**Article:**

eprints@whiterose.ac.uk
https://eprints.whiterose.ac.uk/

1    Linking dimensions of data on global marine animal diversity

2    Thomas J. Webb[1], Bart Vanhoorne[2]

3    1. Department of Animal and Plant Sciences, University of Sheffield, S10 2TN, UK,

4    t.j.webb@sheffield.ac.uk

5    2. Flanders Marine Institute (VLIZ), Ostend, Belgium

6    Abstract

7    Recent decades have seen an explosion in the amount of data available on all aspects of

8    biodiversity, which has led to data-driven approaches to understand how and why diversity

9    varies in time and space. Global repositories facilitate access to various classes of species-

10    level data including biogeography, genetics, and conservation status, which are in turn

11    required to study different dimensions of diversity. Ensuring that these different data sources

12    are interoperable is a challenge as we aim to create synthetic data products to monitor the

13    state of the world's biodiversity. One way to approach this is to link data of different classes,

14    and to inventory the availability of data across multiple sources. Here, we use a

15    comprehensive list of >200,000 marine animal species, and quantify the availability of data

16    on geographic occurrences, genetic sequences, conservation assessments, and DNA

17    barcodes across all phyla and broad functional groups. This reveals a very uneven picture:

18    44% of species are represented by no record other than their taxonomy, but some species

19    are rich in data. Although these data-rich species are concentrated into a few taxonomic and

20    functional groups, especially vertebrates, data is spread widely across marine animals, with

21    members of all 32 phyla represented in at least one database. By highlighting gaps in

22    current knowledge, our census of marine diversity data helps to prioritise future data

23    collection activities, as well as emphasising the importance of ongoing sustained

24    observations and archiving of existing data into global repositories.

26 The explosion in the availability of data describing the natural world has, in recent decades,

27 transformed the kinds of questions that we can now ask as ecologists. Efforts to reconstruct

28 the evolutionary relationships between all living species (e.g. Open Tree of Life; [1,2]) can

29 draw upon over 200M sequences (https://www.ncbi.nlm.nih.gov/genbank/statistics/) from

30 over 170,000 metazoan species stored in GenBank [3,4]. In 2018, the Global Biodiversity

31 Information Facility (GBIF, [5]) passed a billion species occurrence records

32 (https://www.gbif.org/news/5BesWzmwqQ4U84suqWyOQy/big-data-for-biodiversity-gbiforg-

33 surpasses-1-billion-species-occurrences), providing an unparalleled resource for students of

34 biogeography. The conservation status of >116,000 species has now been formally

35 assessed [6]. Significant efforts are underway to collate data biological, physiological,

36 metabolic and thermal traits [7-11] across multiple species, as well as information on animal

37 movement [12,13] and ecological interactions [14].

38

39 Against this background of increased data availability, the oceans are still often

40 characterised as the data-poor relative of the data-rich land. Various autonomous platforms

41 operating throughout the world's oceans do now enable vast quantities of physical and

42 biogeochemical data to be transmitted [15] but marine biodiversity data remain more

43 challenging to collect. In part, the vastness of the oceans precludes routine and casual

44 observation by the citizen scientists who have contributed so much to the collection of

45 terrestrial biodiversity data [16,17], except in some more accessible coastal areas [18-20].

46 However coordinated global initiatives have made enormous progress in collating existing

47 data and promoting systematic new data collection. The Census of Marine Life [21] drove

48 this effort from 2000-2010, and its legacies include the Ocean Biodiversity Information

49 System (OBIS, [22]), which currently holds nearly 60M occurrence records from over

50 120,000 marine species. Initiatives like this have built on sustained observations of marine

51 ecosystems [23], and continue to be developed to deliver the Essential Biodiversity

52    Variables that we need to monitor progress towards Sustainable Development Goals (e.g.

53    [24]). Application of technologies from satellites and drones to biologgers and molecular

54    methods such as eDNA continue to expand the range of data available to marine biodiversity

55    scientists [25]. Crucially, the accumulation of data has proceeded in parallel with massive

56    improvements in data infrastructure, and much better tools (taking advantage of the

57    improved computing power available even to casual users) with which to access and

58    analyse it [26,27]. This is important because the challenge now is to extract meaning from

59    the sea of data, to deliver effective outcomes for marine conservation and monitoring of the

60    state of the global ocean [19,24].

61

62    Although access to biodiversity data of different types is now much improved, to extract full

63    value from existing data requires linking together different datasets that were often collected

64    for different purposes, by different organisations and at different times. This kind of

65    interoperability of diversity data is central to the vision of a 'macroscope' to sample and

66    monitor the entire biosphere [25], and is a fundamental principle of the Bari Manifesto of best

67    practice in biodiversity informatics [28]. Progress towards such interoperability requires

68    comparable coverage across multiple classes of data and dimensions of diversity, as well as

69    parallel measures of the abiotic environment and of human pressures. An exemplar of

70    successful data integration for terrestrial plant communities is the Botanical Information and

71    Ecology Network [29] which combines standardised information on plant distributions, traits,

72    and evolutionary relationships with the computational tools needed to work with them. An

73    important step towards this kind of model is to fully understand the gaps and biases in

74    available data. In the marine environment, key gaps in the overall knowledge of marine

75    biodiversity have been documented [30-32], including estimates of the extent of unknown

76    biodiversity [33] and undocumented extinction risk [34]. Efforts to quantify these gaps across

77    different key variables and data sources have been limited to the regional scale, but have

78    shown for instance that the species and taxonomic groups that we know most in one

79    dimension (e.g. global occurrences) tend to be those that we also know most about in

80    another (e.g. biological traits, extinction risk; [34,35]). To date we lack a global overview of

81    how data (and gaps) are co-distributed across axes of marine diversity, to compare for

82    example with previous global analyses of terrestrial plants [36].

