# Global gap-analysis of amphipod barcode library

**Anna Maria Jażdżewska** [Corresp., 1] , **Anne Helene S. Tandberg** [2] , **Tammy Horton** [3] , **Saskia Brix** [4]

[1] Department of Invertebrate Zoology and Hydrobiology, Faculty of Biology and Environmental Protection, University of Lodz, Lodz, Poland

[2] University Museum, Department of Natural History, University of Bergen, Bergen, Norway

[3] National Oceanography Centre, Southampton, United Kingdom

[4] Department for Marine Biodiversity Research (DZMB), Senckenberg am Meer, Hamburg, Germany

Corresponding Author: Anna Maria Jażdżewska
Email address: anna.jazdzewska@biol.uni.lodz.pl

In the age of global climate change and biodiversity loss there is an urgent need to provide effective and robust tools for diversity monitoring. One of the promising techniques for species identification is the use of DNA barcoding that in Metazoa utilizes the so called 'gold-standard' gene of cytochrome *c* oxidase (COI). However, the success of this method relies on the existence of trustworthy barcode libraries of the species. The Barcode of Life Data System (BOLD) aims to provide barcodes for all existing organisms, and is complemented by the Barcode Index Number (BIN) system serving as a tool for potential species recognition. Here we provide an analysis of all public COI sequences available in BOLD of the diverse and ubiquitous crustacean order Amphipoda, to identify the barcode library gaps and provide recommendations for future barcoding studies. Our gap analysis of 25702 records has shown that although 3835 BINs (indicating putative species) were recognised by BOLD, only 10 % of known amphipod species are represented by barcodes. We have identified almost equal contribution of both records (sequences) and BINs associated with freshwater and with marine realms. Three quarters of records have a complete species-level identification provided, while BINs have just 50%. Large disproportions between identification levels of BINs coming from freshwaters and the marine environment were observed, with three quarters of the former possessing a species name, and less than 40% for the latter. Moreover, the majority of BINs are represented by a very low number of sequences rendering them unreliable according to the quality control system. The geographical coverage is poor with vast areas of Africa, South America and the open ocean acting as "white gaps". Several, of the most species rich and highly abundant families of Amphipoda (e.g. Phoxocephalidae, Ampeliscidae, Caprellidae), have very poor representation in the BOLD barcode library. As a result of our study we recommend stronger effort in identification of already recognised BINs, prioritising the studies of families that are known to be important and abundant

components of particular communities, and targeted sampling programs for taxa coming from geographical regions with the least knowledge.

1 **Global gap-analysis of amphipod barcode library**

2 **Anna M. Jażdżewska[1], Anne Helene S. Tandberg[2], Tammy Horton[3], Saskia Brix[4]**

3 [1] Department of Invertebrate Zoology and Hydrobiology, Faculty of Biology and Environmental

4   Protection, University of Lodz, Lodz, Poland, author ORCID: 0000-0003-2529-0641

5 [2] University of Bergen, University Museum, Department of Natural History, Bergen, Norway,

6   author ORCID: 0000-0003-3470-587X

7 [3] National Oceanography Centre, Southampton, UK, author ORCID: 0000-0003-4250-1068

8 [4] Senckenberg am Meer, Department for Marine Biodiversity Research (DZMB), Hamburg,

9   Germany, author ORCID: 0000-0002-3269-8904

10 Corresponding Author:

11 Anna Jażdżewska[1]
12 12/16 Banacha st., 90–237 Lodz, Poland
13 Email address: anna.jazdzewska@biol.uni.lodz.pl

14

15 **Abstract**

16 In the age of global climate change and biodiversity loss there is an urgent need to provide

17 effective and robust tools for diversity monitoring. One of the promising techniques for species

18 identification is the use of DNA barcoding that in Metazoa utilizes the so called 'gold-standard'

19 gene of cytochrome *c* oxidase (COI). However, the success of this method relies on the existence

20 of trustworthy barcode libraries of the species. The Barcode of Life Data System (BOLD) aims

21 to provide barcodes for all existing organisms, and is complemented by the Barcode Index

22 Number (BIN) system serving as a tool for potential species recognition. Here we provide an

23 analysis of all public COI sequences available in BOLD of the diverse and ubiquitous crustacean

24 order Amphipoda, to identify the barcode library gaps and provide recommendations for future

25 barcoding studies. Our gap analysis of 25702 records has shown that although 3835 BINs

26    (indicating putative species) were recognised by BOLD, only 10% of known amphipod species

27    are represented by barcodes. We have identified almost equal contribution of both records

28    (sequences) and BINs associated with freshwater and with marine realms. Three quarters of

29    records have a complete species-level identification provided, while BINs have just 50%. Large

30    disproportions between identification levels of BINs coming from freshwaters and the marine

31    environment were observed, with three quarters of the former possessing a species name, and

32    less than 40% for the latter. Moreover, the majority of BINs are represented by a very low

33    number of sequences rendering them unreliable according to the quality control system. The

34    geographical coverage is poor with vast areas of Africa, South America and the open ocean

35    acting as "white gaps". Several, of the most species rich and highly abundant families of

36    Amphipoda (e.g. Phoxocephalidae, Ampeliscidae, Caprellidae), have very poor representation in

37    the BOLD barcode library. As a result of our study we recommend stronger effort in

38    identification of already recognised BINs, prioritising the studies of families that are known to be

39    important and abundant components of particular communities, and targeted sampling programs

40    for taxa coming from geographical regions with the least knowledge.

41

42    **Keywords**

43    DNA barcoding, Crustacea, marine, freshwater, semi-terrestrial, taxonomic identification

44

45    **Introduction**

46    Nature in the age of Anthropocene is facing numerous global changes and challenges. One of the

47    drastic results of human associated activities is the acceleration of species extinctions, with one

48    million species estimated to be presently critically endangered (IPBES, 2019). What is more,

49    although the rate of species discovery grows, large numbers of species remain undescribed and it

50    is believed many will not be recognized before they go extinct (Mora et al., 2011, Brix et al.

51    2020). This raises the challenge of efficient environmental monitoring, which is crucial for

52    biodiversity recognition and preservation. Monitoring based on the taxonomic identification of

53    organisms in samples is time-consuming and requires knowledge of the studied group. In the

54    time of the taxonomic impediment (Ebach, Valdecasas & Wheeler, 2011), species identification

55    methods offering an alternative to morphology-based methods are of great interest. Utilization of

56    DNA-barcoding (identifying sequences of individual specimens), metabarcoding (high-

57    throughput identification of bulk samples) and the use of environmental DNA (e-DNA,

58    identifying DNA of taxa directly from water or soil sample, without collection of specimens)

59    have been presented as promising methods in monitoring and ecological studies (e.g., Hajibabaei

60    et al., 2012; Cristescu, 2014; Aylagas et al., 2018; Leese et al., 2018; Bush et al., 2019; Feio et

61    al., 2020). The use of metabarcoding in assessing the status of ecosystems has already received

62    the new term "Biomonitoring 2.0" (Bush et al., 2019). Such approaches require the existence of

63    well-established barcode fragment libraries, which allow accurate recognition of organisms in

64    the environment (Cristescu, 2014; Cowart et al., 2015; Oliveira et al., 2016; Múrria et al., 2020).

65    Recent studies indicate that although the use of barcoding in biomonitoring has great advantages

66    over morphological identification, the current gaps in barcode libraries may hinder their use

67    (Weigand et al., 2019; Duarte, Vieira & Costa, 2020; Feio et al., 2020; Hestetun et al., 2020;

68    Leite et al., 2020; Múrria et al., 2020; Vieira et al. 2021).

69    There are two main repositories where DNA sequences are deposited: NCBI GenBank

70    (www.ncbi.nlm.nih.gov/genbank/, Sayers et al., 2020) and Barcode of Life Data System (BOLD,

71    www.boldsystems.org, Ratnasingham & Hebert, 2007). In contrast to GenBank, which

72    assembles nucleotide data of all genes, the primary aim of BOLD is to store data used for species

73    barcoding, which in the case of Metazoa is the cytochrome *c* oxidase (COI) gene. The

74    development of the BOLD database included the Barcode Index Number (BIN) system

75    implementation (Ratnasingham & Hebert, 2013) that intends to help in biodiversity assessments

76    by providing species-level taxonomic registry. Based on a molecular species delimitation

77    method, each Molecular Operational Taxonomic Unit (MOTU) recognized by BOLD receives a

78    unique alphanumeric code (BIN). Ideally, each BIN is associated with an accurate taxonomic

79    (preferably species) identification and links to the voucher stored in a recognised institution.

80    However, in practice this is not working well, and at the time of system implementation as many

81    as 46% of BINs lacked species names (Ratnasingham & Hebert, 2013). This issue has arisen for

82    a variety of reasons, which we investigate in this study using a particular faunal group, the

83    Amphipoda, as a model.

84    The Order Amphipoda are peracarid crustaceans belonging to the class Malacostraca. They are

85    very diverse components of aquatic environments. According to the World Amphipoda Database

86    (WAD, Horton et al., 2020, accessed on 17-07-2020) there are 10235 accepted amphipod

87    species, the majority of which (78%) inhabit the marine realm, around 20% are freshwater

88    species and just 2% are terrestrial taxa (Horton et al., 2020; Väinölä et al., 2008). The discovery

89    rate of new species has grown steadily since the first amphipod species description and has

90    particularly accelerated in the last six decades (Horton et al., 2020) with mean number of over

91    100 taxa annually described since the 1960s (Coleman, 2015). If the trend from the last sixty

92    years persists, we may expect to have ca. 8000 new species described by 2100. More

93    conservative estimates predict that 6100 new species will be described by that date (Arfianti,

94   Wilson & Costello, 2018). The use of molecular methods in the studies of Amphipoda has

95   revealed very high species diversity (e.g. Knox et al., 2012; Verheye, Backeljau & d'Udekem

96   d'Acoz, 2016; Tempestini, Rysgaard & Dufresne, 2018; Jażdżewska & Mamos, 2019) and

97   revealed the existence of cryptic species complexes within widely distributed taxa (Witt,

98   Therloff & Hebert, 2006; Mamos et al., 2014; Wysocka et al., 2014; Havermans, 2016).

99   Amphipoda are not only a species-rich group, but they also often dominate the crustacean

100  assemblages in which they occur (e.g., Corkum, 1989; Humphries, Davies & Mulcahy, 1996;

101  Vinogradov, Volkov & Semenova, 1996; Jazdzewski et al., 2001; Väinölä et al., 2008; Frutos,

102  Brandt & Sorbe, 2017; Brix et al., 2018; Havermans & Smetacek, 2018). They can be found in

103  both the benthos and the pelagic realm, presenting a variety of states of mobility (from

104  epibenthic clingers to fully mobile swimmers) and, as a result, possess a wide variety of feeding

105  habits including herbivory, detritivory, necrophagy, omnivory, predation and ectoparasitism

106  (Barnard & Karaman, 1991; Vinogradov, Volkov & Semenova, 1996; Dauby, Scailteur & De

107  Broyer, 2001; Väinölä et al., 2008). Being diverse and abundant they are important prey items

108  for other invertebrates and vertebrates, including fish, birds and mammals (e.g. Dalpadado et al.,

109  2001; Dauby, Nyssen & De Broyer, 2003). Certain species of Amphipoda are used in laboratory

110  ecotoxicological studies (Hyne & Everett, 1998; Bundschuh et al., 2013, Major et al., 2013).

111  Some amphipod species are well-adapted to anthropogenic environments such as artificial

112  structures used in coastal protection or are part of fouling communities, and have shown a high

113  invasion potential worldwide (e.g., Bij de Vaate et al., 2002; Kelly et al., 2006; Cabezas et al.,

114  2014; Rewicz et al., 2015; Beermann et al., 2020; Sedano et al., 2020).

115  The combined factors of high diversity and the important role played by amphipods in the

116  aquatic ecosystem highlights the need for accurate species identifications which are required for

117    biological monitoring programs. The use of DNA-barcoding may speed up the identification

118    process, but it will only succeed if the barcode library is well-established and robust. Recent gap-

119    analyses of the barcode libraries in aquatic European environments showed very large

120    differences in the coverage between different taxonomic groups and geographic regions

121    (Weigand et al., 2019; Feio et al., 2020; Hestetun et al., 2020; Leite et al., 2020; Vieira et al.

122    2021). These studies used species lists restricted to particular geographic regions or chosen

123    taxonomic groups. Basic summaries concerning the extent of amphipod data in BOLD identified

124    problems with lack of taxonomic identification or detailed geographic information as well as

125    contamination with human or bacterial DNA and provided recommendations to improve the data

126    (Radulovici & Coleman, 2017; Coleman & Radulovici, 2020). However, to date there are no

127    detailed analyses that have been conducted on a single taxon on a global scale.

128    In this study we have conducted a gap-analysis of the barcode library of a single crustacean

129    order, the Amphipoda, on a global basis. In producing an up-to-date picture of the current state

130    of knowledge, we will provide researchers with a detailed understanding of the both the strengths

131    and the potential limitations of the use of DNA barcodes for identifications. We also propose

132    recommendations for future initiatives that involve molecular data and produce new barcodes to

133    fill the gaps in our knowledge of this taxon.

