjustments to the code. The users
139 may then (frequently if necessary) use the annotated datasets to compile their own
140 personalized and reviewed libraries (e.g. BOLD datasets) and use them for taxonomic
141 assignment of HTS metabarcoding-generated reads.

142 BAGS is composed of four main features which are implemented in sequence (Figure 1): a)
143 data mining and library compilation, b) marine taxa filter (optional), c) library auditing and
144 annotation and d) auditing output and annotation-based library sorting.

145 **2.2. BAGS pipeline**

146 **2.2.1. Data mining and library compilation**

147 BAGS offers the option for library compilation based on a choice of taxa or through a user-
148 provided species list. Records matching the selected taxa or species list will be retrieved and
149 then filtered. All the data is retrieved from BOLD (www.boldsystems.org), using the “bold” R
150 package (Chamberlain, 2019). Therefore, the taxa introduced by the user must be present in
151 BOLD at the time of use. Any taxonomic rank from species to phylum belonging to the kingdom
152 Animalia can be submitted, but it should be noted that some ranks, particularly intermediate
153 ranks, are not implemented in BOLD or may not be available for some species.

154 The mining of the target taxa can be achieved through three options: download all the records
155 available (all taxa), download only records of species occurring in marine habitats (which may
156 include any taxa present in brackish waters) or download the non-marine species' records (i.e.
157 not present in neither marine or brackish water habitats). This marine species selection or
158 exclusion filter is accomplished resorting to the “worms” R package (Holstein, 2018), which
159 checks the habitat type(s) assigned in WoRMS to each species in a query dataset, among the
160 four available (marine, brackish, freshwater or terrestrial).

161 Records are removed if at least one of the following criteria is verified: a) records with
162 sequences shorter than the minimum size chosen by the user (between 300 and 650 bp), or
163 with sequences that have more than 1% ambiguous base calls (Ns); b) records without species
164 name (this includes records identified only by genus or any higher taxonomic rank), or without
165 BIN; c) by default, records without information of the sampling location (either latitude or
166 country of origin), although users can choose to include those records. Records with

167 ambiguous expressions present in the species name (e.g. *sp.*, *complex.*, *etc*; see Appendix 1:
168 <https://doi.org/10.5061/dryad.2rbnzs7kx>) or in the COI sequence (i.e. not IUPAC nucleotide
169 code; see Appendix 1: <https://doi.org/10.5061/dryad.2rbnzs7kx>) are not removed, however,
170 the ambiguous expression is removed.

171 At the end of this procedure, a filtered reference library is downloaded and available for the
172 subsequent auditing and annotation step.

173 **2.2.2. Auditing and annotation**

174 Following the initial quality-filtering steps, the BAGs pipeline subsequently proceeds to the
175 implementation of the auditing and annotation system adapted with modifications from Oliveira
176 et al. (2016). The five annotation grades attributed to each species in a compiled library are
177 defined as follows (Figure 2):

178 Grade A - Consolidated concordance: the morphospecies is assigned to a single BIN, which
179 integrates only members of that species. Additionally, the species is represented by more than
180 10 specimens in the library.

181 Grade B - Basal concordance: the morphospecies is assigned to a single BIN, which integrates
182 only members of that species, but there are 10 or less specimens in the library.

183 Grade C - Multiple BINs: the morphospecies is assigned to more than one BIN, and all of those
184 BINs integrate only members of that species.

185 Grade D - Insufficient data: the species has less than three specimens available in the library
186 and none of the BINs assigned to the species integrates specimens from another species.

187 Grade E - Discordant species assignment: more than one species is assigned to a single BIN.
188 All the records of that species will be assigned to grade E.

189 The BAGs auditing pipeline consists of a series of annotation steps, each comprising data
190 checks with two possible outcomes (Figure 2). Every set of sequences for a given species

191 entering the pipeline will be annotated with a single grade (A to E). Discordant species
192 assignments (grade E) are immediately screened at the front end of the pipeline, followed by
193 records with insufficient data (grade D), then grade C. Grades A or B are attributed last, if the
194 records were not retained in the previous screens. The screening steps involve checking
195 against the full BOLD database, thus not exclusively considering the reference library being
196 downloaded at the time of the annotation, that would limit concordance-checking to the
197 downloaded species' data only.

198 BOLD (like GenBank) limits the number of searches or queries per IP/user to avoid the
199 overload of their webservice. Therefore, to avoid blocking the access to BOLD, we periodically
200 (approximately every two months) download the entire BOLD dataset for animals and protists
201 in order to calculate the number of BINs for each species, as well as the number of species for
202 each BIN. With this solution, BAGS can work faster and without the computational limitations
203 of real-time query searches on BOLD.

204 **2. 2. 3. Output and annotation-based file sorting**

205 The auditing system proceeds then to the annotation of the records with the pre-defined grades
206 to each species in the reference library, following the pipeline described before. In due course
207 the reference library will be created and downloaded in the form of a tabular file containing the
208 following: species name, BIN, COI-5P sequence, country or region of origin, the grade that
209 was attributed to the species, number of base pairs in the sequence, family, order, class,
210 sample ID, process ID, latitude, longitude and in the case of marine taxa libraries, an additional
211 column with the valid species name according to WoRMS. The user has also the option to
212 download the reference library in *fasta* format, giving the choice of which grades to include.
213 The *fasta* files can be download with all grades, combinations of different grades or separately
214 for each grade.

215 Lastly, BAGS summarizes the data regarding the reference library that was created, in the form
216 of a text report plus two bar plots: one displaying the number of specimens for each attributed
217 grade and another displaying the number of species for each attributed grade. In order to
218 repeat the process for additional target libraries, the user must refresh the page and start over
219 again.

220 **2. 3. Informatic implementation**

221 BAGS is an application written entirely in the open-source programming language R, designed
222 using the “shiny” R package framework (Chang, Cheng, Allaire, & Xie, 2019), having therefore
223 an underlying customization with HTML and CSS. It is possible to launch BAGS locally on any
224 environment that has R installed, as well as through any R IDE such as RStudio (RStudio
225 Team, 2016), where it can fully operate as long as there is a stable internet connection and
226 the databases BOLD and WoRMS are functional. The application can be used without any
227 prior knowledge of the R programming language, and the instructions for launching it can be
228 consulted in the “README” file. BAGS is stored at web servers and can also be used remotely,
229 which allows its launching from any web browser (<https://bags.vm.ntnu.no> or [https://tadeu-
230 apps.shinyapps.io/bags](https://tadeu-
230 apps.shinyapps.io/bags); additional links are provided at <https://github.com/tadeu95/BAGS>).
231 The script that allows the application to be run locally without constraints in R, as well as a
232 “README” file, are currently stored at GitHub: <https://github.com/tadeu95/BAGS>.

233 **2.4. Performance assessment**

234 In order to test BAGS performance, two independent tests were performed. First, to
235 understand if the marine and non-marine taxa selection filters were functional and reliable, we
236 downloaded three files, using the “all taxa”, the “marine taxa” and the “non-marine” taxa options
237 for a family of shrimps, Palaemonidae, which comprises species from various aquatic habitats.
238 This was followed by checking the report generated by BAGS and manually checking 30
239 random species from each of the three libraries previously generated.