83

84    Such a task is feasible however, given the availability of a standardised global taxonomy of

85    marine species, the World Register of Marine Species (WoRMS, [37]), which includes links

86    out to other key biodiversity datasets (Table 1). In this paper, we focus on key data sources

87    which, when linked to robust taxonomy, individually or in combination can be used to

88    construct different dimensions of marine diversity. We consider geographic occurrences and

89    nucleotide sequences to be the fundamental building blocks of the spatial and phylogenetic

90    dimensions of diversity, which interact to structure the distribution of key ecological traits

91    across species [38]. A first step to adding the functional dimension of diversity is to classify

92    species into broad ecological guilds, similar to the way in which species can be classified in

93    global theories and models of biodiversity [39,40]. Supplementing these with information on

94    conservation status and molecular taxonomy provides insights into how marine diversity is

95    changing, and how we might efficiently monitor this. Throughout we use open source

96    computational tools to link data across these components of marine diversity to take stock of

97    the current state of data availability, identifying gaps and priorities for future work. In this way

98    we summarise data availability across multiple axes for >200,000 marine animal species

99    from 32 phyla and across broad ecological guilds (e.g. benthos, zooplankton, seabirds), and

100   we assess the extent to which this availability is correlated across different classes of

101   diversity data. Above all, our aim is to highlight the wealth of marine biodiversity data that we

102   have amassed as a community over centuries, and the opportunities that we now have to

103   link different classes of data in order to better understand the dimensions of marine diversity.

104   Methods

105   To provide an overview of the state of knowledge of marine animal biodiversity, we mine the

106   World Register of Marine Species (WoRMS, [37]), the most comprehensive source of

107    taxonomic information on marine species, consisting of over half a million distinct names

108    checked by expert taxonomic editors. We focus our investigation on marine animals, and so

109    filtered the WoRMS database to Kingdom Animalia, retaining only those species considered

110    to be marine by WoRMS (flag `isMarine` is TRUE), and excluding any species only known

111    from fossils. We consider only taxa identified at the species rank, with a current accepted

112    name and valid WoRMS identifier (Aphia ID).

113

114    In addition to taxonomy, WoRMS has aggregated data on species attributes including broad

115    'functional groups'. In reality these are closer to ecological guilds, defining habitat affinity

116    (e.g. benthos, zooplankton) rather than ecological function, but for transparency we retain

117    the terminology employed by WoRMS. We use these attributes to assign each species to a

118    functional group, using a dedicated R function ([https://github.com/tomjwebb/WoRMS-](https://github.com/tomjwebb/WoRMS-)

119    [functional-groups](https://github.com/tomjwebb/WoRMS-functional-groups)) which accesses the WoRMS API using the `worrms` R package [41]. We

120    supplement these functional groups with taxonomic groups to identify fish (using the

121    WoRMS paraphyletic Superclass Pisces; [42]), marine mammals, seabirds, and reptiles. We

122    consolidate functional groups into broad categories for maximum coverage - for example,

123    our 'benthos' group includes all species categorised in WoRMS as endobenthos,

124    epibenthos, hyperbenthos, macrobenthos, meiobenthos, and microbenthos, as well as those

125    originally classified simply as benthos. When separate functional groups are recorded for

126    different life stages, we always use the group for the adult stage. We group together

127    categories with very few species (including meso, macro, neuston) and species with no

128    functional group classification into the single category 'other/unknown'. For fish we include

129    an additional grouping variable based on the broad habitat categories recorded in FishBase

130    [10] accessed using the `rfishbase` package [43], classifying 17,568 of 18,261 species as

131    bathydemersal, bathypelagic, benthopelagic, demersal, pelagic-oceanic, pelagic-nertitic, or

132    reef-associated.

133

134    The WoRMS database includes links to other major biodiversity databases (table 1), and we

135    exploit these to compare the state of biodiversity information availability across axes of

136    biogeography, genetics, conservation, and molecular taxonomy. Specifically, we record for

137    each species its total number of occurrences in the Ocean Biogeographic Information

138    System (OBIS, [22]), and its total number of nucleotide sequences in GenBank. The

139    taxonomy in OBIS is standardised to WoRMS, making these links straightforward, and

140    GenBank's taxonomic information is generally reliable for marine animals [4] meaning that

141    links between WoRMS and GenBank are likely to robustly associate relevant sequences

142    with the correct taxonomic identifier. We also record for each species its IUCN conservation

143    assessment category (if available), and whether or not it has DNA barcodes listed in the