134    **Material and methods**

135    Data for the present study were retrieved from BOLD by searching the "Public Data Portal"

136    using the keyword "Amphipoda". A combined dataset of all records was downloaded as an .xml

137    file on June 24th 2020.

138    All records of the barcoding fragment of the cytochrome c oxidase I (COI-5P in BOLD) were

139    extracted (29016 records). This extracted dataset was used for all further analyses conducted by

140    using various filtering options in an Excel spreadsheet. 2579 records, represented by sequences

141    shorter than 500 bp or having more than 1% ambiguous nucleotides for which BINs were not

142    ascribed, were removed from dataset. Continued analysis of the dataset revealed some duplicate

143    records (1468 records, 734 cases, Supplemental file 1). These derived from data harvested by

144    BOLD from GenBank and seemed to be associated with an update of the records in GenBank. In

145    the dataset, these records had an identical sample ID that referred to a GenBank Accession

146    Number but with an additional '.1' appended (e.g. KP713892 and KP713892.1) and with an

147    identical identification provided. The differences were often linked with more detailed

148    geographical information in the case of one record from the pair. Only the more detailed entry

149    was retained for continued analysis. One sequence of *Niphargus novomestanus* S. Karaman,

150    1952 (KR858496, BOLD:ADD1128) was removed from the dataset because it was deleted from

151    GenBank by its submitter ("This record was removed at the submitter's request because the

152    source organism cannot be confirmed." GenBank website). The resulting dataset contained

153    25702 records (Fig. 1, Supplemental file 2).

154    Each record in the dataset was then further refined by sorting into categories according to the

155    level of taxonomic identification. The following categories were used: order, family, subfamily,

156    genus and species. Where records were provided with a temporary species identification, i.e.

157    they are recognised as separate morphospecies but are not determined to correspond to a known

158    taxon - they were treated as a separate category. In the whole dataset ca. 2.5% of records (596

159    individuals, 145 BINs) had uncertain identification with "cf." or "aff". Because the majority of

160    them (417 records, 101 BINs) were associated with five species of one genus (*Gammarus*) for

161    simplification all such records were treated as final species identifications. However, it is

162    understood that the use of open nomenclature, when applied to identifications, provides an

163    indication of the level of uncertainty, and may be intended to indicate the presence of new

164    species or species complexes.

165    The data in BOLD come from wide variety of projects, some of which involve detailed

166    taxonomic study by specialists, others are focused on monitoring or other topics in which

167    taxonomic specialists are not involved. For the purposes of our analyses it was assumed that the

168    identification accuracy was equal throughout the whole dataset, regardless of its origin. In

169    several cases identification of the specimens within a single BIN varied strongly, with some

170    records remaining at order level while others were determined to the species level. BINs aim to

171    represent a putative species, so in the above example, the most detailed taxonomic information

172    was applied to all records within the single BIN. Sometimes multiple (most often two) species or

173    genus names were associated with a single BIN (87 cases). Each of these cases was checked

174    individually. Sometimes it was an obvious misidentification of a single individual within a large

175    group - if this was noted the misidentified record was added as an additional element to the

176    records identified to the lowest congruent level (e.g. if the genus name matched the BIN genus,

177    the misidentified taxon was added as an additional record identified to the genus level, if the

178    lowest congruent level was family it was added to the family records); and the taxon

179    identification of the majority of records was applied as correct. When it was impossible to judge

180    which name was correct, the name of the identifier was checked and identifications carried out

181    by taxonomists specializing in Amphipoda was prioritised over a that provided by a non-

182    specialist study. Where this process did not give a satisfactory conclusion, the BIN was allocated

183    an identification at a rank that was congruent for the different records. The list of taxa with

184    incongruent identifications together with an explanation of the final decision is presented in

185    Supplemental file 3.

186    Based on the taxonomic identification of the records the associated BINs were divided into the

187    following environmental categories:

188    a) marine

189    b) freshwater

190    c) terrestrial.

191    Taxa that inhabit both marine realm and brackish environments were allocated to the marine

192    category. Taxa from freshwater also occurring in brackish waters were allocated to the

193    freshwater category. All representatives of the family Talitridae were treated as terrestrial taxa.

194    Where taxonomic information was not detailed enough to provide environmental information

195    about the particular BIN, the geographic data (coordinates and/or locality description) of the

196    associated records were used to ascribe a particular BIN to one of the above categories. In some

197    cases, this necessitated checking the original publication. A small number of unallocated BINs

198    (18) and associated records (44) were used only in the first general summary of amphipod

199    barcodes, but they were removed from further analyses (Supplemental file 4).

200    In order to verify the correct environmental allocation of BINs, all BINs with records possessing

201    coordinates were plotted on a map using the software QGIS2.16.1 (QGIS Development Team,

202    2018). Cases where incongruence between the ascribed environment and the geographic position

203    appeared were checked individually. For those records without detailed geographic information

204    the country of origin was taken from either BOLD or the associated publication.

205    In order to verify the barcode coverage within the studied group a list of BINs associated with a

206    species name was compared with the list of accepted amphipod species names available in the

207    World Amphipoda Database (WAD, Horton et al., 2020, accessed on 17-07-2020). A barcode

208    quality assessment of the species represented in BOLD, based on the grading system proposed

209     by Oliveira et al. (2016) and slightly modified by Fontes et al. (2020) was applied. This system

210     consists of five grades: A – consolidated concordance (>10 sequences of a single morphospecies

211     grouped in a single BIN), B – basal concordance (same as grade A but between three and 10

212     sequences available in the library), C – multiple BINs (one morphospecies assigned to more than

213     one BIN), D – insufficient data (single species is assigned to single BIN but it is represented by

214     less than three sequences in the barcode library), E – discordant species assignment (more than

215     one species assigned to a single BIN). Fontes et al. (2020) provide an R-based application

216     (Barcode, Audit & Grad System – BAGS), and uses only those records possessing species

217     names. Since our aim was to focus on all available barcode records (including sequences

218     identified only to higher ranks), the assessment was carried out manually. Additionally, as a

219     result of initial treatment of the dataset, misidentified species records or BINs with unclear

220     species identification, were already removed, so category E (discordant species assignment,

221     Oliveira et al., 2016; Fontes et al., 2020) was not recorded. For the purpose of the present study

222     Lysianassoidea *incertae sedis* was treated as an additional family. The amphipod families were

223     divided into four categories depending on the number of species in each: low species rich

224     families (up to 10 species), moderately species rich families (from 11 to 30 species), species rich

225     families (31-100 species), very species rich families (more than 100 species). This division

226     allowed verification of pattern between the species richness of the family and its representation

227     in BOLD.

228     **Results**

229     Of the 25702 amphipod COI records, 46.5% (11958 records) were freshwater, 43.5% (11169

230     records) were from the marine realm, and 9.8% (2531 records) were terrestrial taxa. Of the 3835

231     recognized BINs in total, 45% (1726 BINs) belonged to freshwater taxa, 50% (1920 BINs) were

232   marine, and 4.5% (171 BINs) were from terrestrial taxa. 44 records (0.2%) and their associated

233   18 BINs (0.5%) could not be ascribed to the above environmental categories and were not

234   considered further (Fig. 2A, B).

235   More than half (57.5%) of the records available in BOLD possessed coordinates, and 20% had

236   information about the country of origin. Geographic information about the remaining 22.5% was

237   provided only in the original publication. Geographic information is more comprehensive for

238   marine taxa, where 71% of records possessed coordinates (compared to 47% for freshwater, and

239   50% for terrestrial taxa). Molecular studies of freshwater Amphipoda are focused mainly in the

240   Northern hemisphere (particularly European countries, Russia and United States) while in the

241   Southern hemisphere, Australia, New Zealand and Argentina are well studied (Fig. 3A). There is

242   a complete lack of records (amphipod sequences) from Brazil, equatorial America and vast areas

243   of Africa. Similar patterns of data coverage were seen for marine amphipods, which have greater

244   numbers of records along European, North American and East Asian coasts. In the Southern

245   hemisphere, Australia, New Zealand and Antarctica had larger numbers of barcode records (Fig.

246   3B). However, vast areas of the deep sea and the Arctic Ocean remain undersampled. Terrestrial

247   Amphipoda in Europe, North America, China, Australia and Chile were the best represented

248   (Fig. 3C), but sampling gaps were seen in the continents of South America and Africa.

249   The majority of records (69.8%, 17922 recs.) had a complete species-level identification. Of the

250   remaining 30.2% of records, 5.6% (1433 recs.) had received temporary names (open

251   nomenclature), 11.3% (2902 recs.) remained identified at the genus level, 0.2% (40 recs.) at

252   subfamily, 5.0% (1285 recs.) at family, and 8.1% (2076 recs.) at the order level. Levels of

253   identification varied according to the environment, with marine taxa having greater proportions

254   of taxa identified only to higher taxonomic ranks (Fig. 4A). The majority of BINs (3817) were

255   associated with species names (55.7%, 2126 BINs). These were followed by BINs identified to

256   the order level (13.3%, 506 BINs), generic or family level (10.7%, 407 BINs each) and those

257   with a temporary name (9.4%, 359). BINs with only a subfamily name constituted just 0.3%

258   (12). Greater variations between environments were seen for the BINs, with 74% (1284) of

259   freshwater BINs having a species level identification, compared to only 39% (751) of marine

260   BINs (Fig. 4B). More than 20% (444, 23%) of the BINs for marine taxa remained identified at

261   the order level.

262   Regardless of the environmental origin, the majority of BINs were represented by a single

263   sequence (Fig. 5). BINs represented by five or fewer sequences constituted around two thirds

264   (67%, 114 terrestrial BINs to three quarters, 78%, 1488 marine BINs) of BINs recorded in a

265   particular environment. Freshwater taxa had 41 BINs (2.4%) represented by more than 50

266   sequences, compared to 28 (1.5%) for marine taxa, and eight (4.7%) for terrestrial taxa. When

267   only those BINs with complete species-level identifications are considered, the proportion of

268   sequences representing a particular MOTU does not change, with freshwater taxa having 78% of

269   BINs (1016) represented by five or fewer sequences. Almost three quarters of marine BINs

270   (71%, 525 BINs) had five or fewer sequences in BOLD, while this proportion was 61% (56

271   BINs) for terrestrial taxa. Freshwater taxa had 35 BINs (3%) represented by more than 50

272   sequences, compared to 27 (4%) for marine taxa, and 6 (7%) for terrestrial taxa. The best

273   represented BIN in BOLD (801 sequences) belonged to the terrestrial species *Orchestoidea*

274   *tuberculata* Nicolet, 1849 (BOLD:ACQ3380), followed by the marine species *Gammarus*

275   *oceanicus* Segerstråle, 1947 (BOLD:AAA1262, 553 sequences), and the freshwater species

276   *Diporeia hoyi* (S.I. Smith, 1874) (BOLD:AAA1473, 512 sequences). A further 26 BINs were

277     represented by more than 100 sequences, including 17 freshwater, seven marine and two

278     terrestrial BINs (Supplemental file 5).

279     Out of the 3817 studied BINs, just over half (55.7%, 2126) were associated with a species-level

280     identification, representing 1001 species. Freshwater BINs with species identification reached

281     1284, associated with 453 species, while 751 marine BINs were determined to 496 species. Of

282     the 91 terrestrial BINs, 52 species were identified. Generally, a single morphological species was

283     associated with each BIN (68%, 680 cases, 288 in freshwater, 359 marine, 33 terrestrial). 17% of

284     the identified species were associated with two different BINs (72 freshwater, 82 marine and 14

285     terrestrial) (Fig. 6). There were however 19 cases when one single morphological species was

286     represented by more than 10 BINs (17 freshwater, one marine and one terrestrial) (Supplemental

287     file 6). The greatest number of BINs was recorded for the freshwater species *Gammarus*

288     *balcanicus* Schäferna, 1923 represented by 143 BINs (45 BINs were identified as "cf." or "aff.")

289     followed by another freshwater taxon *Hyallella azteca* (Saussure, 1858) (62 BINs) and

290     *Gammarus fossarum* Koch, 1836 (51 BINs; 19 BINs identified as "cf." or "aff."). Among

291     terrestrial taxa the highest molecular variation (12 BINs) was recorded for *Morinoia japonica*

292     (Tattersall, 1922) (present in BOLD under former generic name *Platorchestia*), while *Apohyale*

293     *stebbingi* Chevreux, 1888 (with 11 BINs recognized) was the most diverse among marine

294     species.

295     Of the 239 accepted families of Amphipoda (238 families and Lysianassoidea *incertae sedis*),

296     105 (44%) were represented by at least one species in BOLD (Table 1). The largest number of

297     families had up to 20% of species barcoded, while only ten families had more than half of the

298     known species barcoded (Supplemental file 7). Thirteen families lacking barcoded species had at

299   least one barcoded taxon identified at the genus level, a further five families had a taxon

300   identified at the family level.

301   Just under ten percent (999 spp., 9.7%) of the 10330 accepted species of Amphipoda (Horton et

302   al., 2020) had barcodes. Of the nominal species possessing barcodes almost 500 (496 spp.) are

303   marine, 451 spp. are freshwater and 52 spp. are terrestrial taxa. The data coverage of the majority

304   of species, no matter their environmental origin, is not sufficient for the barcodes to be trusted

305   according to the quality control system (Table 2) (Oliveira et al., 2016; Fontes et al., 2020).