240 Second, to assess the accuracy of BAGS' auditing and grade assignment, we selected three
241 trial reference libraries likely to display distinctive features and quality issues: marine
242 Amphipoda (Malacostraca: Crustacea), Chironomidae (Diptera: Insecta), and marine fish
243 (Actinopterygii, Elasmobranchii and Holocephali). These trial libraries include two key
244 invertebrate groups in aquatic monitoring, which are likely to be relevant in metabarcoding
245 applications, and one of the most well-represented groups of vertebrates in BOLD, thereby
246 enabling the screening of a large and diverse number of records and species.. Three reference
247 libraries were downloaded using as input "Amphipoda" (within the marine taxa filter option),
248 "Chironomidae" (all taxa option), and "Actinopterygii,Elasmobranchii,Holocephali" also within
249 the marine taxa filter option. Then, the grade assignment was checked by randomly sampling
250 30 species from each assigned grade, from each compiled library and checking the data
251 manually to assess if the grades were correctly assigned to their specimens. Due to the
252 massive amount of data available for Chironomidae (more than 400,000 sequences accessible
253 on BOLD), the species in the compiled library were matched against a list of European species
254 used for freshwater biomonitoring under the EU Water Framework Directive (BOLD checklist
255 DNAqua-NET: Diptera, code CL-DNADI, 584 spp. Chironomidae) in order to simplify the
256 performance assessment. Neighbour-Joining trees (Saitou & Nei, 1987) of the species
257 assigned to grade C were created on the BOLD workbench, to evaluate the monophyly/non-
258 monophyly of each species. Within grade E, different plausible origins for the discordance were
259 scored for the following categories: synonym; faulty or ambiguous species names;
260 consolidated morphospecies grouped in one BIN; probable misidentification and inconclusive
261 origin.

262

263 **3. RESULTS**

264 **3.1. Marine taxa selection filter**

265 Using the input “Palaemonidae” within the marine filter, the marine taxa library comprised 60
266 species assigned to 73 BINs, and a total of 577 specimens, while the non-marine taxa library
267 comprised 51 species, 67 BINs and a total of 318 specimens. Comparatively, the “all taxa”
268 option library had 123 species, 148 BINs and 1,022 specimens. The 30 species randomly
269 sampled of the marine-filtered library were correctly assigned (i.e. all the 30 species were
270 registered as being from marine or brackish environments when checked manually upon on
271 WoRMS; Appendix 2: <https://doi.org/10.5061/dryad.2rbnzs7kx>). Nonetheless, this included
272 species which were registered simultaneously as occurring in both marine and freshwater
273 habitats. On the other hand, the 30 species manually checked from the non-marine taxa library
274 revealed to be all exclusive from freshwater environments (i.e. not present neither in marine or
275 brackish waters, and therefore not present in the marine library).

276 **3.2. Trial datasets**

277 The marine Amphipoda dataset had a total of 6,385 specimens in the compiled library, 486
278 species and 736 BINs; the Chironomidae dataset consisted of a total of 90,214 specimens,
279 1,113 species and 1,883 BINs; and the marine fishes dataset comprised 107,434 specimens,
280 8,381 species and 9,779 BINs (Appendix 3: <https://doi.org/10.5061/dryad.2rbnzs7kx>). The
281 distributions of the number of species per grade in each of the compiled reference libraries
282 (Figure 3) show that the proportion of possible cases of hidden diversity (grade C) is higher in
283 the two invertebrate libraries (Amphipoda and Chironomidae; around 20%) compared with the
284 marine fish library (less than 10%). Cases of insufficient records, which consist of species with
285 less than three specimens in the BAGS-compiled library (Grade D), are also less prevalent in
286 the marine fishes (~18%) when compared to both invertebrate libraries (40% and 26% for
287 Amphipoda and Chironomidae respectively). On the other hand, cases of apparent
288 discordance (Grade E) are considerably less prevalent in the Amphipoda library (only 12% of
289 the cases) and much more frequent in the marine fish library (44%). The number of species

290 per BIN (grade E) varied between 1 and 49 for Amphipoda, 1 and 12 for Chironomidae and 1
291 and 88 for fish (Figure 4).

292 For the three groups, the 30 randomly sampled species were correctly assigned to the
293 qualitative grades (Appendix 4: <https://doi.org/10.5061/dryad.2rbnzs7kx>). Grade C species
294 (Table 1, Appendix 5: <https://doi.org/10.5061/dryad.2rbnzs7kx>) were mostly monophyletic:
295 between 66% (Chironomidae) and 80% (fish). Discordances or potential errors in grade E
296 annotations had different possible sources (Table 2). Misidentifications (between 37% and
297 67%) and ambiguous species names (between 10% and 33%) contributed the most to the
298 grade E cases, while synonyms the least (overall 3.4%).

299

300 **4. DISCUSSION**

301 While molecular and computational tools have been increasingly providing taxonomists with
302 large volumes of data to analyse, the need for systems which classify and audit that data is
303 now more relevant than ever. This is especially the case when dealing with publicly available
304 DNA barcodes, which can be freely submitted to biological databases and subsequently used
305 by researchers anywhere, at any time (Curry, Gibson, Shokralla, Hajibabaei, & Baird, 2018;
306 Meiklejohn et al., 2019). Moreover, given the establishment of DNA barcoding as one of the
307 primary drivers behind the recent scientific efforts in uncovering and explaining biodiversity
308 (DeSalle & Goldstein, 2019; Pennisi, 2019), our primary goal with BAGS is to facilitate the
309 implementation of curation and quality control measures among taxonomists and biodiversity
310 scientists. Additionally, we seek to do this through a user-friendly and automated platform,
311 removing any need for programming skills in order to audit and annotate a reference library.

312 BAGS differs from the "BIN discordance report" available at BOLD, within the sequence
313 analysis tools. First of all, whereas the BOLD tool is BIN-centred, our approach is
314 morphospecies-centred. This fundamental difference has a number of consequences. While

315 BOLD reports discordant BINs, BAGS reports on discordant morphospecies, meaning that a
316 morphospecies displaying even a single record in a discordant BIN is classified as grade E.
317 The morphospecies-centred approach also enables BAGS to report on species occurring in
318 multiple - but non-discordant - BINS (grade C), therefore serving as a barometer of suspected
319 hidden diversity in reference libraries. Finally, BAGS also takes into consideration the amount
320 of sequences available in the database, providing a grasp of gaps in comprehensiveness of
321 coverage for morphospecies in the reference libraries (grades A, B, and D). From an auditing
322 and taxonomic curation point of view, the morphospecies-centred approach is also more
323 advantageous.

324 Ultimately, we present this application as a way to sort out taxonomic incongruencies and point
325 out possible cases of human error during the generation of the barcodes, as well as uncovering
326 potential cases of hidden diversity among species. Overall, our comprehensive testing of
327 BAGS indicates that it enables researchers with a simple tool for fast screening of the quality
328 status of massive reference libraries, thereby allowing them to sort highly robust records and
329 pinpoint those in need of curation and revision, while unravelling the main issues that may
330 arise during the generation of DNA barcodes for a particular group of organisms.