144    Barcode of Life Data System (BOLD).

145

146    Using our tidy database linking the diversity data sources shown in table 1, we then

147    summarise the availability of biodiversity data across all marine animals as follows. First, we

148    consider the two major quantitative databases, OBIS and GenBank. We calculate the

149    proportion of species within each phylum with records in each of these databases, and the

150    distribution of records between species within each phylum. To derive an indication of

151    relative data availability across functional groups, highlighting groups that are particularly

152    highly likely (or unlikely) to occur in the dataset, and those which tend to have more records

153    when they are present, we model data availability across functional groups. We apply a two-

154    step hurdle process, because of the high degree of zero-inflation in our data [44]. To assess

155    whether certain functional groups were better represented in the databases than others, we

156    model presence of species in OBIS or GenBank using a binomial GLM of the form species

157    presence ~ functional group, and we model the distribution of counts (OBIS records or

158    GenBank nucleotides) between functional groups, for those species present in the data

159    source, using a zero-truncated negative binomial GLM. These hurdle models are

160    implemented using the `hurdle` function in the `pscl` package [44,45]. For visualisation, we

161    plot the exponentiated binomial coefficients from the zero component of the model, which

162  shows the ratio of the probability of getting a non-zero to a zero observation within a

163  functional group. We also plot the predicted counts for the subset of species in each

164  functional group with non-zero counts.

165

166  To assess whether data availability is correlated across data sources, we use categorical

167  scales of numbers of records per species in both OBIS and GenBank, using categories

168  bounded by upper limits of 0, 1, 10, 100, 1,000, 10,000, 100,000, and a final category of

169  >100,000 records. We use mosaic plots [46], created using the `ggmosaic` R package [47],

170  to illustrate the distribution of GenBank count categories for each OBIS count category. We

171  also consider how IUCN conservation assessments are distributed across species in

172  different functional groups, and between species present and absent in OBIS, and we

173  compare the number of OBIS occurrence records between species in different IUCN

174  categories. To simplify this analysis, we aggregate to the following IUCN assessment

175  categories: Not Assessed, Data Deficient (i.e., formally assessed but insufficient data to

176  assign the species to a threat category), Threatened (formally assessed as Vulnerable,

177  Endangered, Critically Endangered, Conservation Dependent, Extinct in the Wild, or Extinct),

178  and Non-threatened (formally assessed as Near Threatened or Least Concern). We perform

179  a similar analysis comparing species presence or absence in the Barcode of Life database

180  with presence in OBIS and number of OBIS records.

181

182  All data and links were extracted from WoRMS on 2020-01-11 and the statistics we report

183  are correct as of that date. Manipulation, visualisation, and analysis is performed in R 3.6.2

184  [48] using RStudio 1.2.5033 [49] and the `tidyverse` suite of packages [50] as well as

185  `worrms` [41] to access the WoRMS API and `rfishbase` {Boettiger:2012bz} to access

186  FishBase, and the plotting packages `ggmosaic [47]`, `ggbeeswarm` [51] and `patchwork`

187  `[52]`. All data used in this article are publicly available via WoRMS. The processed

188  summary data we use for our analysis is openly available under a Creative Commons

189  Attribution 4.0 International License in the Marine Data Archive

190  (https://doi.org/10.14284/417). R code to replicate our analyses and figures is available

191  via https://github.com/tomjwebb/linking_marine_diversity_data and is archived on Figshare

192  via the University of Sheffield's Online Research Data repository

193  here: https://doi.org/10.15131/shef.data.12833891.


194  Results

195  Our final dataset consisted of 206,849 valid marine animal species, from 32 phyla and 89

196  classes. Of these, 106,213 (51%) have at least one occurrence record listed in OBIS (table

197  2). Of these, 18,869 (18% of species in OBIS, 9% of all species) are represented by just a

198  single occurrence record (table 2), while one species (Atlantic Cod, Gadus morhua) has over

199  a million occurrence records (1,108,463). Overall, there are 45,974,726 OBIS occurrence

200  records across all species. 36,094 (17%) of all species have at least one nucleotide

201  recorded in GenBank, while 8 species (five fish, the Antarctic Minke Whale Balaenoptera

202  bonaerensis, the tunicate Ciona intestinalis and the California Sea Hare Aplysia californica)

203  have more than a million. Overall the species in our database total 56,846,294 GenBank

204  nucleotides. Furthermore, 13,179 species have had their conservation status assessed by

205  the IUCN, and 25,272 have at least one DNA barcode in the Barcode of Life database.