306   Additionally, a large group of taxa is represented by multiple BINs; only 10% of species

307   represent consolidated concordance of available barcodes.

308   The breakdown of amphipod families according to the assigned categories of richness and their

309   respective representation in BOLD can be seen in Table 3. Almost every one of the very species

310   rich families had at least one species barcoded (31 families out of 32), and 22 of 30 species rich

311   families are represented in BOLD. For both moderately low and low species rich families 26

312   possessed at least one representative in BOLD constituting respectively 48% and 21% of all

313   families each (Supplemental file 7). The mean coverage of barcodes for species in each of the

314   above groups was around 10% with the highest observed for low species rich families (12%) and

315   the lowest (8%) recorded for families grouping from 30 to 100 species. However, if the families

316   without any molecular information were removed from the study these numbers considerably

317   change. The low species rich families (1-10 spp.) had a barcode coverage at the level of 49%,

318   moderately species rich families (11-30 spp.) reached 21% of coverage, while the rich and very

319   rich amphipod families (more than 30 spp.) had only 9-10% of species studied.

320   A third of families (34) have at least one species characterized by consolidated concordance of

321   available barcodes (category A of the quality grading system). Another third of families (38) do

322   not have any species in categories A or B, indicating that the species already studied represent a

323   potential cryptic diversity or the available data are insufficient (Table 4).

324   Within the very species rich families, the best representation in BOLD was recorded for

325   Niphargidae (36.5% of known species represented with a barcode), Gammaridae (31%) and

326   Crangonyctidae (16%). Only the family Stegocephalidae did not have any representative with

327   species level identification (although barcodes belonging to this family but identified at genus

328   level were present). The least studied families within this group (but having at least one species

329   barcode) were: Phoxocephalidae (1% of the species with a barcode), Dexaminidae,

330   Liljeborgiidae and Maeridae (ca. 2% of the species with a barcode). Among species rich families

331   41% of the species from Pseudoniphargidae had barcodes, while the Epimeriidae and

332   Pontogammaridae had 20% and 19%, respectively. The best represented moderately species rich

333   families were Metacrangonyctidae, Oxycephalidae and Hyperiidae with 55%, 50% and 48% of

334   the associated species represented with a barcode respectively. Within low species rich families

335   four (Baikalogammaridae, Crymostygidae, Cyllopodidae and Tryphanidae) had all known

336   species represented with barcodes, but other than Cyllopodidae (two species) the families are

337   monotypic (Supplemental file 7).

338   **Discussion**

339   **Extent of barcode library of Amphipoda**

340   One of the aims of establishing the BOLD database was to store and publish barcodes, based on

341   records uploaded by its users and supplemented by the data harvested from GenBank

342   (Ratnasingham & Hebert, 2007). Together with the BIN system, that groups similar sequences in

343   clusters representing putative species (Ratnasingham & Hebert, 2013), the BOLD database aids

344   in recognising and quantifying biodiversity. The extent of data in BOLD expresses the activity of

345  researchers studying particular groups using molecular methods. The number of available

346  sequences of Amphipoda in BOLD is comparatively large. At the time of download (end of June

347  2020) they were represented by almost 26 000 records (3835 BINs), and by the end of August

348  there were more than 34 000 public sequences (3914 BINs) (BOLD accessed on 20-08-2020),

349  indicating the great intensity of molecular studies involving this crustacean group, and that the

350  data in BOLD are actively growing. Among other crustacean groups only Decapoda is

351  represented by a higher number of records (64 281 records). Copepoda are represented by 18

352  511, Thecostraca by 15 554, Isopoda by 13 858 and Branchiopoda by 12 326 sequences. The

353  large number of identified BINs within the Amphipoda also places this group second only to

354  Decapoda (with 6056 BINs). Isopods and copepods are represented by 1853 and 1804 BINs,

355  respectively, while 969 BINs were identified within Branchiopoda. Within Thecostraca only 545

356  BINs were identified (boldsystems.org, accessed on 20-08-2020).

357  When the BIN system was implemented, Ratnasingham & Hebert (2013) indicated that 12% of

358  all the sequences available in BOLD lacked a family name, 19% a genus name and 40% a

359  species name. A comparison of these numbers with the present data on Amphipoda looks

360  optimistic, where only 8% of sequences are without family indication, 13% are without genus

361  and 29% lack a species identification. However, the global analysis of Ratnasingham & Hebert

362  (2013) identified 10% of BINs lacking family names, almost 24% lacking generic names and

363  46% lacking species names. These numbers are almost identical for amphipod BINs known

364  presently (13%, 23%, 43% of BINs lacking family, genus and species information, respectively).

365  Among all known species of Amphipoda, almost 80% of species are marine, some 20% live in

366  freshwaters, while 2% may be considered as terrestrial (Horton et al., 2020; Väinölä et al., 2008).

367  The above proportions are expressed neither in the number of records nor the number of

368    recognized BINs that are more or less evenly distributed between freshwater and marine taxa.

369    This demonstrates that in terms of amphipod crustaceans freshwater taxa are much better studied

370    than the marine taxa. These disproportions are even more striking when the level of

371    identification of sequences and BINs is considered. Although the majority of data in BOLD

372    possess species-level identifications, marine amphipods are less thoroughly identified. This is

373    especially clear for marine BINs, of which only 39% had species-level identifications, while as

374    much as one fifth are identified only as "Amphipoda". The fact that freshwater amphipods are

375    better studied is not surprising considering the easier access to this environment. In the case of

376    marine fauna, obtaining samples suitable for molecular analysis can be challenging, especially

377    when extreme habitats (polar regions, deep-sea, hydrothermal vents etc.) are considered (e.g.,

378    Riehl et al., 2014; Jażdżewska & Mamos, 2019). Additionally, rarity is a common feature of

379    numerous marine species (particularly in the deep-sea environment, see Kaiser et al., [2007]),

380    where many taxa are known only from their original descriptions and type localities (Jażdżewska

381    & Mamos, 2019). The question of how many of the BINs not associated with a species

382    identification actually belong to already known species is also of concern. In these cases, it is

383    highly advisable to put every effort to identify the already available material – this will relatively

384    efficiently improve data usability. Taxa that are associated with a BIN, yet are known to be new

385    to science are another cause for concern. This is particularly evident for marine taxa collected

386    during recent deep-sea exploration programs (e.g., Brandt et al., 2007; Jażdżewska, 2015; Brandt

387    et al., 2019; Brix et al., 2020). It is imperative that full scientific descriptions of new species are

388    produced to reduce the current proliferation of 'dark taxa' (Page, 2016).

389    The geographic distribution of available amphipod sequence records shows clear sampling gaps.

390    In particular the African continent, the northern part of South America and the Coral Triangle in

391  Asia are complete "white spots" when freshwater and terrestrial taxa are considered. For marine

392  species, the coasts of Africa and South America, the Coral Triangle, and large parts of the deep

393  sea of all oceans, lack coverage. Considering the known high species diversity of these regions it

394  will be necessary to establish targeting sampling programs before we can consider that we have

395  adequate global coverage of the molecular diversity of the Amphipoda.

396  Our study shows that globally the barcoding coverage of amphipod species is only about 10%. In

397  comparison, just over 20% of all species registered in the European Register of Marine Species

398  (ERMS) and almost 50% of species listed in the AZTI Marine Biotic Index (AMBI) have been

399  barcoded (Weigand et al., 2019). The percentage of barcoded European freshwater invertebrates

400  used in environmental monitoring reaches 64.5%, and when considering only Peracarida, 24% of

401  ERMS species, 45% of AMBI and 82% of freshwater monitored taxa have been barcoded

402  (Weigand et al., 2019). It has to be emphasized however, that only ERMS lists all marine

403  invertebrates from European region, while both other datasets studied by Weigand et al. (2019)

404  consists of a subset of species from this area. More specific studies of Iberian macroinvertebrates

405  revealed that ca. 40% of amphipod species possess barcodes (Leite et al., 2020; Múrria et al.,

406  2020). Hestetun et al. (2020) conducted a barcode library gap-analysis of the benthic macrofauna

407  of one region of the North Sea, which indicated the barcode coverage varying from 42.4% to

408  61% (depending on the calculation method) while Vieira et al. (2021) found that in Macaronesia

409  34.2% to 72.6% of macroinvertebrate species have barcode representation with much better

410  coverage of non-indigenous taxa in comparison to the native ones. This indicates that for smaller

411  subset of taxa and specified geographic region it is much easier to produce good barcode

412  coverage. It can be concluded that although Amphipoda are an actively studied taxonomic group

413    where scientists increasingly use molecular methods, this diverse and abundant macrofaunal

414    taxon is still insufficiently represented in the BOLD barcode library.

415    **Quality of amphipod barcode library**

416    In order to provide a trusted barcode for a particular species, at least one good quality sequence

417    associated with a species-level identification provided by taxonomic specialist is required as an

418    absolute minimum. However, a single sequence cannot provide information about intraspecific

419    variation, and overlooked contamination of the sample will mean the sequence cannot be

420    validated. As such, it is advisable to provide a small number of sequences to characterise each

421    taxon. The recently proposed barcode quality auditing system suggests providing at least three

422    sequences to enable proper barcode evaluation (Oliveira et al., 2016; Fontes et al., 2020).

423    Unfortunately, as we have shown in the case of Amphipoda, globally more than half of BINs are

424    represented by only 1-2 sequences in BOLD. This low number of sequences places them in

425    category D of the Oliveira et al. (2016) system, indicating the existing data is insufficient for use

426    as trusted barcodes. Similar observations for a restricted amphipod dataset are made by Fontes et

427    al. (2020).

428    Due to methodological differences it is impossible to make direct comparisons of our data with

429    the results of the gap analysis of aquatic organisms in European waters (Weigand et al., 2019),

430    but re-calculation of their data shows much improved barcode coverage. Among all freshwater

431    invertebrates 65% of taxa barcoded are represented by more than five sequences, while this

432    percentage rises to 77% when considering only freshwater Peracarida. This proportion of high

433    quality datasets diminishes when marine taxa are considered; with 52% of the marine species

434    from the AMBI list and 45% those listed in ERMS having at least five barcodes available. These

435    numbers do not change when considering only marine Peracarida (52% and 46% of the ones

436   presented in AMBI and ERMS lists, respectively). Our analysis of Amphipoda shows opposite

437   pattern with about 1/4-1/3 of BINs represented by more than five sequences but the good

438   barcode coverage observed by Weigand et al. (2019) may be biased by the fact that they targeted

439   the species used in water quality assessment programs. Because of their practical use such taxa

440   receive more scientific interest and it may be assumed that their barcoding is prioritized by

441   different institutions.

442   The amphipod BINs that have the largest numbers of sequences in BOLD are often the result of

443   detailed studies of targeted species, which have produced large numbers of sequences as a

444   secondary aim of the study. For example, 750 out of the 801 sequences in BOLD of terrestrial

445   *Orchestoidea tuberculata* come from a single study by Brante et al. (2019); 406 records out of

446   411sequences in BOLD of freshwater *Dikerogammarus haemobaphes* (Eichwald, 1841) come

447   from Jażdżewska et al. (2020); while 232 records of 235 sequences in BOLD of marine *Caprella*

448   *scaura* Templeton, 1836 come from Cabezas et al. (2014). The disproportional representation

449   between the few species that are very thoroughly studied and the remaining majority of species

450   that are represented only by a single, or a low number of sequences emphasises the need for

451   more targeted sampling of less common species.

452   **Best studied families and cryptic diversity**

453   Almost half of the 239 known amphipod families are represented in BOLD. However, only ten

454   of these families have more than 50% of their associated species sequenced. It is important to

455   underline that there are 18 families in BOLD that do not have species-level identifications, but

456   have records left at the family or genus level. A small effort to provide trusted species-level

457   identifications for these taxa will greatly improve barcode coverage of the Amphipoda,

458   particularly if they represent species already known to science.

459    Another concern that has arisen as part of this study relates to the format of temporary names in

460    GenBank and BOLD, the different requirements by users for their input, and how this has

461    changed following development of the databases. In GenBank, the incorporation of temporary

462    names or codes is allowed (referred to as placeholder names in GenBank). In 2010, a large

463    amount of COI data was incorporated into the BOLD database. The identifications associated

464    with each of these imported sequences were included verbatim from GenBank. BOLD users,

465    however, were originally able to use temporary names in the database only in private

466    projects/dataset and when opening their data for public they were expected to provide the

467    identification to the lowest taxonomic level possible (e.g. genus) and to provide the temporary

468    name (e.g. incorporating "cf." or "aff.") as a taxonomy note (that has happened to the authors of

469    the present paper). However, in BOLD a taxonomy note is only visible when the specimen page

470    is open, and not in a general search. Recently, we have learnt that open nomenclature identifiers

471    (such as 'cf.' and 'aff.') are accepted by BOLD, but it may be assumed that numerous records

472    remain at a higher taxonomic level, with more detailed identifications available that are hidden

473    from general searches. This discrepancy in dealing with temporary names has become apparent

474    when analysing the whole dataset as part of this study. In particular, the inconsistent use of

475    temporary names in these databases mean that it is very difficult to differentiate between

476    temporary names which are being used to refer to species that are new to science, and those

477    which have remained at a higher taxonomic level because they were simply not identified further

478    (which could be for a variety of reasons). Molecularly well-defined temporary names for new

479    species are likely to become more abundant and therefore critical to our knowledge of

480    biodiversity in the coming years, and we need to ensure they are managed carefully and

481    consistently. Recommendations for the use of open nomenclature have been proposed recently to

482    attempt to standardise and overcome these issues (Sigovini, Keppel & Tagliapietra, 2016; Horton

483    et al., 2021) and it is hoped that these standard formats will be considered for use in both BOLD

484    and GenBank.