331 **4.1. BAGS performance assessment tests and some considerations**

332 The different efficiency tests performed with BAGS (either marine/non-marine and grade
333 annotation) allowed to verify the correct performance of this application. The different manual
334 tests (i.e. non-automated; Appendixes 2 and 4: <https://doi.org/10.5061/dryad.2rbnzs7kx>) and
335 the ongoing tests performed by us and colleagues during beta tests, did not bring to light any
336 errors of the application in the filtering or the auditing and annotation steps.

337 It is important to point out that some transitional marine species (i.e. present in estuaries) are
338 registered in WoRMS as being from brackish habitats, which can include both typical marine
339 or freshwater species (e.g. *Phoxinus* Rafinesque, 1820). These species should not be

340 excluded from marine reference libraries as they may be also detected in metabarcoding
341 studies in fully marine environments. If the goal is to gather as much barcode compliant records
342 as possible in the final dataset, regardless of the habitat, it is advisable to use the “all taxa”
343 option. However, we consider the “marine” and “non-marine” options of BAGS as a useful
344 resource if the user wishes to use a customized and size-amenable reference library targeting
345 preferentially only marine or non-marine organisms.

346 By using three distinct taxa important in biomonitoring studies, BAGS allowed us to promptly
347 understand the differences in the level of congruency of their available DNA barcodes and in
348 the quality of their respective reference libraries. Recent initiatives (e.g. deWaard et al., 2019;
349 Hobern & Hebert, 2019; Leese et al., 2016) have been striving to increase the taxonomic
350 coverage of universal databases, however, DNA barcodes are still missing for many species
351 (e.g. Weigand et al., 2019) or are poorly represented (high prevalence of grade D species here
352 observed; Figure 3), reinforcing the continuous need for the completion of reference libraries.

353 BAGS performance tests allowed to spot a high proportion of grade C species (multiple BINs),
354 reaching around 20% in Chironomidae and Amphipoda, but less prevalent in marine fish
355 (Figure 3). Species with multiple BINs may occur for a number of reasons, starting with the
356 non-optimal BIN splitting (Ratnasingham & Hebert, 2013), *Wolbachia*-related artefacts
357 (especially terrestrial arthropods; Smith et al., 2012), or may simply reflect phylogeographic
358 differentiation within the same species. However, they also often suggest undescribed or
359 cryptic diversity. Indeed, a fair amount of cases of cryptic diversity have been reported in the
360 literature for marine amphipods (e.g. Hyalidae, Desiderato et al., 2019; Gammaridae, Hupalo
361 et al., 2019), while the family Chironomidae belongs to an order (Diptera) notorious for
362 incorporating large numbers of hidden species (Ekrem, Stur, & Hebert, 2010; Lin, Stur, &
363 Ekrem, 2015). In marine fish on the other hand, detection of cryptic species has been less
364 reported (Knebelsberger et al., 2014; Oliveira et al., 2016), maybe due to the fact that their
365 taxonomy is possibly more updated, morphological differentiation is more rigorously

366 established for most species, or the fact that their high mobility may reduce the likelihood of
367 genetic divergence between populations over larger distances. A number of studies have been
368 addressing the curation of marine invertebrate's DNA barcodes, including Amphipoda (e.g.
369 Lobo et al. 2016; Radulovici et al. 2019; Raupach et al. 2015), which may explain the lowest
370 proportion of possible discordances (Grade E) out of the three groups analysed (Figure 3).
371 Contrarily, the marine fishes' reference library showed a prominently high proportion (~44%)
372 of grade E species (Figure 3), mainly due to misidentifications, consolidated morphospecies
373 aggregated in one BIN or faulty species names lexicon (Table 2). There are some extreme
374 cases which greatly contribute to this scenario, as for instance, BINs BOLD:AAC8034 and
375 BOLD:AAB3926, consisting of 40 and 88 species respectively (Figure 4). In the latter case, out
376 of 88 species, only one is spelled correctly ("*Pseudanthias squamipinnis*"), while the remaining
377 were named "Unknown" or "*Pseudanthias* sp." followed by different alphanumeric
378 designations. Since these ambiguous species names, possibly interim names, are not properly
379 standardized, BAGS considers them different species for the purpose of comparison against
380 BOLD database and grade assignment, even though it does remove the ambiguous
381 expressions and specimens assigned only to genus, in the compiled libraries. Considering this
382 and other possible grade E scenarios, we hold the view that this grade should serve as an
383 incentive for a close examination of that particular species' records, and not as a definitive
384 signalling of unreliability. Indeed, the detailed inspection of grade E cases after BAGS
385 annotation revealed that most of them are likely pseudo-discordances and, if eventually
386 clarified, could lead to an estimated overall reduction of 80% in grade E species.

387 **4.3. BAGS limitations**

388 Although this current version of BAGS has its own merits and stands on its own as a complete
389 tool, filling a gap in the current DNA barcoding research landscape that we identified, there are
390 still limitations that we would like to address in future versions. Currently, BAGS does not have
391 the ability to flag gross sequence mismatches, such as bacterial sequences mistakenly

392 assigned to animals, as it has been previously reported (Siddall et al., 2009). Although these
393 might be rare events, it would be useful to fully discriminate these cases so that the congruency
394 of the reference library is increased, and more errors are subsequently flagged. Additionally,
395 in its current version, BAGS cannot distinguish grade C's monophyletic from non-monophyletic
396 species, nor can it recognize synonyms and other apparent discordances, such as faulty or
397 interim species names, in species graded E. Moreover, since BAGS implements grades which
398 are defined based on the BIN/morphospecies matches, the limitations associated with the
399 accuracy of the BIN clustering algorithm may emerge in some results or particular groups of
400 organisms. This could be possibly improved in future versions with the introduction of
401 customized OTU clustering algorithms that may be useful to complement the BIN-based
402 auditing, opening possibilities for its application beyond COI sequences and the BOLD
403 database.

404 Many databases (e.g. BOLD, GenBank, WoRMS) have systems that detect excessive calls by
405 the same user (i.e., too many searches or queries) that might overload their webservice, and
406 therefore, they either limit the number of calls or block the user's IP address for a period of
407 time. Since BAGS relies on multiple searches on BOLD, this restriction would limit its efficiency.
408 To overcome this constraint, part of the data necessary to implement the grade annotation
409 system is regularly downloaded by us from BOLD, and used for comparison. However, since
410 the full species name and BIN dataset is locally stored for this purpose, the grade attribution
411 can potentially change every time new barcode records and BINs are added to BOLD.

412

413 **5. FINAL REMARKS**

414 We can envision several prospective improvements that may be considered in future versions
415 of BAGS. One such key improvement would be to introduce the capability to detect cases of
416 deep discordance which may in fact appear concordant (hence pseudo-concordances), such

417 as the cases of bacterial DNA inadvertently amplified from metazoan DNA during PCR, further
418 included in public genetic repositories assigned to metazoan species (Siddall et al., 2009).
419 Introduction of a phylogenetic placement auditing tool would constitute a possible solution to
420 detect such events, and it would also be essential to discriminate cases of monophyly and
421 non-monophyly in grade C-assigned species. Additional improvements to BAGS may include
422 implementation of alternative clustering algorithms and customized filtering thresholds, making
423 it prone for future implementations using other DNA-barcode sequence systems and
424 databases. Finally, the inclusion of a subsidiary tool to perform a detailed revision of grade E
425 records, in order to signal, for example, pseudo-discordances generated by synonyms or
426 ambiguous species designations, possibly using machine learning and artificial intelligence
427 systems. Eventually, some discordances may require individual professional judgement that
428 cannot be accomplished with automated procedures.