206

207  The distribution of OBIS and GenBank records across animal phyla and functional groups is

208  shown in Fig 1. At least one species from every phylum has records in either OBIS or

209  GenBank, with all phyla except Loricifera (which has just 29 species) represented in both

210  databases (Fig 1A). Across all phyla, just over half (55%) of all species are represented in

211  one or other database. Most species that are present in OBIS have only a few occurrence

212  records, with median values of records ranging from 1 to 92 across phyla (Fig 1B). A similar

213  pattern is observed for GenBank nucleotides (fig 1C), with median values between 1 and 94

214  except in phyla Orthonectida and Placozoa, both of which have only two species

215  represented in GenBank, one of which has several thousand nucleotides (in Orthonectida,

216  Intoshia linei has 3,522, in Placozoa, Trichoplax adhaerens has 29,176).

217

218  Data availability is variable across functional groups (fig 1B, C; fig 2). Modelling the presence

219  or absence of species in OBIS in a binomial GLM shows that species of fish, mammal, bird,

220  and reptile are much more likely to have occurrences in OBIS than are benthic or

221  zooplankton species, with nekton falling in between, and species with unknown or other

222  functional group classification the least likely to have occurrence records (fig 2A). A broadly

223  similar pattern holds when modelling the number of occurrence records for those species

224  with at least 1 (fig 2B), with the vertebrate taxa again tending to have most records, although

225  distinctions between vertebrates and other groups are less stark. Benthic invertebrates

226  typically have few OBIS records, but zooplankton that do occur in OBIS tend to have more

227  records than nekton. In GenBank, birds, reptiles and mammals are most likely to be present

228  in the database, followed by fish, nekton, and zooplankton, with benthos and other/unknown

229  functional groups least likely to be represented (fig 2C). The rank order changes somewhat

230  when considering number of nucleotides across species present in GenBank (fig 2D), with

231  most records from mammals and reptiles, followed but birds and fish. Nekton tend to have

232  fewest records, but there is considerable variability within all major groups. Data availability

233  in both major databases is broadly similar across fish habitat groupings (figure S1, S2).

234

235  Considering the joint distribution of species across OBIS and GenBank categorical scales,

236  93,519 (45%) species have no records in either database (table 2, fig 3A). In general,

237  species with more records in OBIS also tend to have more nucleotides in GenBank (table 2,

238  fig 3), indicating that these different biodiversity data aggregators have similar biases in

239  terms of the known marine biodiversity that they encompass. There are exceptions though:

240  in particular several species have many (>100,000) GenBank nucleotides but very few (if

241  any) OBIS records (table 3).

242

243 A similar pattern is evident when examining the distribution of OBIS records across different

244 IUCN assessment categories. In general, and across functional groups, the proportion of

245 species with records in OBIS is higher in assessed species (threatened and non-threatened)

246 than it is in unassessed or data-deficient species: overall, 84% of threatened and 94% of

247 non-threatened species have occurrence records in OBIS, compared to 75% of data-

248 deficient and 49% of unassessed species (table 4A). Considering only those species with

249 records in OBIS, there is considerable variation within and between IUCN categories in the

250 number of occurrence records per species, but a general tendency is apparent in all

251 functional groups for species in threatened and non-threatened categories to have more

252 occurrence records than those in data-deficient and unassessed categories (fig 4A).

253

254 Species with DNA barcodes are disproportionately likely to also have occurrence records in

255 OBIS: 45% of species with no record in the Barcode of Life database have at least one

256 occurrence record in OBIS, compared to 89% of species with a barcode (table 4B). In

257 addition, in all functional groups, species with barcodes tend to have more OBIS records

258 than those which do not (fig 4B).

259 Discussion

260 Using the taxonomic backbone of the World Register of Marine Species [37] we have

261 summarised data availability across axes of biogeography, genetics, molecular taxonomy,

262 and conservation status for 206,849 marine animal species. This presents a mixed picture.

263 One the one hand, 91,828 (44%) species have no records in any of these databases, and

264 are represented only by their name. This is considerably higher than the 27% of plant

265 species with no information other than their name [36], although of course the marine

266 environment represents far larger habitable volume [53] and marine animals are a much

267 more diverse taxonomic group. Only 6,688 marine animal species (3%) have records in all

268 four of the datasets that we consider – again, rather lower than the 18% of broadly-covered

269 plant species [36]. At the same time, it is important to remember that presence in a dataset

270 does not imply extensive knowledge: among the 106,203 species with records in OBIS, for

271 example, the median number of recorded occurrences is just 7, and 18% of these species

272 (18,869 species) are known from only a single occurrence. Nonetheless, the distribution of

273 biogeographic and genetic information across the animal tree of life is extensive, with all

274 animal phyla represented in at least one database (fig 1). Data availability tends to be biased

275 towards well-known taxa and functional groups (especially vertebrates; figs 1, 2, 4), in

276 agreement with previous assessments (e.g. [32]), but the subset of 225 species with >1,000

277 occurrences in OBIS and >1,000 nucleotides in GenBank is drawn from 10 phyla and 27

278 classes, representing all major functional groups, and most of them have a barcode in BOLD

279 (214 species), and have been assessed by the IUCN as something other than data deficient

280 (102 non-threatened, 23 threatened species). For these diverse marine animal species then,

281 it is reasonable to propose that the information available across multiple sources can be

282 translated into knowledge about their distribution, evolutionary relationships, and

283 conservation status.

284

285 The broad positive correlation between data availability across different sources (table 2,

286 table 4, fig 3) reinforces previous findings that species with good information on one facet of

287 their biology and ecology tend to be well represented in other databases too, both in plants

288 [36] and in marine species [35]. These information-rich species are likely to be those most

289 easily and frequently observed, or those of high economic or cultural value, and so will not

290 be a random subset of all species. However, the consequences of biases towards data

291 availability from these common species will vary depending on the specific question of

292 interest. For instance, ecosystem function may be driven largely by just those common

293 species that tend to be so well known [54]; but rare species will clearly be of great interest to

294 conservationists, and may indeed sometimes contribute unique trait combinations to marine

295 communities [55].