485    Barcode coverage of families varies depending on the species richness. For species rich families

486    it is around 10%, while coverage is increased for moderate and low species rich taxa. This is not

487    surprising considering it is much easier to receive better coverage for monotypic families or

488    those represented by only a few species. The best studied families are the ones that remain under

489    the interest of large working groups who focus on studying specific families (e.g. Hou, Fu & Li,

490    2007; Mamos et al., 2014; Wysocka et al., 2014; Delić et al., 2017a; Delić et al., 2017b; Fišer et

491    al., 2017; Copilaş-Ciocianu, Sidorov & Gontcharov, 2019). It is worth noting that providing

492    barcodes is generally more a "by-product" of other analyses than the goal per se. Another issue

493    that should be emphasized is that species rich families are proportionally under studied. This is

494    important because they usually do not only group many species but very often the species from

495    these families constitute the majority of amphipods characterizing different assemblages. This is

496    clearly shown by the Phoxocephalidae (1% of the 367 known species are barcoded),

497    Ampeliscidae (7% of 312 of the known species barcoded) or Oedicerotidae (10% of the 246

498    known species barcoded), all constitute very large and important components of marine benthic

499    communities worldwide (Brandt, 1993; Weisshappel & Svavarsson, 1998; Frutos & Sorbe, 2017;

500    Brix et al., 2018). Another example is provided by the Caprellidae (6% of the 443 known species

501    are barcoded) which are an important part of many fouling communities (e.g., Ros, Vázquez-

502    Luis & Guerra-García, 2013; Ros et al., 2013) and where proper species identifications are

503    crucial in the context of growing transport with their resulting potential alien species invasions

504 (op. cit.). The studies of these families should be prioritized in order to support marine

505 monitoring programs based on barcode libraries.

506 The analysis of the amphipod BINs with a species-level identification showed that there were

507 only a few cases where multiple names were associated with a single molecular unit. A quarter

508 of these cases resulted from the misidentification of single individuals within a taxon. In some

509 cases different names were associated with the description of new species (present in the

510 database under both former and newly established name). Problems with morphological

511 identification of cryptic species and the lack of well-established diagnostic characters within

512 closely related species may also be the reason of the presence of multiple names for single BIN.

513 The above problems have been recognized within *Gammarus ochridensis* Schäferna, 1926

514 species complex that is the group of morphologically very similar species of which two

515 *Gammarus cryptosalemaai* Grabowski, Wysocka & Mamos, 2017 and *Gammarus*

516 *cryptoparechiniformis* Grabowski, Wysocka & Mamos, 2017 are recognizable only based on

517 molecular data (Wysocka et al., 2013; Grabowski, Wysocka & Mamos, 2017). This indicates that

518 generally BOLD may be considered a trusted tool for species identification. Our analyses

519 showed that in the majority of cases, a single BIN was characterising a single species, which is

520 congruent with the results of other similar studies (Fontes et al., 2020; Leite et al., 2020). Some

521 morphologically identified species were represented by two or even three BINs, which can

522 indicate overlooked diversity. It has been noted however, that sometimes due to the methodology

523 used during BIN-identification and the threshold used (2% of similarity, Ratnasingham &

524 Hebert, 2013) some valid species may be split into two or more BINs (Lörz, Jażdżewska &

525 Brandt, 2018; Jażdżewska & Mamos, 2019). This happens more frequently when the sample size

526 is small and the intraspecific variation range cannot be adequately assessed. In such cases, the

527    use of additional genes or other data analysing methods may help to decide the proper species

528    delineation. The present study revealed 19 morphological species that were represented by 11 or

529    more BINs. This multi-BIN representation was much more common in freshwater environments,

530    where 17 species with potential cryptic diversity were observed. The existence of such high

531    cryptic diversity especially in European waters was recognized by authors of the original works

532    (Witt, Threloff & Hebert, 2006; Bauzà-Ribot et al., 2011; Major et al., 2013; Mamos et al., 2014;

533    Wysocka et al., 2014; Delić et al., 2017a; Delić et al., 2017b, Fišer et al., 2017; Tomikawa et al.,

534    2017) most recently confirmed also by Wattier et al. (2020). A detailed study of the available

535    barcodes and cryptic diversity of the Gammaridae and other representatives of the superfamily

536    Gammaroidea is in preparation (Mamos pers. comm.). The large representation of freshwater

537    taxa forming cryptic species complexes (especially in Europe) can be partly explained by the

538    geological events that shaped the European freshwater system (Wysocka et al., 2014; Mamos et

539    al., 2016; Wattier et al., 2020). Presence of marine cryptic or pseudo-cryptic species have also

540    been reported (Havermans, 2016; Verheye, Backeljau & d'Udekem d'Acoz, 2016), but the extent

541    of molecular studies of amphipods from this realm is much smaller and as a result cryptic species

542    may have been overlooked. A study of the marine genus *Apohyale* showed high diversification

543    of species within the genus, and confirms that more studies are required to correctly identify

544    species diversity and uncover cryptic diversity in marine taxa (Desiderato et al., 2019). Cases of

545    highly diverse nominative species usually come from studies based in a single research group

546    that was already aware of the high diversity within the taxon. There are, however, cases where

547    multiple BINs have received the same identification but this was carried out by different authors

548    at different times (without comparison of the material) and it is difficult to judge if the observed

549    diversity is a result of the existence of a cryptic species or of misidentification of the species. In

550    such cases it is impossible to decide which of the BINs represents the known species and which

551    are cryptic/new species that require more detailed study (Jażdżewska et al., 2018). A detailed

552    analysis of the species represented by several BINs was not the focus of the present study, but it

553    should be a priority for BOLD to identify such cases and inform users about the presence of

554    possible cryptic taxa. Users of BOLD who seek to obtain identification for their own sequence

555    should be notified that the specimen they have may belong to a group of cryptic species so that

556    the taxonomic identification can be treated with caution. Some initiatives to improve the curation

557    of BOLD data have already begun (Radulovici et al. 2021) and it is  highly recommended that

558    mistakes or problematic issues that are found in the database are corrected and published e.g. the

559    case of *Hyperiella antarctica* Bovallius, 1887/*H. dilatata* Stebbing, 1888 which was recently

560    clarified by Havermans et al. (2019).

561    **Conclusions and Future recommendations**

562    We have conducted a gap-analysis of the barcode library using a single crustacean order, the

563    Amphipoda, as a model. The high diversity and the important role played by amphipods in the

564    aquatic ecosystem combine to highlight the need for accurate species identifications which are

565    required for biological monitoring programs. DNA-barcoding may speed up the identification

566    process, but success is dependent on the barcode library coverage and quality. Our gap analysis

567    has shown that although a large number of BINs (indicating putative species) was recognized by

568    BOLD still only 10 % of the amphipod species are represented by barcodes. Moreover, the

569    majority of BINs is represent by a very low number of sequences that make them unreliable

570    according to the quality control system. The geographical coverage is poor with vast areas of

571    Africa, South America and the open ocean acting as "white gaps", also the level of barcoding

572    effort is skewed depending on the environment.

573 As such, we make the following recommendations (in order of priority), which will improve the

574 data currently held within BOLD, and we outline steps that are needed to provide a more equal

575 coverage of the sequence data within the Amphipoda, and thus improve the utility of the

576 database for a variety of applications, including species identification and biomonitoring.

577 1. Morphological identification of the already recognised BINs (that are missing species ID) if

578 the voucher specimens are available.

579 2. Analysis of the nominal species that are represented by more than one BIN, especially if

580 identifications represented by different BINs were produced by separate working teams.

581 3. Prioritised barcoding of representatives from families that are known to be important and

582 abundant components of communities; Phoxocephalidae, Ampeliscidae, Oedicerotidae, and

583 Caprellidae should be prioritised.

584 4. Targeted sampling programs for taxa coming from geographical regions with the least

585 knowledge.

586 5. Targeted sampling to obtain more sequences for taxa present in BOLD but represented by

587 small numbers of sequences (especially singletons), from different parts of the species' range if

588 possible.

589 6. Targeted programs to sequence type specimens stored in musea or to collect and study fresh

590 individuals from type localities if types are unsuitable for analyses.

591

592 **References**

593 **Arfianti T, Wilson S, Costello MJ. 2018.** Progress in the discovery of amphipod crustaceans.

594    *PeerJ* **6**:e5187. https://doi.org/10.7717/peerj.5187

595    **Aylagas E, Borja A, Muxika I, Rodríguez-Ezpeleta N. 2018.** Adapting metabarcoding based

596    benthic biomonitoring into routine marine ecological status assessment networks. Ecological

597    Indicators **95**:194–202. https://dx.doi.org/10.1016/j.ecolind.2018.07.044

598    **Barnard JL, Karaman GS. 1991.** The families and genera of marine gammaridean Amphipoda

599    (except marine gammaroids). Part 1 & 2. *Records of the Australian Museum Supplement*

600    **13**:1–866.

601    **Bauzà-Ribot MM, Jaume D, Fornós JJ, Juan C, Pons J. 2011.** Islands beneath islands:

602    phylogeography of a groundwater amphipod crustacean in the Balearic archipelago. BMC

603    *Evolutionary Biology* **11(1)**:221. https://doi.org/10.1186/1471-2148-11-221

604    **Beermann J, Hall-Mullen AK, Havermans C, Coolen JW, Crooijmans RP, Dibbits B, Held**

605    **C, Desiderato A. 2020.** Ancient globetrotters—connectivity and putative native ranges of

606    two cosmopolitan biofouling amphipods. *PeerJ* **8**:e9613 https://doi.org/10.7717/peerj.9613

607    **Bij de Vaate A, Jazdzewski K, Ketelaars HAM, Gollasch S, Van der Velde G. 2002.**

608    Geographical patterns in range extension of Ponto-Caspian macroinvertebrate species in

609    Europe. *Canadian Journal of Fisheries and Aquatic Science* **59**:1159–1174.

610    https://doi.org/10.1139/f02-098

611    **Brandt A. 1993.** Composition, abundance, and diversity of peracarid crustaceans on a transect of

612    the Kolbeinsey Ridge, north of Iceland. *Polar Biology* **13**:565–576.

613    https://doi.org/10.1007/BF00236399

614    **Brandt A, Alalykina I, Brix S, Brenke N, Błażewicz M, Golovan OA, Johannsen N, Hrinko**

615    **AM, Jażdżewska AM, Jeskulke K, Kamenev GM, Lavrenteva AV, Malyutina MV,**

616    **Riehl T, Lins L. 2019.** Depth zonation of deep-sea macrofauna of the Northwest Pacific.

617    *Progress in Oceanography* https://doi.org/10.1016/j.pocean.2019.102131

618    **Brandt A, De Broyer C, De Mesel I, Ellingsen K, Gooday A, Hilbig B, Linse K, Thomson**

619    **MR, Tyler P. 2007.** The biodiversity of the deep Southern Ocean benthos. *Philosophical*

620    *Transactions of the Royal Society B Biological Sciences* **362**:39–66.

621    https://doi.org/10.1098/rstb.2006.1952

622    **Brante A, Guzmán-Rendón G, Barría EM, Guillemin ML, Vera-Escalona I, Hernández**

623    **CE. 2019.** Post-disturbance genetic changes: the impact of the 2010 mega-earthquake and

624    tsunami on Chilean sandy beach fauna. *Scientific Reports* **9(1)**:1–12.

625    https://doi.org/10.1038/s41598-019-50525-1

626    **Brix S, Lörz A-N, Jażdżewska A, Hughes L, Tandberg AH, Pabis K, Stransky B, Krapp-**

627    **Schickel T, Sorbe J-C, Hendrycks E, Vader WJM, Frutos I, Horton T, Jażdżewski K,**

628    **Peart R, Beermann J, Coleman CO, Buhl-Mortensen L, Corbari L, Havermans C, Tato**

629    **R, Jimenez Campean A. 2018.** Amphipod family distributions around Iceland. *Zookeys*

630    **731**:41–53. https://doi.org/10.3897/zookeys.731.19854

631    **Brix S, Osborn KJ, Kaiser S, Truskey SB, Schnurr SM, Brenke N, Malyutina M, Martinez**

632    **Arbizu P. 2020.** Adult life strategy affects distribution patterns in abyssal isopods –

633    implications for conservation in Pacific nodule areas. *Biogeosciences* **17**: 6163–6184.

634    https://doi.org/10.5194/bg-17-6163-2020

635    **Bundschuh M, Zubrod JP, Klemm P, Elsaesser D, Stang C, Schulz R. 2013.** Effects of peak

636    exposure scenarios on *Gammarus fossarum* using field relevant pesticide mixtures.