429 It is our goal that BAGs can facilitate and stimulate the much-needed revision and curation of
430 reference libraries. We urge all users to contribute to this critical task for the sake of the quality
431 of the libraries and ultimately the soundness of the research that depends on it.

432

433 **Acknowledgments**

434 This study was supported by the project NextSea [NORTE-01-0145-FEDER-000032], under
435 the PORTUGAL 2020 Partnership Agreement, through the European Regional Development
436 Fund (ERDF). PV was supported through the Portuguese Foundation for Science and
437 Technology (FCT, I.P.) in the scope of the project NIS-DNA [PTDC/BIA-BMA/29754/2017].

438 This manuscript contributes to the COST Action DNAqua-Net CA15219 goals, in particular
439 Work Group 1 (WG1), and benefited from comments and suggestions over an incipient version
440 of BAGS from participants in DNAqua-Net WG1 workshop in Limassol, Cyprus.

441 Finally, we would like to thank Sofia Duarte for testing BAGS and to the anonymous reviewers
442 for their suggestions to improve the manuscript.

443

444 **References**

- 445 Bengtsson-Palme, J., Richardson, R. T., Meola, M., Wurzbacher, C., Tremblay, E. D., Thorell, K., ...
446 Nilsson, R. H. (2018). Metaxa2 Database Builder: enabling taxonomic identification from
447 metagenomic or metabarcoding data using any genetic marker. *Bioinformatics*, **34**(23), 4027–
448 4033. <https://doi.org/10.1093/bioinformatics/bty482>
- 449 Cariani, A., Messinetti, S., Ferrari, A., Arculeo, M., Bonello, J. J., Bonnici, L., ... Tinti, F. (2017).
450 Improving the conservation of mediterranean chondrichthyans: The ELASMOMED DNA barcode
451 reference library. *PLoS ONE*, **12**(1), e0170244. <https://doi.org/10.1371/journal.pone.0170244>
- 452 Chamberlain, S. (2019). bold: Interface to Bold Systems API. R package version 0.9.0. [https://cran.r-](https://cran.r-project.org/package=bold)
453 [project.org/package=bold](https://cran.r-project.org/package=bold)
- 454 Costa, F. O., Landi, M., Martins, R., Costa, M. H., Costa, M. E., Carneiro, M., ... Carvalho, G. R. (2012).
455 A ranking system for reference libraries of DNA barcodes: application to marine fish species from
456 Portugal. *PloS One*, **7**(4), 1–9. <https://doi.org/10.1371/journal.pone.0035858>
- 457 Costa, F. O., & Antunes, P. M. (2012). The contribution of the Barcode of Life initiative to the discovery
458 and monitoring of biodiversity. In: Mendonca, A., Cunha, A., & Chakrabarti, R. (eds.), *Natural*
459 *Resources, Sustainability and Humanity: A Comprehensive View*. Springer, Dordrecht, 37–68.
460 https://doi.org/10.1007/978-94-007-1321-5_4
- 461 Chang, W., Cheng, J., Allaire, J. J., Xie, Y., McPherson, J. (2019). shiny: Web Application Framework
462 for R. R package version 1.4.0. <https://cran.r-project.org/package=shiny>
- 463 Curry, C. J., Gibson, J. F., Shokralla, S., Hajibabaei, M., & Baird, D. J. (2018). Identifying north American
464 freshwater invertebrates using DNA barcodes: Are existing COI sequence libraries fit for purpose?
465 *Freshwater Science*, **37**(1), 178–189. <https://doi.org/10.1086/696613>
- 466 DeSalle, R., & Goldstein, P. (2019). Review and interpretation of trends in DNA Barcoding. *Frontiers in*
467 *Ecology and Evolution*, **7**, 302. <https://doi.org/10.3389/fevo.2019.00302>
- 468 Desiderato, D., Costa, F. O., Serejo, C., Abiatti, M., Queiroga, H., Vieira, P. E. (2019). Macaronesian
469 islands as promoters of diversification in amphipods: the remarkable case of the family Hyalidae
470 (Crustacea, Amphipoda). *Zoologica Scripta*, **48**(3), 359-375. <https://doi.org/10.1111/zsc.12339>.
- 471 deWaard, J. R., Ratnasingham, S., Zakharov, E. V. Borisenko, A. V., Steinke, D., Telfer, A. C., ... Hebert,
472 P. D. N. (2019). A reference library for Canadian invertebrates with 1.5 million barcodes, voucher
473 specimens, and DNA samples. *Scientific Data*, **6**, 308. [https://doi.org/10.1038/s41597-019-0320-](https://doi.org/10.1038/s41597-019-0320-2)
474 [2](https://doi.org/10.1038/s41597-019-0320-2)
- 475 Ekrem, T., Willassen, E., & Stur, E. (2007). A comprehensive DNA sequence library is essential for
476 identification with DNA barcodes. *Molecular Phylogenetics and Evolution*, **43**(2), 530-542.
477 <https://doi.org/10.1016/j.ympev.2006.11.021>.
- 478 Ekrem, T., Stur, E., & Hebert, P. D. N. (2010). Females do count: Documenting chironomidae (Diptera)
479 species diversity using DNA barcoding. *Organisms Diversity and Evolution*, **10**(5), 397–408.
480 <https://doi.org/10.1007/s13127-010-0034-y>
- 481 Hanner, B. R. (2005). Proposed Standards for BARCODE Records in INSDC (BRIs). Technical report,
482 Database Working Groups, Consortium for the Barcode of Life, 2009.