296

297    In terrestrial conservation, considerable concern has been expressed over the likely

298    conservation status of species too poorly known to formally assess, as they tend to have

299    characteristics (rarity, small ranges, occurring in poorly studied regions) which will

300    predispose them to be at risk [56]. For some marine taxa this appears to be the case too,

301    with high rates of extinction risk predicted for European sharks and rays formally assessed

302    as Data Deficient [57], and low levels of conservation assessment in poorly-known marine

303    groups may contribute to low overall documented levels of extinction risk [58]. On the other

304    hand, the fact that the biggest data gaps in marine biodiversity tend to be in remote habitats

305    largely inaccessible to humans (e.g. the deep pelagic ocean; [59]), and the highest rates of

306    discoveries of new species and habitats are also in the deep sea [60,61], provides some

307    contrast with the terrestrial situation, and may insulate these poorly-known species

308    somewhat from human pressures. However, some patterns still hold in the deep sea, such

309    as the tendency for widespread species to be encountered and described first [62], meaning

310    that many of the species not yet present in major databases may be genuinely rare. Given

311    the acceleration of human activities into previously unexploited regions of the oceans [63],

312    with new threats including deep sea mining [64] and exploitation of the mesopelagic [65], it

313    seems unwise to assume that the large fraction of marine biodiversity that remains poorly

314    known is not at risk. Given the fact that Data Deficient conservation assessments are twice

315    as frequent in marine versus non-marine taxa [34], data-driven predictive conservation

316    assessments [57,66,67] which rely on some of the kinds of data we consider here (spatial

317    distribution, evolutionary relationships, ecological guilds) combined with biological traits may

318    prove to be especially valuable tools.

319

320    An aim of this study was to flag priorities for future work. One important point is that the

321    major publicly available databases on which we draw do not constitute the sum total of data

322    ever collected on marine species. This is particularly the case for occurrence data, as

323    globally researchers have yet to adopt the routine deposition of species occurrences in OBIS

324    as a cultural norm, in the way that genetic sequence data is deposited in GenBank. To this

325   end, improving incentives for researchers to add their data to global repositories in an

326   important goal [25], while data archaeology and rescue initiatives can help to ensure that

327   historical data are captured [68]. Equally, it remains vital that ongoing survey schemes are

328   properly valued [69], even as novel exploration is planned. At the same time, our

329   quantification exercise can help to identify groups of species where a little additional

330   research effort in one area would quickly result in a more valuable dataset. One candidate

331   set of species might be those that are frequently observed but poorly represented in other

332   databases. For instance, 1,216 species have >1,000 OBIS records but <10 GenBank

333   nucleotides; and over half of the 3,533 species with >1,000 OBIS occurrences are either not

334   assessed by the IUCN (1,876 species) or data-deficient (82 species). The fact that almost

335   90% (3,163) of these species have DNA barcodes in BOLD is encouraging, however,

336   suggesting considerable potential for an increasing role for molecular studies to address a

337   wide range of questions in marine ecology [70].

338

339   Mining the spatial information already present in other databases also has potential for

340   supplementing existing occurrence datasets. In this study we relied on existing links between

341   WoRMS and GenBank and BOLD, which simply summarise the number of nucleotides or

342   barcodes present for each species. The spatial meta-data stored in the sequence databases

343   provides an additional source of information, although in GenBank this data is relatively

344   unstructured. Searching the GenBank nucleotide database, we found just 1,437 records for

345   animals which contained a lat-lon field; matching this to our list of marine animals reduced

346   this further to 183 records from 42 species. Nonetheless, even from this small set of species,

347   21 do not have occurrence records in OBIS, suggesting that mining GenBank for spatial data

348   would likely add valuable information for a small number of species. Various methods have

349   been developed to attempt this, based around mining spatial information from the full text of

350   associated publications [71,72] with initiatives such as the Genomic Observatories

351   MetaDatabase (GEOME, https://geome-db.org) also seeking to simplify access to meta-data

352   from sequence datasets.

353

354 BOLD typically does store spatial data for individual specimens in a well-structured manner,

355 only some of which has been harvested by OBIS. In our dataset, 3,117 species have BOLD

356 barcodes but no OBIS records. Several of these are parasites, which we know are not well

357 recorded in OBIS (e.g. Schistocephalus solidus, 718 barcodes; Anguillicoloides crassus, 508

358 barcodes) but there are free-living marine species too, such as the Gastropod mollusc

359 Littoraria sinensis (257 barcodes) and the Copepod Calanoides natalis (183 barcodes).

360 Accessing the specimen data from BOLD using the `bold` R package [73] for these two

361 species reveals that none of the L. sinensis have information in the latitude and longitude

362 fields, but full geographic information is available for 227 specimens for Calanoides natalis.

363 Although none of these locations are currently recorded in OBIS, some are in GBIF,

364 highlighting the often complex pipelines from data providers to global data aggregators.

365 Improving pipelines from genetic databases to occurrence databases is currently a priority

366 for OBIS (W. Appeltans, OBIS Project Manager, pers. comm.).