637    *Ecotoxicological Environment Safety* **95**:137–143

638    https://doi.org/10.1016/j.ecoenv.2013.05.025

639    **Bush A, Compson ZG, Monk WA, Porter TM, Steeves R, Emilson EJ, Gagne N,**

640    **Hajibabaei M, Roy M, Baird DJ. 2019.** Studying ecosystems with DNA metabarcoding:

641    lessons from biomonitoring of aquatic macroinvertebrates. *Frontiers in Ecology and*

642    *Evolution* **7** (434) http://dx.doi.org/10.3389/fevo.2019.00434

643    **Cabezas MP, Xavier R, Branco M, Santos AM, Guerra-Garcia JM. 2014.** Invasion history of

644    *Caprella scaura* Templeton, 1836 (Amphipoda: Caprellidae) in the Iberian Peninsula: multiple

645    introductions revealed by mitochondrial sequence data. *Biological Invasions* **16**:2221–2245.

646    https://doi.org/10.1007/s10530-014-0660-y

647    **Coleman CO. 2015.** Taxonomy in Times of the Taxonomic Impediment – Examples from the

648    Community of Experts on Amphipod Crustaceans. *Journal of Crustacean Biology* **35(6)**:729–

649    740. https://doi.org/10.1163/1937240X-00002381

650    **Coleman CO, Radulovici AE. 2020.** Challenges for the future of taxonomy: talents, databases

651    and knowledge growth. *Megataxa* **001(1)**:028–034. https://doi.org/10.11646/megataxa.1.1.5

652    **Copilaş-Ciocianu D, Sidorov D, Gontcharov A. 2019.** Adrift across tectonic plates: molecular

653    phylogenetics supports the ancient Laurasian origin of old limnic crangonyctid amphipods.

654    *Organisms, Diversity and Evolution* **19**:191–207. https://doi.org/10.1007/s13127-019-00401-

655    7

656    **Corkum LD. 1989.** Patterns of benthic invertebrate assemblages in rivers of northwestern North

657    America. *Freshwater Biology* **21**:191–205. https://doi.org/10.1111/j.1365-

658    2427.1989.tb01358.x

659    **Cowart DA, Pinheiro M, Mouchel O, Maguer M, Grall J, Miné J, Arnaud-Haond S. 2015.**

660    Metabarcoding Is Powerful yet Still Blind: A Comparative Analysis of Morphological and

661    Molecular Surveys of Seagrass Communities. *PLoS ONE* **10(2)**:e0117562.

662    https://doi.org/10.1371/journal.pone.0117562

663    **Cristescu ME. 2014.** From barcoding single individuals to metabarcoding biological

664    communities: Towards an integrative approach to the study of global biodiversity. *Trends in*

665    *Ecology and Evolution* **29(10)**:566–571. http://dx.doi.org/10.1016/j.tree.2014.08.001

666    **Dalpadado P, Borkner N, Bogstad B, Mehl S. 2001.** Distribution of *Themisto* (Amphipoda)

667    spp. in the Barents Sea and predator-prey interactions. *ICES Journal of Marine Sciences*

668    **58**:876–895. https://doi.org/10.1006/jmsc.2001.1078

669    **Dauby P, Nyssen F, De Broyer C. 2003.** Amphipods as food sources for higher trophic levels in

670    the Southern Ocean: a synthesis. In: Huiskes AHL, Gieskes WWC, Rozema J, Schorno

671    RML, van der Vies SM, Wolff WJ, editors. Antarctic Biology in a Global Context. Backhuys

672    Publishers, Leiden, the Netherlands, 129–134.

673    **Dauby P, Scailteur Y, De Broyer C. 2001.** Trophic diversity within the eastern Weddell Sea

674    community. *Hydrobiologia* **443**:69–86. https://doi.org/10.1023/A:1017596120422

675    **Delić T, Švara V, Coleman CO, Trontelj P, Fišer C. 2017a.** The giant cryptic amphipod

676    species of the subterranean genus *Niphargus* (Crustacea, Amphipoda). *Zoologica Scripta*

677    **46(6)**:740–752. https://doi.org/10.1111/zsc.12252

678    **Delić T, Trontelj P, Rendoš M, Fišer C. 2017b.** The importance of naming cryptic species and

679    the conservation of endemic subterranean amphipods. *Scientific Reports* **7(1)**:1–12.

680    https://doi.org/10.1038/s41598-017-02938-z

681    **Desiderato A, Costa FO, Serejo CS, Abbiati M, Queiroga H, Vieira PE. 2019.** Macaronesian

682    islands as promoters of diversification in amphipods: The remarkable case of the family

683    Hyalidae (Crustacea, Amphipoda). *Zoologica Scripta* **48**:359– 375.

684    https://doi.org/10.1111/zsc.12339

685    **Duarte S, Vieira PE, Costa FO. 2020.** Assessment of species gaps in DNA barcode libraries of

686        non-indigenous species (NIS) occurring in European coastal regions. *Metabarcoding and*

687        *Metagenomics* **4**:35–46. https://doi.org/10.3897/mbmg.4.55162

688    **Ebach MC, Valdecasas AG, Wheeler QD. 2011.** Impediments to taxonomy and users of

689        taxonomy: accessibility and impact evaluation. *Cladistics* **27(5)**:550–557.

690        https://doi.org/10.1111/j.1096-0031.2011.00348.x

691    **Feio MJ, Filipe AF, Garcia-Raventos A, Ardura A, Calapez AR, Pujante AM, Mortágua A,**

692        **Múrria C, Diaz-de-Quijano D, Martins FMS, Duarte S, Sáinz Bariáin M, Cordeiro R,**

693        **Rivera SF, Väisänen LOS, Fonseca A, Gonçalves V, Garcia-Vazquez E, Vieites**

694        **Rodríguez D, Ivanova EA, Costa FO, Barquín J, Rojo V, Vierna J, Fais M, Suarez M,**

695        **Nieminen M, Hammers-Wirtz M, Kolmakova OV, Trusova MY, Beja P, González R,**

696        **Planes S, Almeida SFP. 2020.** Advances in the use of molecular tools in ecological and

697        biodiversity assessment of aquatic ecosystems. *Limnetica* **39(1)**:419–440.

698        https://doi.org/10.23818/limn.39.27

699    **Fišer C, Konec M, Alther R., Švara V, Altermatt F. 2017.** Taxonomic, phylogenetic and

700        ecological diversity of *Niphargus* (Amphipoda: Crustacea) in the Hölloch cave system

701        (Switzerland). *Systematics and Biodiversity* **15(3)**: 218-237

702        https://doi.org/10.1080/14772000.2016.1249112

703    **Fontes JT, Vieira PE, Ekrem T, Soares P, Costa FO. 2020.** BAGS: an automated Barcode,

704        Audit & Grade System for DNA barcode reference libraries. *Molecular Ecology Resources*.

705        https://doi.org/10.1111/1755-0998.13262

706    **Frutos I, Brandt A, Sorbe JC. 2017.** Deep-Sea Suprabenthic Communities: The Forgotten

707        Biodiversity. In: Rossi S, Bramanti L, Gori A, Orejas C, editors. Marine Animal Forests.

708        Springer International Publishing, Cham, p.475–503. https://doi.org/10.1007/978-3-319-

709        21012-4_21

710    **Frutos I, Sorbe JC. 2017.** Suprabenthic assemblages from the Capbreton area (SE Bay of Biscay):

711        faunal recovery after a canyon turbiditic disturbance. *Deep-Sea Research I* **130**:36–46.

712        https://doi.org/10.1016/j.dsr.2017.10.007

713    **Grabowski M, Wysocka A, Mamos T. 2017.** Molecular species delimitation methods provide

714        new insight into taxonomy of the endemic gammarid species flock from the ancient Lake

715        Ohrid. *Zoological Journal of the Linnean Society-London* **181(2)**:272–285.

716        https://doi.org/10.1093/zoolinnean/zlw025

717    **Hajibabaei M, Spall JL, Shokralla S, van Konynenburg S. 2012.** Assessing biodiversity of a

718        freshwater benthic macroinvertebrate community through non-destructive environmental

719        barcoding of DNA from preservative ethanol. *BMC Ecology* **12**:28

720        http://dx.doi.org/10.1186/1472-6785-12-28

721    **Havermans C. 2016.** Have we so far only seen the tip of the iceberg? Exploring species

722        diversity and distribution of the giant amphipod *Eurythenes*. *Biodiversity* **17**:12–25

723        https://doi.org/10.1080/14888386.2016.1172257

724    **Havermans C, Hagen W, Zeidler W, Held C, Auel H. 2019.** A survival pack for escaping

725        predation in the open ocean: amphipod–pteropod associations in the Southern Ocean. *Marine*

726        *Biodiversity* **49(3)**:1361–1370. https://doi.org/10.1007/s12526-018-0916-3

727  **Havermans C, Smetacek V. 2018.** Bottom-up and top-down triggers of diversification: A new

728    look at the evolutionary ecology of scavenging amphipods in the deep sea. *Progress in*

729    *Oceanography* **164**:37–51. https://doi.org/10.1016/j.pocean.2018.04.008

730  **Hestetun JT, Bye-Ingebrigtsen E, Nilsson RH, Glover AG, Johansen PO, Dahlgren TG.**

731    **2020.** Significant taxon sampling gaps in DNA databases limit the operational use of marine

732    macrofauna metabarcoding. *Marine Biodiversity* **50(5)**:1–9. https://doi.org/10.1007/s12526-

733    020-01093-5

734  **Horton T, Lowry J, De Broyer C, Bellan-Santini D, Coleman CO, Corbari L, Costello M J,**

735    **Daneliya M, Dauvin J-C, Fišer C, Gasca R, Grabowski M, Guerra-García JM,**

736    **Hendrycks E, Hughes L, Jaume D, Jazdzewski K, Kim Y-H, King R, Krapp-Schickel T,**

737    **LeCroy S, Lörz A-N, Mamos T, Senna AR, Serejo C, Sket B, Souza-Filho JF, Tandberg**

738    **AH, Thomas J, Thurston M, Vader W, Väinölä R, Vonk R, White K, Zeidler W. 2020.**

739    World Amphipoda Database. Accessed at http://www.marinespecies.org/amphipoda on 17-

740    07-2020.

741  **Horton, T. Marsh, L., Bett, B.J., Gates, A.R., Jones, D.O.B., Benoist, N.M.A., Pfeifer, S.**

742    **Simon-Lledó, E. Durden, J., Vandepitte, L., Appeltans, W. 2021.** Recommendations for

743    the standardisation of open taxonomic nomenclature for image-based identifications.

744    *Frontiers in Marine Sciences, Deep-Sea Environments and Ecology*

745    https://doi.org/10.3389/fmars.2021.620702

746  **Hou Z, Fu J, Li S. 2007.** A molecular phylogeny of the genus *Gammarus* (Crustacea:

747    Amphipoda) based on mitochondrial and nuclear gene sequences. *Molecular Phylogenetics*

748    *and Evolution* **45**:596–611. https://doi.org/10.1016/j.ympev.2007.06.006

749 **Humphries P, Davies PE Mulcahy ME. 1996.** Macroinvertebrate assemblages of littoral

750　　　 habitats in the Macquarie and Mersey rivers, Tasmania: Implications for the management of

751　　　 regulated rivers. *Regulated Rivers: Research & Management* **12**:99–122.

752　　　 doi:10.1002/(SICI)1099-1646(199601)12:1<99::AID-RRR382>3.0.CO;2-1

753 **Hyne RV, Everett DA. 1998.** Application of a benthic euryhaline amphipod, *Corophium* sp., as

754　　　 a sediment toxicity testing organism for both freshwater and estuarine systems. *Archives of*

755　　　 *Environmental Contamination and Toxicology* **34(1)**:26–33

756　　　 https://doi.org/10.1007/s002449900282

757 **IPBES. 2019.** Summary for policymakers of the global assessment report on biodiversity and

758　　　 ecosystem services of the Intergovernmental Science-Policy Platform on Biodiversity and

759　　　 Ecosystem Services. Díaz S, Settele J, Brondízio ES, Ngo HT, Guèze M, Agard J, Arneth A,

760　　　 Balvanera P, Brauman KA, Butchart SHM, Chan KMA, Garibaldi LA, Ichii K, Liu J,

761　　　 Subramanian SM, Midgley GF, Miloslavich P, Molnár Z, Obura D, Pfaff A, Polasky S,

762　　　 Purvis A, Razzaque J, Reyers B, Roy Chowdhury R, Shin YJ, Visseren-Hamakers IJ, Willis

763　　　 KJ, Zayas CN, editors. IPBES secretariat, Bonn, Germany; 2019.

764 **Jażdżewska A. 2015.** Kuril–Kamchatka deep sea revisited – insights into the amphipod abyssal

765　　　 fauna. *Deep-Sea Research II* **111**:294–300. https://doi.org/10.1016/j.dsr2.2014.08.008

766 **Jażdżewska AM, Corbari L, Driskell A, Frutos I, Havermans C, Hendrycks E, Hughes L,**

767　　　 **Lörz A-N, Stransky B, Tandberg AHS, Vader W, Brix S. 2018.** A genetic fingerprint of

768　　　 Amphipoda from Icelandic waters – the baseline for further biodiversity and biogeography

769　　　 studies. *Zookeys* **731**:55–73. https://doi.org/10.3897/zookeys.731.19931

770 **Jażdżewska AM, Mamos T. 2019.** High species richness of Northwest Pacific deep-sea

771      amphipods revealed through DNA barcoding. *Progress in Oceanography*

772      https://doi.org/10.1016/j.pocean.2019.102184

773 **Jażdżewska AM, Rewicz T, Mamos T, Wattier R, Bącela-Spychalska K, Grabowski M.**

774      **2020.** Cryptic diversity and mtDNA phylogeography of the invasive demon shrimp,

775      *Dikerogammarus haemobaphes* (Eichwald, 1841), in Europe. *NeoBiota* **57**:53–86.