- 483 Harris, T. W., Lee, R., Schwarz, E., Bradnam, K., Lawson, D., Chen, W., ... Stein, L. D. (2003).
484 WormBase: A cross-species database for comparative genomics. *Nucleic Acids Research*, **31**(1),
485 133–137. <https://doi.org/10.1093/nar/gkg053>
- 486 Heller, P., Casaletto, J., Ruiz, G., & Geller, J. (2018). A database of metazoan cytochrome c oxidase
487 subunit I gene sequences derived from GenBank with CO-ARBitrator. *Scientific Data*, **5**, 180156.
488 <https://doi.org/10.1038/sdata.2018.156>
- 489 Hobern, D. & Hebert, P. D. N. (2019). BIOSCAN - Revealing Eukaryote Diversity, Dynamics, and
490 Interactions. *Biodiversity Information Science and Standards*, **3**, e37333.
491 <https://doi.org/10.3897/biss.3.37333>.
- 492 Holstein, J. (2018). worms: Retriving Aphia Information from World Register of Marine Species. R
493 package version 0.2.2. <https://cran.r-project.org/package=worms>
- 494 Hupało, K., Teixeira, M. A. L., Rewicz, T., Sezgin, M., Iannilli, V., Karaman, G. S., ...Costa, F. O. (2019).
495 Persistence of phylogeographic footprints helps to understand cryptic diversity detected in two
496 marine amphipods widespread in the Mediterranean basin. *Molecular Phylogenetics and*
497 *Evolution*, **132**, 53-66. <https://doi.org/10.1016/j.ympev.2018.11.013>.
- 498 Keller, A., Hohlfeld, S., Kolter, A., Schultz, J., Gemeinholzer, B., & Ankenbrand, M. J. (2020).
499 BCdatabaser: on-the-fly reference database creation for (meta-)barcoding. *Bioinformatics*, **36**(8),
500 2630-2631. <https://doi.org/10.32942/osf.io/cmfu2>
- 501 Knebelberger, T., Landi, M., Neumann, H., Kloppmann, M., Sell, A. F., Campbell, P. D., ... Costa, F.
502 O. (2014). A reliable DNA barcode reference library for the identification of the North European
503 shelf fish fauna. *Molecular Ecology Resources*, **14**(5), 1060–1071. <https://doi.org/10.1111/1755-0998.12238>
- 505 Leese, F., Altermatt, F., Bouchez, A., Ekrem, T., Hering, D., Meissner, K., ... Zimmermann, J. (2016).
506 DNAqua-Net: Developing new genetic tools for bioassessment and monitoring of aquatic
507 ecosystems in Europe. *Research Ideas and Outcomes*, **2**, e11321.
508 <https://doi.org/10.3897/rio.2.e11321>
- 509 Leese, F., Bouchez, A., Abarenkov, K., Altermatt, F., Borja, Á., Bruce, K., ... Weigand, A. M. (2018).
510 Why We Need Sustainable Networks Bridging Countries, Disciplines, Cultures and Generations
511 for Aquatic Biomonitoring 2.0: A Perspective Derived From the DNAqua-Net COST Action.
512 *Advances in Ecological Research*, **58**, 63–99. <https://doi.org/10.1016/bs.aecr.2018.01.001>
- 513 Lin, X. L., Stur, E., & Ekrem, T. (2015). Exploring Genetic Divergence in a Species-Rich Insect Genus
514 Using 2790 DNA Barcodes. *PLoS ONE*, **10**(9), e0138993.
515 <https://doi.org/10.1371/journal.pone.0138993>
- 516 Lin, X. L., Stur, E., & Ekrem, T. (2018). DNA barcodes and morphology reveal unrecognized species in
517 Chironomidae (Diptera). *Insect Systematics and Evolution*, **49**(4), 329–398.
518 <https://doi.org/10.1163/1876312X-00002172>
- 519 Lobo, J., Ferreira, M. S., Antunes, I. C., Teixeira, M. A., Borges, L. M., Sousa, R., ...Costa, F. O. (2017).
520 Contrasting morphological and DNA barcode-suggested species boundaries among shallow-water
521 amphipod fauna from the southern European Atlantic coast. *Genome*, **60**(2), 147-157.
522 <https://doi.org/10.1139/gen-2016-0009>.
- 523 Machida, R., Leray, M., Ho, S., & Knowlton, N. (2017). Metazoan mitochondrial gene sequence
524 reference datasets for taxonomic assignment of environmental samples. *Scientific Data*, **4**,
525 170027. <https://doi.org/10.1038/sdata.2017.27>
- 526 Meiklejohn, K. A., Damaso, N., & Robertson, J. M. (2019). Assessment of BOLD and GenBank – Their
527 accuracy and reliability for the identification of biological materials. *PLoS ONE*, **14**(6), e0217084.
528 <https://doi.org/10.1371/journal.pone.0217084>
- 529 Mioduchowska, M., Czyz, M. J., Gołdyn, B., Kur, J., & Sell, J. (2018). Instances of erroneous DNA
530 barcoding of metazoan invertebrates: Are universal *cox1* gene primers too “universal”? *PLoS ONE*,

531 **13**(6), e0199609. <https://doi.org/10.1371/journal.pone.0199609>

532 Nilsson, R. H., Larsson, K-H., Taylor, A. F. S., Bengtsson-Palme, J., Jeppesen, T. S., Schigel, D., ...
533 Abarenkov, K. (2018). The UNITE database for molecular identification of fungi: handling dark taxa
534 and parallel taxonomic classifications. *Nucleic Acids Research*, **47**(D1), D259–D264.
535 <https://doi.org/10.1093/nar/gky1022>

536 Nugent, C. M., Elliott, T. A., Ratnasingham, S., & Adamowicz, S. J. (2020). coil: An R package for
537 cytochrome C oxidase I (COI) DNA barcode data cleaning, translation, and error evaluation.
538 [published online ahead of print, 2020 May 14] *Genome*. <https://doi.org/10.1139/gen-2019-0206>

539 Oliveira, L. M., Knebelsberger, T., Landi, M., Soares, P., Raupach, M. J., & Costa, F. O. (2016).
540 Assembling and auditing a comprehensive DNA barcode reference library for European marine
541 fishes. *Journal of Fish Biology*, **89**(6), 2741–2754. <https://doi.org/10.1111/jfb.13169>

542 Packer, L., Gibbs, J., Sheffield, C., & Hanner, R. (2009). DNA barcoding and the mediocrity of
543 morphology. *Molecular Ecology Resources*, **9**(Suppl s1), 42–50. <https://doi.org/10.1111/j.1755-0998.2009.02631.x>

545 Pennisi, E. (2019). DNA barcodes jump-start search for new species. *Science*, **364**(6444), 920–921.
546 <https://doi.org/10.1126/science.364.6444.920>

547 Pentinsaari, M., Ratnasingham, S., Miller, S. E., & Hebert, P. D. N. (2020). BOLD and GenBank revisited
548 – Do identification errors arise in the lab or in the sequence libraries? *PLoS ONE*, **15**(4), e0231814.
549 <https://doi.org/10.1371/journal.pone.0231814>

550 Porter, T. M., & Hajibabaei, M. (2018). Over 2.5 million COI sequences in GenBank and growing. *PLoS*
551 *ONE*, **13**(9), e0200177. <https://doi.org/10.1371/journal.pone.0200177>

552 Radulovici, A. E., Costa, F. O. and the Hackathon Participants (2019). New avenues for data curation:
553 hackathon on marine invertebrates. *Genome*, **62**(6), 422. <https://doi.org/10.1139/gen-2019-0083>.

554 Ratnasingham, S., & Hebert, P. D. N. (2007). BOLD: The Barcode of Life Data System
555 (<http://www.barcodinglife.org>). *Molecular Ecology Notes*, **7**(3), 355–364.
556 <https://doi.org/10.1111/j.1471-8286.2006.01678.x>

557 Ratnasingham, S., & Hebert, P. D. N. (2013). A DNA-Based Registry for All Animal Species: The
558 Barcode Index Number (BIN) System. *PLoS ONE*, **8**(7), e66213.
559 <https://doi.org/10.1371/journal.pone.0066213>

560 R Development Core Team. (2019). R: A language and environment for statistical computing. R
561 Foundation for Statistical Computing. Vienna, Austria. <http://www.R-project.org>.