367

368 Finally, the dimensions of diversity that we summarise in this study are somewhat limited.

369 We did not consider the traits of species, for instance, beyond functional groups that indicate

370 habitat affiliation in very broad terms (e.g. benthic vs planktonic). These groupings are

371 already useful as global patterns of diversity are known to differ between them [39], and they

372 can also be used to refine methods of matching species occurrences to global sea

373 temperature datasets [74], helping to predict species responses to climate change [75].

374 Beyond these coarse functional groups, however, traits data remain scarce even in

375 reasonably common marine species in well-studied regions [35], and despite many efforts at

376 collating traits – including within WoRMS; [76] - there is still no widely-adopted central

377 standard [77]. Certain groups are well covered by existing initatives (e.g. FishBase [10], the

378 Coral Trait Database [11]), and whether a single overarching portal to cover the immense

379 diversity of marine lifeforms is possible – or even desirable – remains open for discussion.

380 However, it is certainly the case that multiple smaller-scale projects collect valuable traits

381    data for a subset of species which is typically made available (if at all) via supplementary

382    material or bespoke web portals, at risk of being lost to the community. A wider adoption of

383    principles embedded in initiatives like the Open Traits Network [7] would ensure

384    interoperability of these small, project-specific traits datasets, maximising the availability of

385    information on key traits for the largest possible fraction of marine diversity. Readily

386    availabilie information on even just a few traits (e.g. body size, longevity, fecundity,

387    planktonic larval duration) would help to test predictions from biodiversity models, embed life

388    history theory into marine conservation, and predict the consequences of human activities on

389    marine diversity [39,78-80].

390

391    The stocktake of marine biodiversity data availability that we have undertaken here adds to

392    previous efforts focused on occurrence data [19,32,81]. While we reveal a similar story of

393    gaps and biases across other data sources, there is considerable overlap in coverage too,

394    and overall the potential to link dimensions of marine animal diversity is now high. The

395    priority now should be to build on the substantial community-built foundations and to improve

396    the pipeline from raw data to interoperable data products, both as a resource for

397    fundamental macroecological research and to facilitate effective stewardship of our blue

398    planet.

408 References

409 1.      1. Redelings, B. D. & Holder, M. T. 2017 A supertree pipeline for summarizing
410         phylogenetic and taxonomic information for millions of species. PeerJ **5**, e3058.
411         (doi:10.7717/peerj.3058)

412 2.      Hinchliff, C. E. et al. 2015 Synthesis of phylogeny and taxonomy into a
413         comprehensive tree of life. P Natl Acad Sci USA **112**, 12764–12769.
414         (doi:10.1073/pnas.1423041112)

415 3.      Clark, K., Karsch-Mizrachi, I., Lipman, D. J., Ostell, J. & Sayers, E. W. 2016
416         GenBank. Nucleic Acids Res. **44**, D67–D72. (doi:10.1093/nar/gkv1276)

417 4.      Leray, M., Knowlton, N., Ho, S.-L., Nguyen, B. N. & Machida, R. J. 2019
418         GenBank is a reliable resource for 21st century biodiversity research. P Natl
419         Acad Sci USA **116**, 22651–22656. (doi:10.1073/pnas.1911714116)

420 5.      GBIF 2020 GBIF Home Page. https://www.gbif.org.

421 6.      IUCN 2020 The IUCN Red List of Threatened Species. Version 2020-1.
422         httpswww.iucnredlist.org.

423 7.      Gallagher, R. et al. 2019 The Open Traits Network: Using Open Science
424         principles to accelerate trait-based science across the Tree of Life.
425         ecoevorxiv.org. (doi:https://doi.org/10.32942/osf.io/kac45)

426 8.      Bennett, J. M. et al. 2018 GlobTherm, a global database on thermal tolerances
427         for aquatic and terrestrial organisms. Sci. Data **5**, 180022.
428         (doi:10.1038/sdata.2018.22)

429 9.      Makarieva, A., Gorshkov, V. & LI, B. 2005 Biochemical universality of living
430         matter and its metabolic implications. Funct Ecol

431 10.     Froese, R. & Pauly, D. 2019 FishBase. World Wide Web electronic publication.
432         version (12/2019). www.fishbase.org.

433 11.     Madin, J. S. et al. 2016 The Coral Trait Database, a curated database of trait
434         information for coral species from the global oceans. Sci. Data **3**, 178–22.
435         (doi:10.1038/sdata.2016.17)

436 12.     Kranstauber, B., Cameron, A., Weinzerl, R., Fountain, T., Tilak, S., Wikelski, M.
437         & Kays, R. 2011 The Movebank data model for animal tracking. Environmental
438         Modelling & Software **26**, 834–835. (doi:10.1016/j.envsoft.2010.12.005)

439 13.     Wikelski, M., Davidson, S. C. & Kays, R. 2020 Movebank: archive, analysis and
440         sharing of animal movement data. www.movebank.org.