776      https://doi.org/10.3897/neobiota.57.46699

777 **Jazdzewski K, De Broyer C, Pudlarz M, Zielinski D. 2001.** Seasonal fluctuations of vagile

778      benthos in the uppermost sublittoral of a maritime Antarctic fjord. *Polar Biology* **24**:910–917.

779      https://doi.org/10.1007/s003000100299

780 **Kaiser S, Barnes DK, Brandt A. 2007.** Slope and deep-sea abundance across scales: Southern

781      Ocean isopods show how complex the deep sea can be. *Deep-Sea Research II* **54**:1776–1789.

782      https://doi.org/10.1016/j.dsr2.2007.07.006

783 **Kelly DW, Muirhead JR, Heath DD, MacIsaac HJ. 2006.** Contrasting patterns in genetic

784      diversity following multiple invasions of fresh and brackish waters. *Molecular Ecology*

785      **15**:3641–3655. https://doi.org/10.1111/j.1365-294X.2006.03012.x

786 **Knox MA, Hogg ID, Pilditch CA. 2011.** The role of vicariance and dispersal on New Zealand's

787      estuarine biodiversity: the case of *Paracorophium* (Crustacea: Amphipoda). *Biological*

788      *Journal of the Linnean Society* **103(4)**:863–874. https://doi.org/10.1111/j.1095-

789      8312.2011.01675.x

790 **Knox MA, Hogg ID, Pilditch CA, Lörz AN, Hebert PDN, Steinke D. 2012.** Mitochondrial

791      DNA (COI) analyses reveal that amphipod diversity is associated with environmental

792    heterogeneity in deep-sea habitats. *Molecular Ecology* **21**:4885–4897

793    https://doi.org/10.1111/j.1365-294X.2012.05729.x

794    **Leese F, Bouchez A, Abarenkov K, Altermatt F, Borja A, Bruce K, Torbjřrn E, Čiampor**

795    **F, Čiamporová Zaťovičová Z, Costa FO, Duarte S, Elbrecht V, Fontaneto D, Franc A,**

796    **Geiger MF, Hering D, Kahlert M, Stroil BK, Kelly M, Keskin E, Liska I, Mergen P,**

797    **Meissner K, Pawlowski J, Penev L, Reyjol Y, Rotter A, Steinke D, van der Wal B,**

798    **Vitecek S, Zimmermann J, Weigand AM. 2018.** Why we need sustainable networks

799    bridging countries, disciplines, cultures and generations for aquatic biomonitoring 2.0: a

800    perspective derived from the DNAqua-Net COST action. *Advances in Ecological Resources*

801    **58**:63–99 http://dx.doi.org/10.1016/bs.aecr.2018.01.001

802    **Leite BR, Vieira PE, Teixeira MAL, Lobo-Arteaga J, Hollatz C, Borges LMS, Duarte S,**

803    **Troncoso JS, Costa FO. 2020.** Gap-analysis and annotated reference library for supporting

804    macroinvertebrate metabarcoding in Atlantic Iberia. *Regional Studies in Marine Science*

805    101307. https://doi.org/10.1016/j.rsma.2020.101307

806    **Lörz A-N, Jażdżewska AM, Brandt A. 2018.** A new predator connecting the abyssal with the

807    hadal in the Kuril-Kamchatka Trench, NW Pacific. *PeerJ* **6**:e4887.

808    https://doi.org/10.7717/peerj.4887

809    **Major K, Soucek DJ, Giordano R, Wetzel MJ, Soto-Adames F. 2013.** The common

810    ecotoxicology laboratory strain of *Hyalella azteca* is genetically distinct from most wild

811    strains sampled in eastern North America. *Environmental toxicology and chemistry*

812    **32(11)**:2637–2647. https://doi.org/10.1002/etc.2355

813    **Mamos T, Wattier R, Burzyński A, Grabowski M. 2016.** The legacy of a vanished sea: a high

814    level of diversification within a European freshwater amphipod species complex driven by 15

815     My of Paratethys regression. *Molecular Ecology* **3**:795–810.

816     https://doi.org/10.1111/mec.13499

817     **Mamos T, Wattier R, Majda A, Sket B, Grabowski M. 2014.** Morphological vs. molecular

818     delineation of taxa across montane regions in Europe: the case study of *Gammarus*

819     *balcanicus* Schäferna, 1922 (Crustacea: Amphipoda). *Journal of Zoological Systematics and*

820     *Evolutionary Research* **52(3)**:237–248. https://doi.org/10.1111/jzs.12062

821     **Mora C, Tittensor DP, Adl S, Simpson AGB, Worm B. 2011.** How Many Species Are There

822     on Earth and in the Ocean? *PLoS Biology* **9(8)**:e1001127.

823     https://doi.org/10.1371/journal.pbio.1001127

824     **Múrria C, Väisänen LO, Somma S, Wangensteen OS, Arnedo MA, Prat N. 2020.** Towards

825     an Iberian DNA barcode reference library of freshwater macroinvertebrates and fishes.

826     *Limnetica* **39(1)**:73–92. http://dx.doi.org/10.23818/limn.39.06

827     **Oliveira LM, Knebelsberger T, Landi M, Soares P, Raupach MJ, Costa FO. 2016.**

828     Assembling and auditing a comprehensive DNA barcode reference library for European

829     marine fishes. *Journal of Fish Biology* **89(6)**:2741–2754. http://dx.doi.org/10.1111/jfb.13169

830     **Page MJ, McKenzie JE, Bossuyt PM, Boutron I, Hoffmann TC, Mulrow CD, Shamseer L,**

831     **Tetzlaff JM, Akl EA, Brennan SE, Chou R, Glanville J, Grimshaw JM, Hróbjartsson**

832     **A, Lalu MM, Li T, Loder EW, Mayo-Wilson E, McDonald S, McGuinness LA, Stewart**

833     **LA Thomas J, Tricco AC, Welch VA, Whiting P, Moher D. 2021.** The PRISMA 2020

834     statement: an updated guideline for reporting systematic reviews. *BMJ* **372**:n71.

835     http://doi.org/10.1136/bmj.n71

836     **Page RDM. 2016.** DNA barcoding and taxonomy: dark taxa and dark texts. *Philosophical*

837     *Transactions of the Royal Society B* **371**:20150334. https://doi.org/10.1098/rstb.2015.0334

838    **QGIS Development Team. 2018.** QGIS Geographic Information System. Open Source

839    Geospatial Foundation Project. http://qgis.osgeo.org.

840    **Radulovici AE, Coleman CO. 2017.** Reconciling large molecular datasets, bioinformatics and

841    taxonomy: prospects for Amphipoda. *Biodiversity Journal* **8(2)**:633–634.

842    **Radulovici AE, Vieira PE, Duarte S, Teixeira MAL, Borges LMS, Deagle B, Majaneva S,**

843    **Redmond N, Schultz JA, Costa FO. 2021.** Revision and annotation of DNA barcode

844    records for marine invertebrates: report of the 8th iBOL conference hackathon. bioRxiv

845    2021.03.07.434272. https://doi.org/10.1101/2021.03.07.434272

846    **Ratnasingham S, Hebert P. 2007.** BOLD: The Barcode of Life Data System. *Molecular*

847    *Ecology Notes* **7(3)**:355–364 https://doi.org/10.1111/j.1471-8286.2007.01678.x

848    **Ratnasingham S, Hebert PDN. 2013.** A DNA-based registry for all animal species: the barcode

849    index number (BIN) system. *PLoS ONE* **8**:e66213

850    https://doi.org/10.1371/journal.pone.0066213

851    **Rewicz T, Wattier R, Grabowski M, Rigaud T, Bącela-Spychalska K. 2015.** Out of the Black

852    Sea: Phylogeography of the Invasive Killer Shrimp *Dikerogammarus villosus* across Europe.

853    *PLoS ONE* **10(2)**:e0118121. https://doi.org/10.1371/journal.pone.0118121

854    **Riehl T, Brenke N, Brix S, Driskell A, Kaiser S, Brandt A. 2014.** Field and Laboratory

855    Methods for DNA Studies on Deep-sea Isopod Crustaceans. *Polish Polar Research* **35**:203–

856    224. https://doi.org/10.2478/popore-2014-0018

857    **Ros M, Guerra-García JM, González-Macías M, Saavedra Á, López-Fe CM. 2013.**

858    Influence of fouling communities on the establishment success of alien caprellids (Crustacea:

859    Amphipoda) in Southern Spain. *Marine Biology Research* **9(3)**:261–273.

860    https://doi.org/10.1080/17451000.2012.739695

861 **Ros M, Vázquez-Luis M, Guerra-García JM. 2013.** The role of marinas and recreational

862      boating in the occurrence and distribution of exotic caprellids (Crustacea: Amphipoda) in the

863      Western Mediterranean: Mallorca Island as a case study. *Journal of Sea Research* **83**:94–103.

864      https://doi.org/10.1016/j.seares.2013.04.004

865 **Sayers EW, Cavanaugh M, Clark K, Ostell J, Pruitt KD, Karsch-Mizrachi I. 2020.**

866      GenBank. *Nucleic Acids Research* **47**:D94–D93. https://doi.org/10.1093/nar/gkz956.

867 **Sedano F, Navarro-Barranco C, Guerra-García JM, Espinosa F. 2020.** From sessile to

868      vagile: Understanding the importance of epifauna to assess the environmental impacts of

869      coastal defence structures. *Estuarine, Coastal and Shelf Sciences* **235**:106616.

870      https://doi.org/10.1016/j.ecss.2020.106616

871 **Sigovini M., Keppel E., Tagliapietra D. 2016.** Open nomenclature in the biodiversity era.

872      *Methods in Ecology and Evolution* **7**:10. https://doi.org/10.1111/2041-210X.12594

873 **Tempestini A, Rysgaard S, Dufresne F. 2018.** Species identification and connectivity of

874      marine amphipods in Canada's three oceans. *PLoS ONE* **13(5)**: e0197174.

875      https://doi.org/10.1371/journal.pone.0197174

876 **Tomikawa K, Kyono M, Kuribayashi K, Nakano T. 2017.** The enigmatic groundwater

877      amphipod *Awacaris kawasawai* revisited: synonymisation of the genus *Sternomoera*, with

878      molecular phylogenetic analyses of *Awacaris* and *Sternomoera* species (Crustacea:

879      Amphipoda: Pontogeneiidae). *Invertebrate Systematics* **31(2)**:125–140.

880      https://doi.org/10.1071/IS16037

881 **Väinölä R, Witt JDS, Grabowski M, Bradbury JH, Jażdżewski K, Sket B. 2008.** Global

882      diversity of amphipods (Amphipoda; Crustacea) in freshwater. *Hydrobiologia* **595**:241–255.

883      https://doi.org/10.1007/s10750-007-9020-6

884 **Verheye ML, Backeljau T, d'Udekem d'Acoz C. 2016.** Looking beneath the tip of the iceberg:

885 diversification of the genus *Epimeria* on the Antarctic shelf (Crustacea, Amphipoda). *Polar*

886 *Biology* **39(5)**:925–945. https://doi.org/ 10.1007/s00300-016-1910-5

887 **Vieira PE, Lavrador AS, Parente MI, Parretti P, Costa AC, Costa FO, Duarte S. 2021**.

888 Gaps in DNA sequence libraries for Macaronesian marine macroinvertebrates imply decades

889 till completion and robust monitoring. *Diversity and Distributions* **00**: 1–13.

890 https://doi.org/10.1111/ddi.13305

891 **Vinogradov ME, Volkov AF, Semenova TN. 1996.** Hyperiid Amphipods (Amphipoda,

892 Hyperiidea) of the World Oceans. Smithsonian Institution Libraries, Washington, DC.

893 **Wattier, R., Mamos, T., Copilaş-Ciocianu, D., Jelić M., Ollivier A., Chaumot A., Danger M,**

894 **Felten V., Piscart C., Žganec K., Rewicz T., Wysocka A., Rigaud T., Grabowski M.**

895 **2020.** Continental-scale patterns of hyper-cryptic diversity within the freshwater model taxon

896 *Gammarus fossarum* (Crustacea, Amphipoda). *Scientific Reports* **10**:16536.