562 Raupach, M. J., Barco, A., Steinke, D., Beermann, J., Laakmann, S., Mohrbeck, I., ... Knebelsberger, T.
563 (2015). The Application of DNA barcodes for the identification of marine crustaceans from the
564 North Sea and adjacent regions. *PLoS ONE*, **10**(9), e0139421.
565 <https://doi.org/10.1371/journal.pone.0139421>

566 Rimet, F., Chaumeil, P., Keck, P., Kermarrec, L., Vasselon, V., Kahlert, M., ... Bouchez, A. (2016). R-
567 Syst::diatom: an open-access and curated barcode database for diatoms and freshwater
568 monitoring. *Database*, **2016**, 1-21. <https://doi.org/10.1093/database/baw016>

569 RStudio Team (2016). RStudio: Integrated Development for R. RStudio, Inc., Boston, MA.
570 <http://www.rstudio.com>.

571 Rulik, B., Eberle, J., von der Mark, L., Thormann, J., Jung, M., Köhler, F., ... Ahrens, D. (2017). Using
572 taxonomic consistency with semi-automated data pre-processing for high quality DNA barcodes.
573 *Methods in Ecology and Evolution*, **8**(12), 1878–1887. <https://doi.org/10.1111/2041-210X.12824>

574 Saitou, N., & Nei, M. (1987). The neighbor-joining method: a new method for reconstructing phylogenetic
575 trees. *Molecular Biology and Evolution*, **4**(4), 406–425.
576 <https://doi.org/10.1093/oxfordjournals.molbev.a040454>

- 577 Sayers, E. W., Cavanaugh, M., Clark, K., Ostell, J., Pruitt, K. D., & Karsch-Mizrachi, I. (2019). GenBank.
578 *Nucleic Acids Research*, **47**(D1), D94–D99. <https://doi.org/10.1093/nar/gky989>
- 579 Siddall, M. E., Fontanella, F. M., Watson, S. C., Kvist, S., & Erséus, C. (2009). Barcoding bamboozled
580 by bacteria: Convergence to metazoan mitochondrial primer targets by marine microbes.
581 *Systematic Biology*, **58**(4), 445–451. <https://doi.org/10.1093/sysbio/syp033>
- 582 Smith, M. A., Bertrand, C., Crosby K., Eveleigh E. S., Fernandez-Triana, J., Fisher, B. L., ... Zhou, X.
583 (2012). *Wolbachia* and DNA barcoding insects: patterns, potential, and problems. *PLoS ONE*, **7**(5),
584 e36514. <https://doi.org/10.1371/journal.pone.0036514>
- 585 Vilgalys, R. (2003). Taxonomic misidentification in public DNA databases. *New Phytologist*, **160**(1), 4-
586 5. <https://doi.org/10.1046/j.1469-8137.2003.00894.x>
- 587 Weber, A. A. T., Stöhr, S., & Chenuil, A. (2019). Species delimitation in the presence of strong
588 incomplete lineage sorting and hybridization: Lessons from Ophioderma (Ophiuroidea:
589 Echinodermata). *Molecular Phylogenetics and Evolution*, **131**, 138–148.
590 <https://doi.org/10.1016/j.ympev.2018.11.014>
- 591 Weigand, A. M., Jochum, A., Pfenninger, M., Steinke, D., & Klussmann-Kolb, A. (2011). A new approach
592 to an old conundrum-DNA barcoding sheds new light on phenotypic plasticity and morphological
593 stasis in microworms (Gastropoda, Pulmonata, Carychiidae). *Molecular Ecology Resources*, **11**(2),
594 255–265. <https://doi.org/10.1111/j.1755-0998.2010.02937.x>
- 595 Weigand, H., Beermann, A. J., Čiampor, F., Costa, F. O., Csabai, Z., Duarte, S., ... Ekrem, T. (2019).
596 DNA barcode reference libraries for the monitoring of aquatic biota in Europe: Gap-analysis and
597 recommendations for future work. *Science of the Total Environment*, **678**, 499–524.
598 <https://doi.org/10.1016/j.scitotenv.2019.04.247>
- 599 WoRMS Editorial Board (2020). World Register of Marine Species. Available from
600 <http://www.marinespecies.org> at VLIZ. <https://doi.org/10.14284/170>

601

602 **Data accessibility**

603 BAGS script and associated README file can be consulted at:

604 <https://github.com/tadeu95/BAGs>.

605 To run BAGS remotely, please use <https://bags.vm.ntnu.no>, <https://tadeu->

606 [apps.shinyapps.io/bags](https://tadeu-apps.shinyapps.io/bags) or any additional link available at <https://github.com/tadeu95/BAGS>).

607 All the Appendixes can be found at: <https://doi.org/10.5061/dryad.2rbnzs7kx>.

608

609 **Author contributions**

610 JTF, PEV, PS and FOC designed the research plan. JTF and PEV wrote the BAGs script and
611 developed the BAGS application. JTF, PEV and TE performed the different assessment tests.

612 All the authors wrote the manuscript, contributed with suggestions to the manuscript structure
613 and reviewed the manuscript final version.

614 **Tables and figures**

615 Table 1 - Percentage of monophyletic or non-monophyletic species assigned to grade C of each tested taxonomic
 616 group, according to their position in the Neighbour-Joining trees constructed.

	Monophyletic	Non-monophyletic
Marine Amphipoda	76.7%	23.3%
Chironomidae	66.7%	33.3%
Marine fish	80.0%	20.0%
Overall	74.4%	25.6%