441 14.     Poelen, J. H., Simons, J. D. & Mungall, C. J. 2014 Global biotic interactions: An
442         open infrastructure to share and analyze species-interaction datasets.
443         Ecological Informatics **24**, 148–159. (doi:10.1016/j.ecoinf.2014.08.005)

444 15.     Tanhua, T. et al. 2019 Ocean FAIR Data Services. Front. Mar. Sci. **6**, 92.
445         (doi:10.3389/fmars.2019.00440)

446    16.    Silvertown, J. 2009 A new dawn for citizen science. Trends Ecol Evol **24**, 467–
447           471. (doi:10.1016/j.tree.2009.03.017)

448    17.    Chandler, M. et al. 2017 Contribution of citizen science towards international
449           biodiversity monitoring. Biological Conservation **213**, 280–294.
450           (doi:10.1016/j.biocon.2016.09.004)

451    18.    Hyder, K., Townhill, B., Anderson, L. G., Delany, J. & Pinnegar, J. K. 2015 Can
452           citizen science contribute to the evidence-base that underpins marine policy?
453           Marine Policy **59**, 112–120. (doi:10.1016/j.marpol.2015.04.022)

454    19.    Edgar, G. J., Bates, A. E., Bird, T. J., Jones, A. H., Kininmonth, S., Stuart-Smith,
455           R. D. & Webb, T. J. 2015 New Approaches to Marine Conservation Through
456           Scaling Up of Ecological Data. Annual Review of Marine Science **8**,
457           150807173619006. (doi:10.1146/annurev-marine-122414-033921)

458    20.    Edgar, G. J. & Stuart-Smith, R. D. 2014 Systematic global assessment of reef
459           fish communities by the Reef Life Survey program. Sci. Data **1**, 1–8.
460           (doi:10.1038/sdata.2014.7)

461    21.    Snelgrove, P. V. R. 2010 Discoveries of the Census of Marine Life: Making
462           Ocean Life Count. 1st edn. Cambridge University Press.

463    22.    OBIS 2020 Ocean Biodiversity Information System. www.iobis.org.

464    23.    Mieszkowska, N., Sugden, H., Firth, L. B. & Hawkins, S. J. 2014 The role of
465           sustained observations in tracking impacts of environmental change on marine
466           biodiversity and ecosystems. Philos T R Soc A **372**, 20130339–20130339.
467           (doi:10.1098/rsta.2013.0339)

468    24.    Miloslavich, P. et al. 2018 Essential ocean variables for global sustained
469           observations of biodiversity and ecosystem changes. Global Change Biol **105**,
470           10456. (doi:10.1111/gcb.14108)

471    25.    Dornelas, M. et al. 2019 Towards a macroscope: Leveraging technology to
472           transform the breadth, scale and resolution of macroecological data. Global Ecol
473           Biogeogr **28**, 1937–1948. (doi:10.1111/geb.13025)

474    26.    Basset, A. & Los, W. 2012 Biodiversity e-Science: LifeWatch, the European
475           infrastructure on biodiversity and ecosystem research. Plant Biosystems - An
476           International Journal Dealing with all Aspects of Plant Biology **146**, 780–782.
477           (doi:10.1080/11263504.2012.740091)

478    27.    La Salle, J., Williams, K. J. & Moritz, C. 2016 Biodiversity analysis in the digital
479           era. Philos T R Soc B **371**. (doi:10.1098/rstb.2015.0337)

480    28.    Hardisty, A. R. et al. 2019 The Bari Manifesto: An interoperability framework for
481           essential biodiversity variables. Ecological Informatics **49**, 22–31.
482           (doi:10.1016/j.ecoinf.2018.11.003)

483    29.    Maitner, B. S. et al. 2018 The bien r package: A tool to access the Botanical
484           Information and Ecology Network (BIEN) database. Methods in Ecology and
485           Evolution **9**, 373–379. (doi:10.1111/2041-210X.12861)

486  30.  Costello, M., Coll, M., Danovaro, R., Halpin, P. & Ojaveer, H. 2010 A Census of
487      Marine Biodiversity Knowledge, Resources, and Future Challenges. PLoS ONE

488  31.  Snelgrove, P. et al. 2016 Global Patterns in Marine Biodiversity. In The First
489      Global Integrated Marine Assessment World Ocean Assessment, un.org.

490  32.  Miloslavich, P. et al. 2016 Extent of Assessment of Marine Biological Diversity.
491      In The First Global Integrated Marine Assessment World Ocean Assessment
492      (eds L. Inniss & A. Simcock), United Nations.

493  33.  Appeltans, W. et al. 2012 The magnitude of global marine species diversity.
494      Curr. Biol. **22**, 2189–2202. (doi:10.1016/j.cub.2012.09.036)

495  34.  Webb, T. J. & Mindel, B. L. 2015 Global Patterns of Extinction Risk in Marine
496      and Non-marine Systems. Current Biology **25**, 506–511.
497      (doi:10.1016/j.cub.2014.12.023)

498  35.  Tyler, E. H. M., Somerfield, P. J., Berghe, E. V., Bremner, J., Jackson, E.,
499      Langmead, O., Palomares, M. L. D. & Webb, T. J. 2012 Extensive gaps and
500      biases in our knowledge of a well-known fauna: implications for integrating
501      biological traits into macroecology. Global Ecol Biogeogr **21**, 922–934.
502      (doi:10.1111/j.1466-8238.2011.00726.x)

503  36.  Cornwell, W. K., Pearse, W. D., Dalrymple, R. L. & Zanne, A. E. 2019 What we
504      (don't) know about global plant diversity. Ecography, ecog.04481.
505      (doi:10.1111/ecog.04481)

506  37.  WoRMS Editorial Board 2020 World Register of Marine Species.
507      http://www.marinespecies.org.