897 https://doi.org/10.1038/s41598-020-73739-0

898 **Weigand H, Beermann AJ, Čiampor F, Costa FO, Csabai Z, Duarte S, Geiger MF,**

899 **Grabowski M, Rimet F, Rulik B, Strand M, Szucsich N, Weigand AM, Willassen E,**

900 **Wyler SA, Bouchez A, Borja A, Čiamporová Zaťovičová Z, Ferreira S, Dijkstra KB,**

901 **Eisendle U, Freyhof J, Gadawski P, Graf W, Haegerbaeumer A, van der Hoorn BB,**

902 **Japoshvili B, Keresztes L, Keskin E, Lesse F, Macher JN, Mamos T, Paz G, Pešić V,**

903 **Pfannkuchen DM, Pfannkuchen MA, Price BW, Rinkevich B, Teixeira MAL, Várbíró**

904 **G, Ekrem T. 2019.** DNA barcode reference libraries for the monitoring of aquatic biota in

905 Europe: Gap-analysis and recommendations for future work. *Science of the Total*

906 *Environment* **678**:499–524. http://dx.doi.org/10.1016/j.scitotenv.2019.04.247

907   **Weisshappel JBF, Svavarsson J. 1998.** Benthic amphipods (Crustacea: Malacostraca) in

908       Icelandic waters: diversity in relation to faunal patterns from shallow to intermediate deep

909       Arctic and North Atlantic Oceans. *Marine Biology* **131**:133–142.

910   **Witt JD, Threloff DL, Hebert PD. 2006.** DNA barcoding reveals extraordinary cryptic

911       diversity in an amphipod genus: implications for desert spring conservation. *Molecular*

912       *Ecology* **15(10)**:3073–3082. https://doi.org/10.1111/j.1365-294X.2006.02999.x

913   **Wysocka A, Grabowski M, Sworobowicz L, Burzyński A, Kilikowska A, Kostoski G, Sell J.**

914       **2013.** A tale of time and depth: intralacustrine radiation in endemic *Gammarus* species flock

915       from the ancient Lake Ohrid. *Zoological Journal of the Linnean Society-London* **167(3)**:345–

916       359. https://doi.org/10.1111/j.1096-3642.2012.00878.x

917   **Wysocka A, Grabowski M, Sworobowicz L, Mamos T, Burzyński A, Sell J. 2014.** Origin of

918       the Lake Ohrid gammarid species flock: ancient local phylogenetic lineage diversification.

919       *Journal of Biogeography* **41(9)**:1758–1768 https://doi.org/10.1111/jbi.12335

920   **Figure captions**

921   Figure 1. PRISMA 2020 work-flow diagram (Page et al., 2021). Summary of the data download,

922       identification and screening before analysis. All record removals were done by the leading

923       author of the paper.

924   Figure 2. Environmental origin of the amphipod records (A) and BINs (B) in BOLD database.

925   Figure 3. Geographic distribution of amphipod records expressed by sequences present in BOLD

926       (A – freshwater, B – marine, C – terrestrial). Dots indicate records with exact coordinates, for

927       records without latitude and longitude the country of origin was checked. Background color

928       of the country indicates this number per country.

929     Figure 4. Proportion of records (A) and BINs (B) with different level of identification within

930       freshwater, marine and terrestrial amphipod taxa.

931     Figure 5. Number of BINs represented by given number of sequences. Upper set (A, B, C) – all

932       BINs, lower set (D, E, F) – only BINs with complete species-level identification considered.

933       A, D – freshwater, B, E – marine, C, F – terrestrial taxa.

934     Figure 6. Number of nominal species represented by given number of BINs.

935

**Supplemental files list**

936     **Supplemental files list**

937

938     **Supplemental file 1**

939     - File format .xlsx

940     - Title: List of the doubled records indentified in the original dataset

941     - Description: The file consists of records that appeared to be doubled in the BOLD database.

942

943     **Supplemental file 2**

944     - File format .xlsx

945     - Title: The list of records analysed (after removal of doubled records).

946     - Description: The file consists of all records that was basis of the present study. The colors of

947       the cell of recordID indicates the environment: green - marine (including brakishwater and

948       fully marine taxa), red - freshwater (including freshwater and brakishwater taxa), yellow -

949       terrestrial, blue - environment not recorded.

950

951     **Supplemental file 3**

952     - File format .xlsx

953     - Title: List of BINs possessing more than one ID variant with notes on the identification

954     - Description: The file presents BINs that have received different identifications with details of

955       the ID and comments concerning the final identification used in the study.

956

957 **Supplemental file 4**

958 - File format .xlsx

959 - Title: List of BINs for which the environment was not able to be assessed.

960 - Description: The file presents BINs for which the available data did not allow to specify the

961    environment.

962

963 **Supplemental file 5**

964 - File format .xlsx

965 - Title: List of BINs with the largest number of records (blue - freshwater, green - marine, yellow

966    - terrestrial taxa).

967 - Description: The file presents BINs that are represented by the largest number of records. The

968    information about environmental origin of associated species are provided.

969

970 **Supplemental file 6**

971 - File format .xlsx

972 - Title: Nominal species with the largest number of BINs identified.

973 - Description: The file presents nominal species for which the largest number of BINs has been

974    identified. Environmental origin of species is also provided.

975

976 **Supplemental file 7**

977 - File format .xlsx

978 - Title: List of amphipod families with number of accepted and barcoded species as well as

979    information of the barcoding coverage within family. Families within each category with the

980    highest barcoding coverage indicated in bold.

981

982 **Supplemental file 8**

983 - File format .docx

984 - PRISMA 2020 checklist

**Table 1**(on next page)

Representation of amphipod families in BOLD.

* in parentheses the number of families without barcoded species but with at least one BIN identified to the genus (g) or family (f) level.

1  Table 1. Representation of amphipod families in BOLD. * in parentheses the number of families

2  without barcoded species but with at least one BIN identified to the genus (g) or family (f) level.

| Number of families | |
|---|---|
| without any barcoded species | 117 (+ 13g, 5f)* |
| with up to 10% barcoded species | 47 |
| with 11-20% barcoded species | 24 |
| with 21-50% barcoded species | 24 |
| with >50% barcoded species | 10 |

3

4

# Table 2 (on next page)

Number of amphipod species in each realm with indication of their barcode quality according to grading system from Fontes et al. (2020).

A – consolidated concordance, B – basal concordance, C – multiple BINs for single

morphospecies, D – insufficient data; for more detailed explanation of grading system, see

Material and methods section.

1   Table 2. Number of amphipod species in each realm with indication of their barcode quality

2   according to grading system from Fontes et al. (2020). A – consolidated concordance, B – basal

3   concordance, C – multiple BINs for single morphospecies, D – insufficient data; for more

4   detailed explanation of grading system, see Material and methods section.

|                   | A   | B   | C   | D   | all species |
|-------------------|-----|-----|-----|-----|-------------|
| All species       | 100 | 155 | 276 | 468 | 999         |
| Freshwater spp.   | 31  | 55  | 140 | 225 | 451         |
| Marine spp.       | 58  | 92  | 120 | 226 | 496         |
| Terrestrial spp.  | 11  | 8   | 16  | 17  | 52          |

5

6

**Table 3**(on next page)

Number of accepted families and species of Amphipoda (according to WAD accessed on 17-07-2020), number of families with representation in BOLD, number of species present in BOLD and mean coverage of barcodes in amphipod families represented in BOLD.

1 Table 3. Number of accepted families and species of Amphipoda (according to WAD accessed

2 on 17-07-2020), number of families with representation in BOLD, number of species present in

3 BOLD and mean coverage of barcodes in amphipod families represented in BOLD.

| | No. of families | No. of species | No. of families with species representation in BOLD | No. of species present in BOLD | Mean barcode coverage [%] of those families with representation in BOLD |
|---|---|---|---|---|---|
| Very species rich families (>100 spp.) | 33 | 7302 | 32 | 714 | 8 |
| Species rich families (31-100 spp.) | 30 | 1633 | 22 | 127 | 10 |
| Moderately species rich families (11-30 spp.) | 53 | 979 | 26 | 107 | 21 |
| Low species rich families (<10 spp.) | 123 | 416 | 26 | 51 | 49 |

4

5

**Table 4**(on next page)

Percent of families with species belonging to different quality grading categories (Fontes et al., 2020).

A – consolidated concordance, B – basal concordance, C – multiple BINs for single morphospecies, D – insufficient data; for more detailed explanation of grading system, see Material and methods section.

1    Table 4. Percent of families with species belonging to different quality grading categories

2    (Fontes et al., 2020). A – consolidated concordance, B – basal concordance, C – multiple BINs

3    for single morphospecies, D – insufficient data; for more detailed explanation of grading system,

4    see Material and methods section.

| | % of families | | | | |
|---|---|---|---|---|---|
| | All families | Very species rich families (>100 spp.) | Species rich families (31-100 spp.) | Moderately species rich families (11-30 spp.) | Low species rich families (<10 spp.) |
| At least one sp. in the category A | 32.4 | 65.6 | 31.8 | 16 | 7.7 |
| At least one sp. in the category B | 31.4 | 28.1 | 36.4 | 28 | 34.6 |
| At least one sp. in the category C | 10.5 | 0 | 13.6 | 24 | 7.7 |
| At least one sp. in the category D | 25.7 | 6.3 | 18.2 | 32 | 50 |
| Number of families | 105 | 32 | 22 | 25 | 26 |

5

6

# Figure 1

PRISMA 2020 work-flow diagram (Page et al., 2021).

Summary of the data download, identification and screening before analysis. All record removals were done by the leading author of the paper.
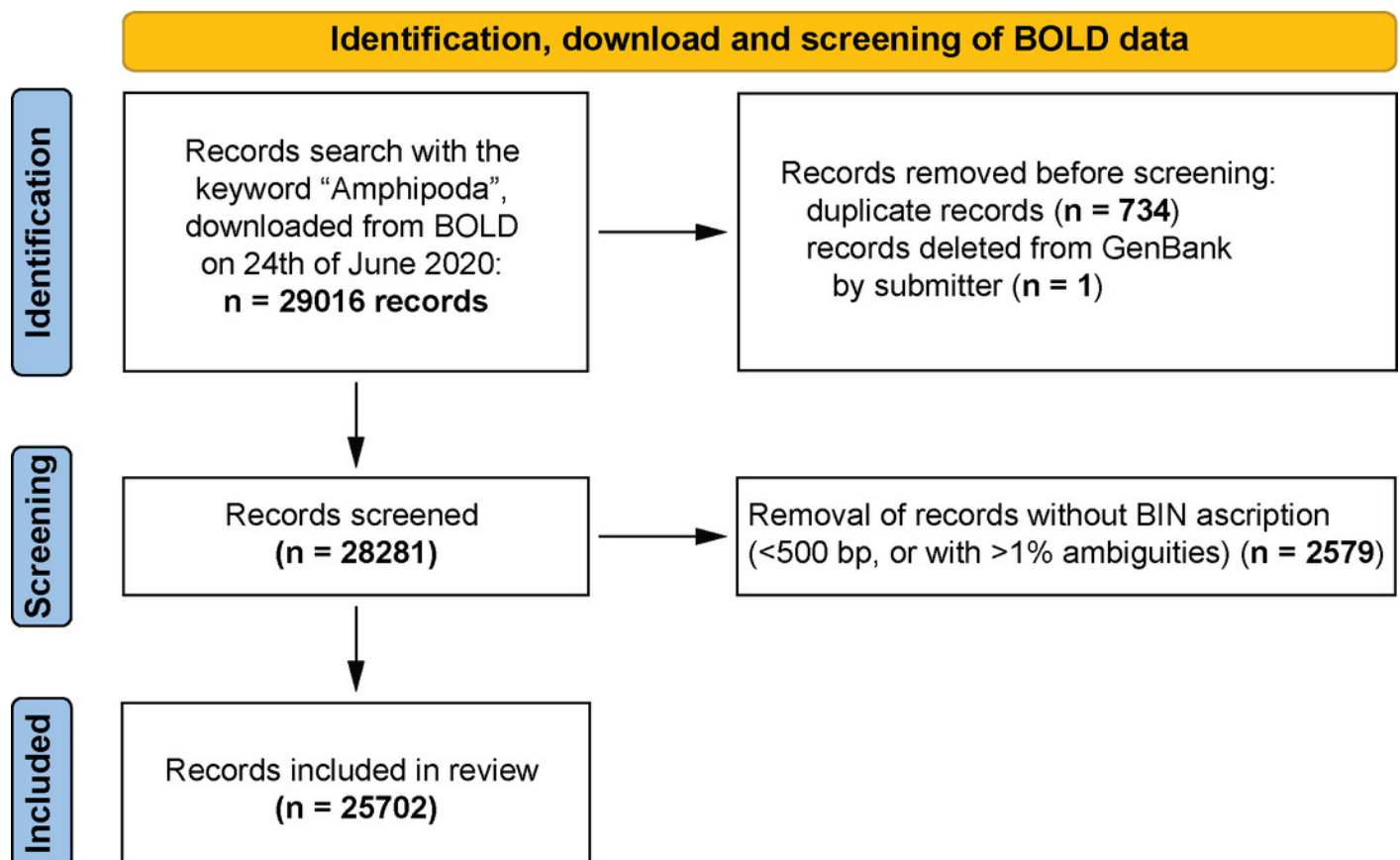
# Figure 2