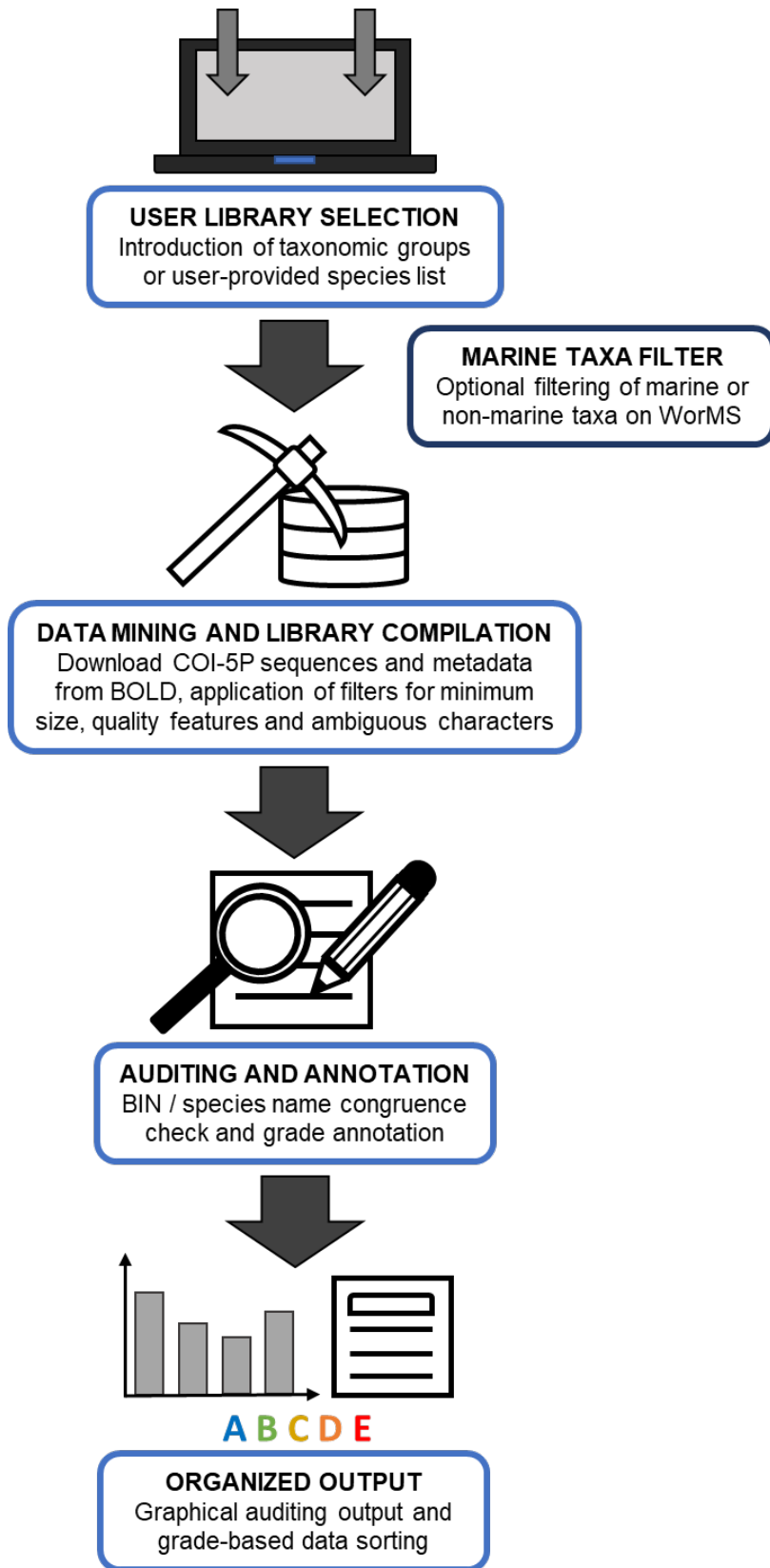
617

618

619 Table 2 - Percentage of the different plausible origins for the assignment of grade E to species in for each tested
 620 taxonomic group.

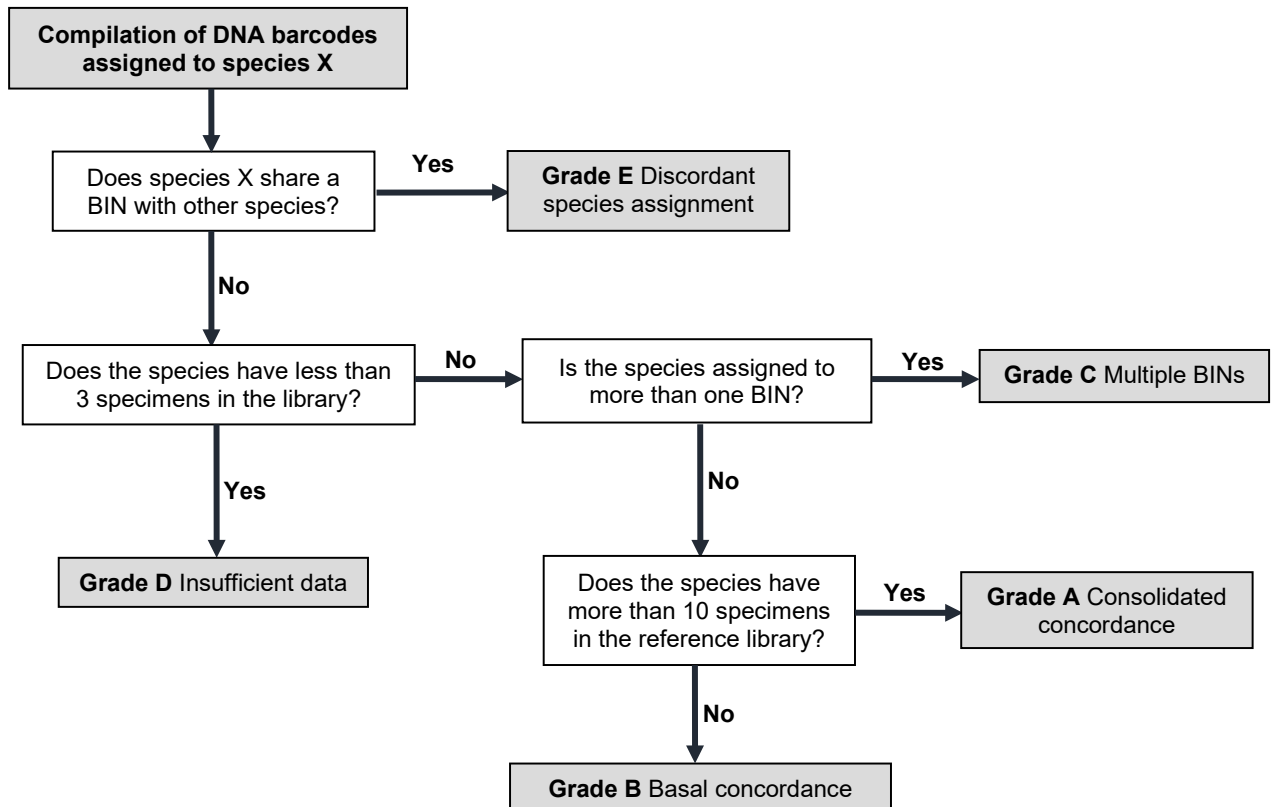
	Synonym	Ambiguous species names	Consolidated morphospecies aggregated in one BIN	Misidentification	Inconclusive
Marine Amphipoda	0.0%	30.0%	0.0%	66.7%	3.3%
Chironomidae	0.0%	33.3%	10.0%	50.0%	6.7%
Marine fish	10.0%	10.0%	26.6%	36.7%	16.7%
Overall	3.4%	24.4%	12.2%	51.1%	8.9%

621



622

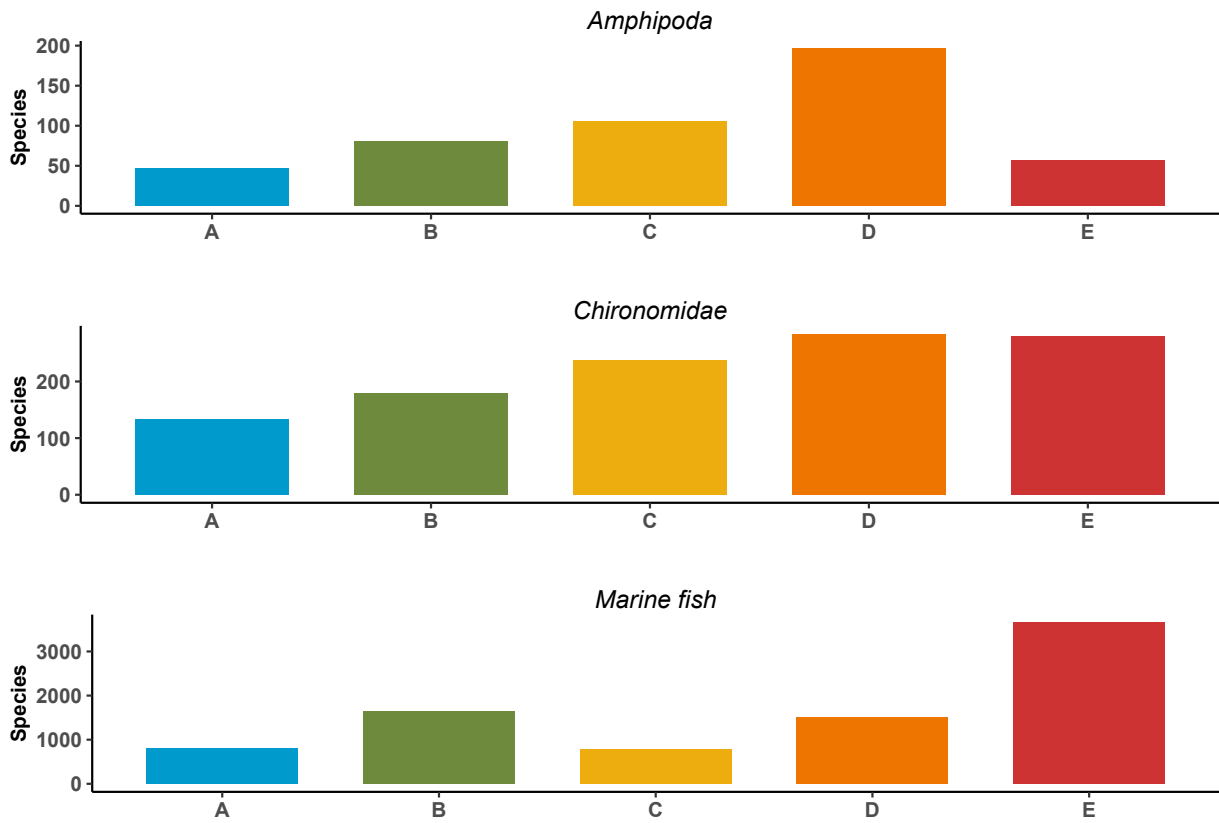
623 Figure 1 – Overview of BAGS' four main features and their arrangement along the informatics pipeline.



624

625 Figure 2 – Workflow for automated auditing and annotation of qualitative grades to each species in a BAGS-
 626 compiled reference library (adapted from Oliveira et al. 2016).

627

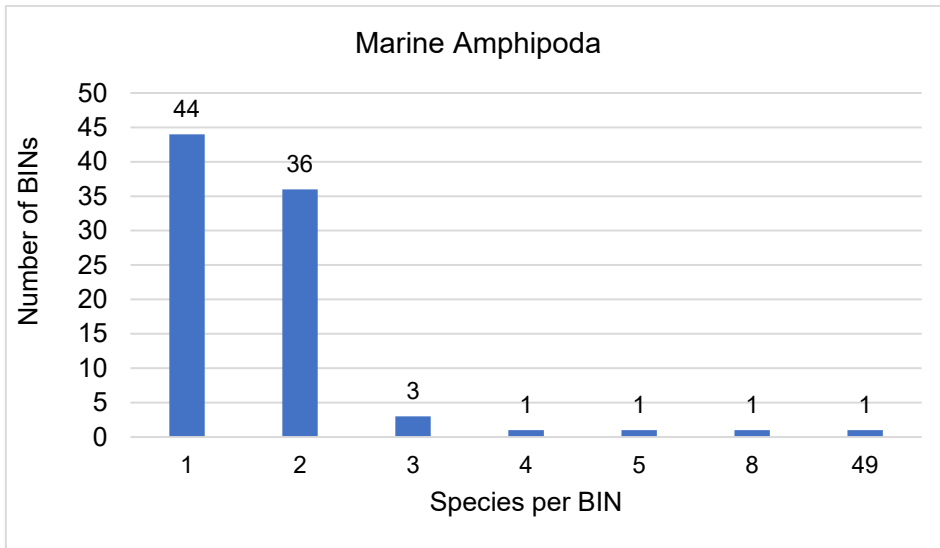


628

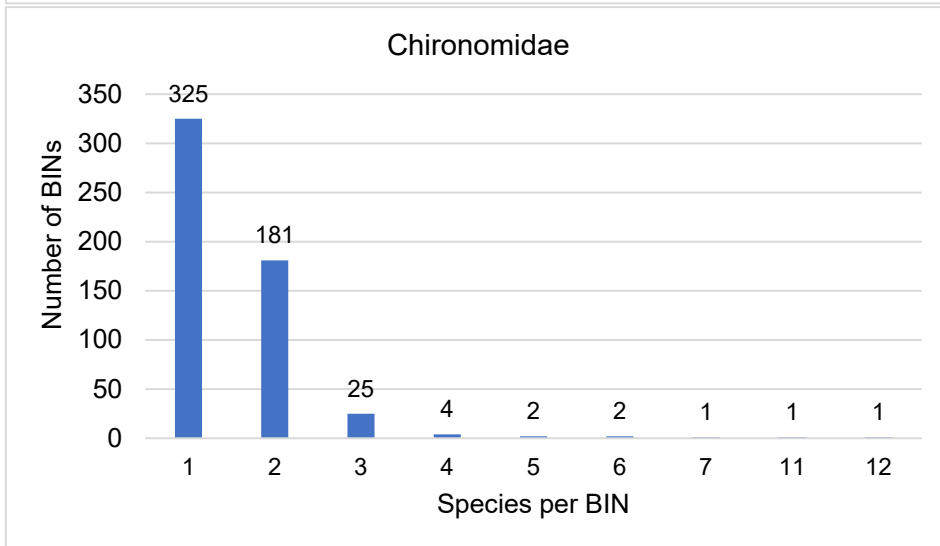
629 Figure 3 - Barplots displaying the distribution of the number of species assigned to each qualitative grade for the
 630 three taxonomic groups tested. From top to bottom: marine Amphipoda, Chironomidae and marine fish
 631 (Actinopterygii, Elasmobranchii and Holocephali).

632

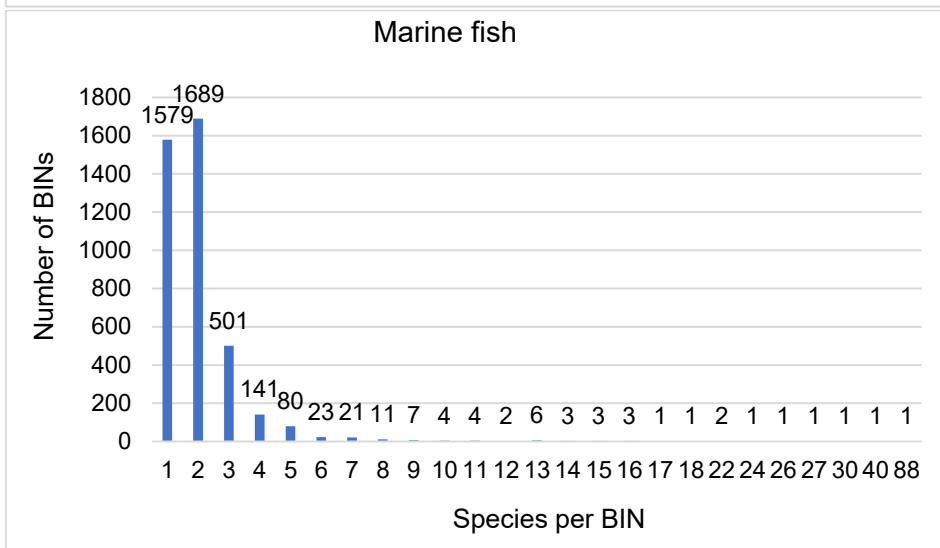
633



634



635



636 Figure 4 – Number of species per BIN in the grade E dataset generated through BAGS for each tested taxonomic
637 group (marine Amphipoda, Chironomidae, marine fish).