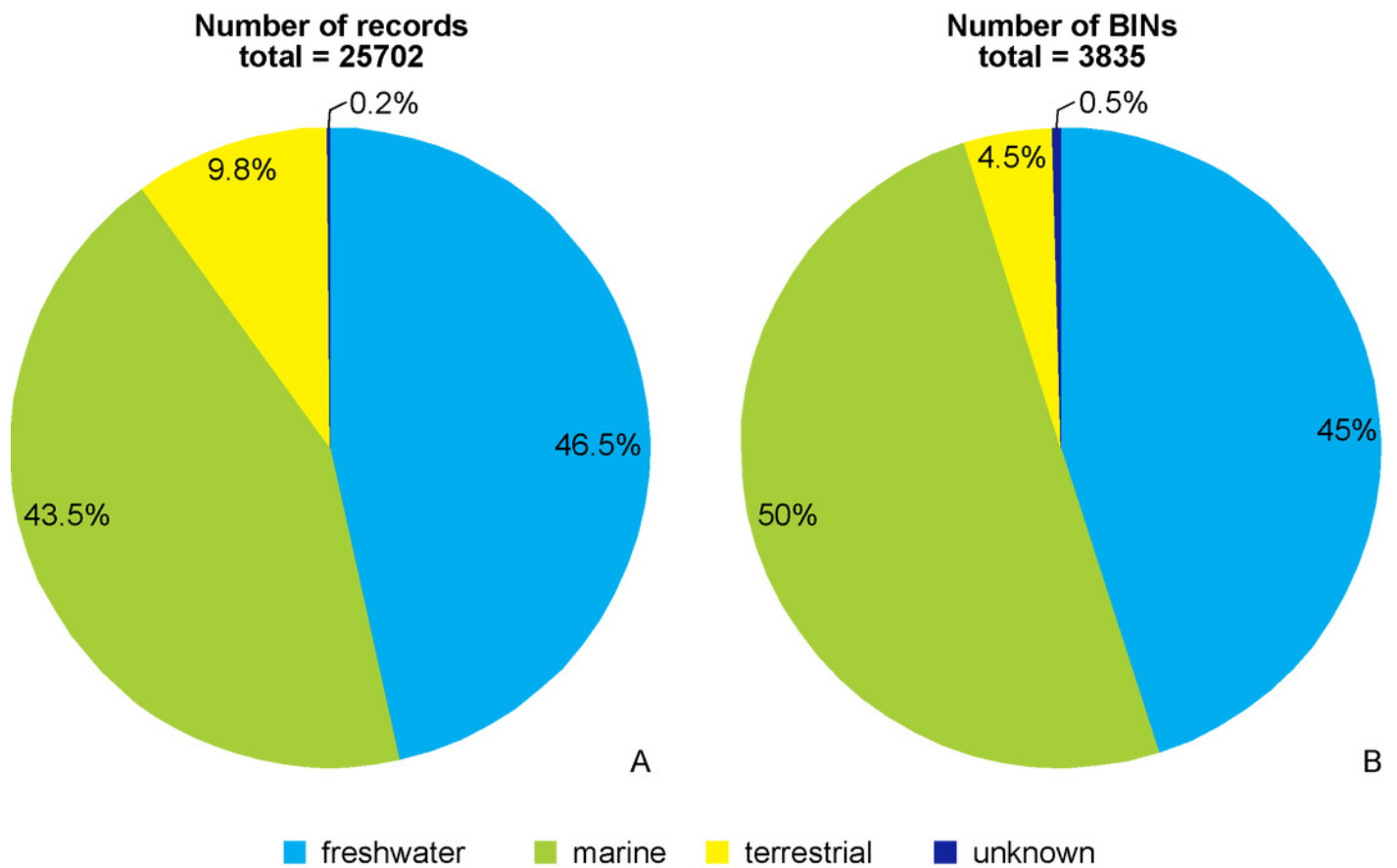Environmental origin of the amphipod records (A) and BINs (B) in BOLD database.



**Number of records**
**total = 25702**

0.2%
9.8%
46.5%
43.5%

A

**Number of BINs**
**total = 3835**

0.5%
4.5%
45%
50%

B

freshwater  marine  terrestrial  unknown

# Figure 3

Geographic distribution of amphipod records expressed by sequences present in BOLD (A – freshwater, B – marine, C – terrestrial).

Dots indicate records with exact coordinates, for records without latitude and longitude the country of origin was checked. Background color of the country indicates this number per country.
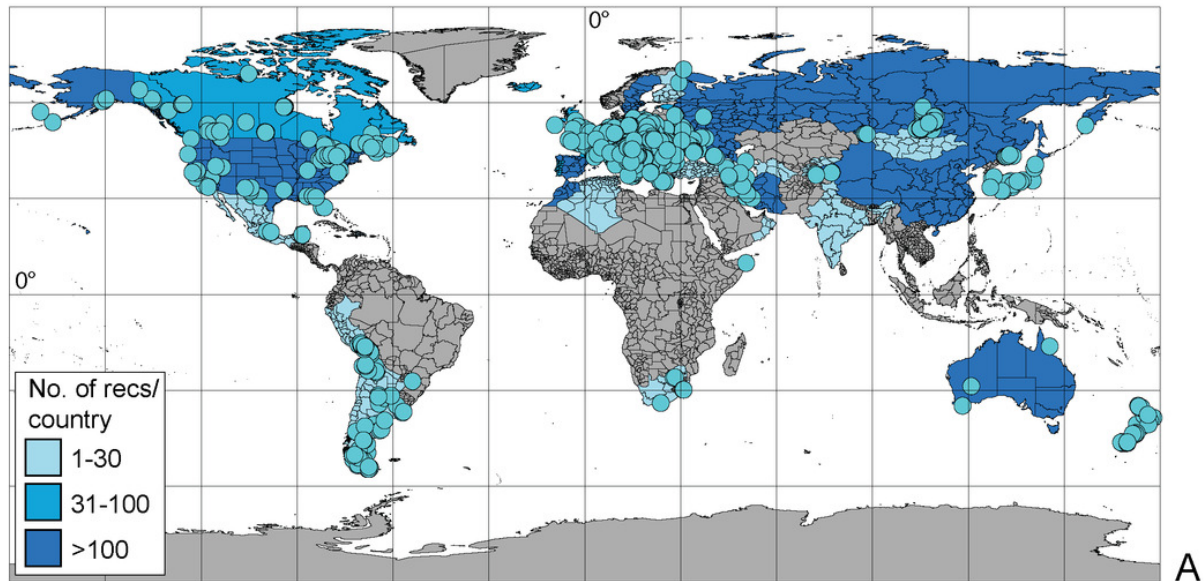
# Figure 4

Proportion of records (A) and BINs (B) with different level of identification within freshwater, marine and terrestrial amphipod taxa.
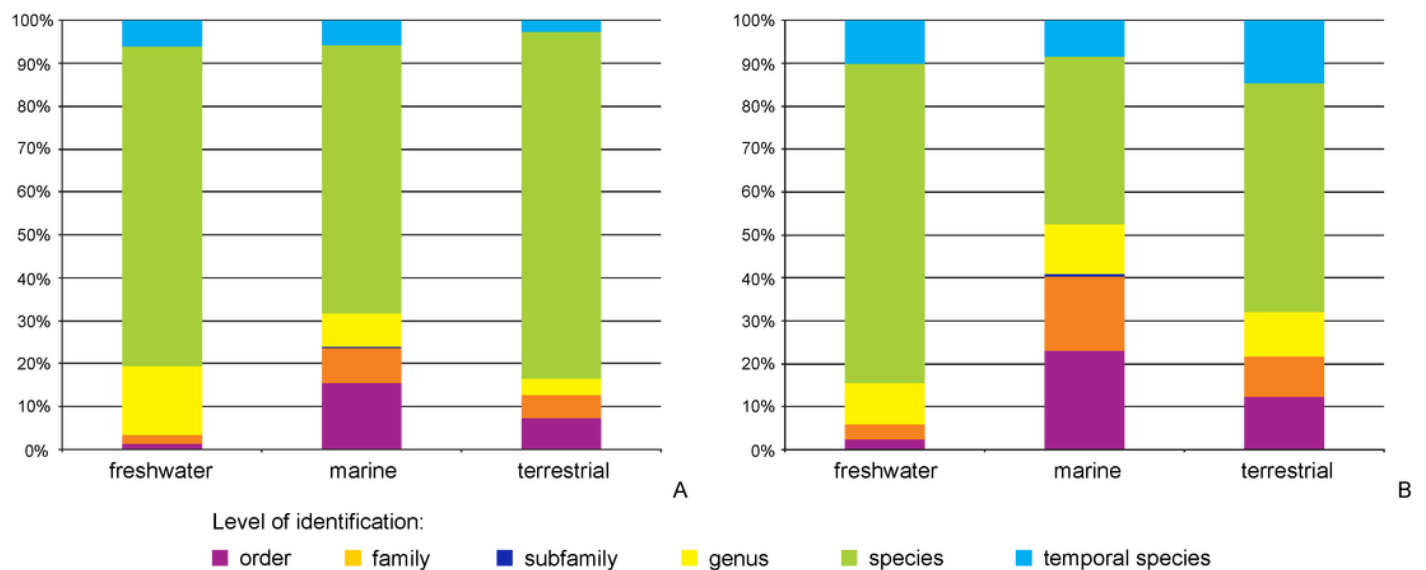
# Figure 5

Number of BINs represented by given number of sequences.

Upper set (A, B, C) – all BINs, lower set (D, E, F) – only BINs with complete species-level
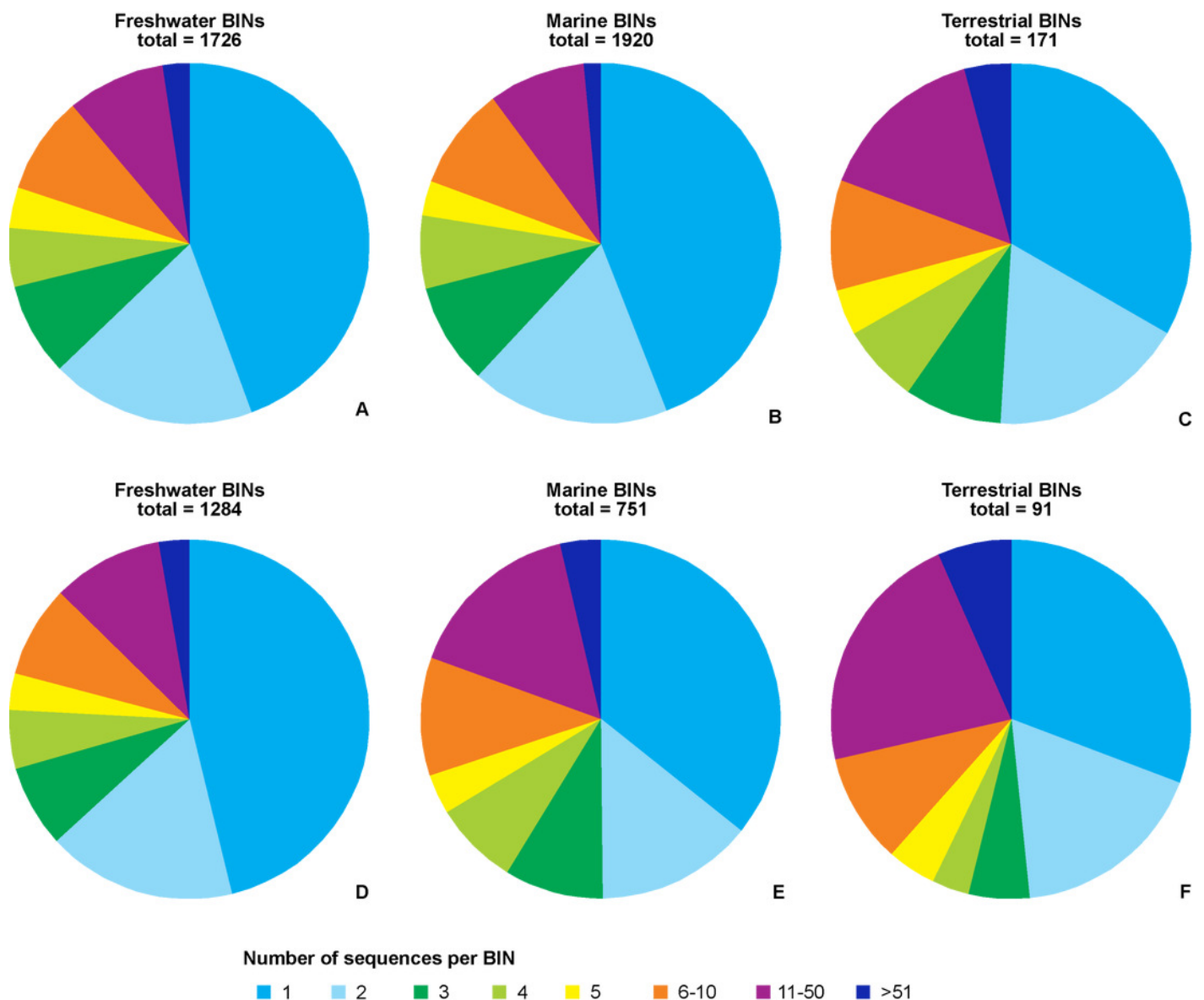identification considered. A, D – freshwater, B, E – marine, C, F – terrestrial taxa.

# Figure 6

Number of nominal species represented by given number of BINs.