508  38.  Freckleton, R. P. & Jetz, W. 2009 Space versus phylogeny: disentangling
509      phylogenetic and spatial signals in comparative data. P R Soc B **276**, 21–30.
510      (doi:10.1098/rspb.2008.0905)

511  39.  Worm, B. & Tittensor, D. P. 2018 A Theory of Global Biodiversity. Princeton and
512      Oxford: Princeton University Press. (doi:10.2307/j.ctt1zkjz6q)

513  40.  Harfoot, M. B. J., Newbold, T., Tittensor, D. P., Emmott, S., Hutton, J.,
514      Lyutsarev, V., Smith, M. J., Scharlemann, J. P. W. & Purves, D. W. 2014
515      Emergent Global Patterns of Ecosystem Structure and Function from a
516      Mechanistic General Ecosystem Model. PLoS Biol **12**, e1001841.
517      (doi:10.1371/journal.pbio.1001841)

518  41.  Chamberlain, S. 2019 World Register of Marine Species (WoRMS) Client. R
519      package worrms version 0.4.0. httpsCRAN.R-project.orgpackageworrms.

520  42.  WoRMS 2020 Pisces. httpswww.marinespecies.orgaphia.phpptaxdetailsid.

521  43.  Boettiger, C., Lang, D. T. & Wainwright, P. C. 2012 rfishbase: exploring,
522      manipulating and visualizing FishBase data from R. Journal of Fish Biology **81**,
523      2030–2039. (doi:10.1111/j.1095-8649.2012.03464.x)

524  44.  Zeileis, A., Kleiber, C. & Jackman, S. 2008 Regression models for count data in
525      R. Journal of Statistical Software **27**, 1–25.

526   45.   Jackman, S. 2020 pscl: Classes and Methods for R Developed in the Political
527         Science Computational Laboratory. United States Studies Centre, University of
528         Sydney. Sydney, New South Wales, Australia. R package version 1.5.5.
529         httpsgithub.comatahkpscl.

530   46.   Hofmann, H. 2008 Mosaic Plots and Their Variants. In Handbook of Data
531         Visualisation (eds C.-H. Chen W. Härdle & A. Unwin), Berlin Heidelberg.

532   47.   Jeppson, H., Hofmann, H. & Cook, D. 2018 ggmosaic: Mosaic Plots in the
533         'ggplot2' Framework. R package version 0.2.0. https://CRAN.R-
534         project.orgpackageggmosaic.

535   48.   R Core Team 2019 R: A language and environment for statistical computing.
536         https://www.R-project.org.

537   49.   RStudio Team 2019 RStudio: Integrated Development for R.
538         httpwww.rstudio.com.

539   50.   Wickham, H. et al. 2019 Welcome to the Tidyverse. Journal of Open Source
540         Software **4**, 1686. (doi:10.21105/joss.01686)

541   51.   Clarke, E. & Sherrill-Mix, S. 2017 ggbeeswarm: Categorical Scatter (Violin
542         Point) Plots. R package version 0.6.0. httpsCRAN.R-
543         project.orgpackageggbeeswarm.

544   52.   Pedersen, T. L. 2019 patchwork: The Composer of Plots. R package version
545         1.0.0. httpsCRAN.R-project.orgpackagepatchwork.

546   53.   Dawson, M. N. 2012 Species richness, habitable volume, and species densities
547         in freshwater, the sea, and on land. Frontiers of Biogeography **4**, fb_12675.

548   54.   Gaston, K. & Fuller, R. 2008 Commonness, population depletion and
549         conservation biology. Trends Ecol Evol **23**, 14–19.

550   55.   Mouillot, D. et al. 2013 Rare Species Support Vulnerable Functions in High-
551         Diversity Ecosystems. PLoS Biol **11**, e1001569.
552         (doi:10.1371/journal.pbio.1001569)

553   56.   Scheffers, B. R., Joppa, L. N., Pimm, S. L. & Laurance, W. F. 2012 What we
554         know and don't know about Earth's missing biodiversity. Trends Ecol Evol **27**,
555         501–510. (doi:10.1016/j.tree.2012.05.008)

556   57.   Walls, R. H. L. & Dulvy, N. K. 2019 Predicting the conservation status of
557         Europe's Data Deficient sharks and rays. bioRxiv **276**, 614776.
558         (doi:10.1101/614776)

559   58.   Mindel, B. L., Webb, T. J., Neat, F. C. & Blanchard, J. L. 2016 A trait-based
560         metric sheds new light on the nature of the body size–depth relationship in the
561         deep sea. J Anim Ecol **85**, 427–436. (doi:10.1111/1365-2656.12471)

562   59.   Webb, T. J., Vanden Berghe, E. & O'Dor, R. 2010 Biodiversity's big wet secret:
563         the global distribution of marine biological records reveals chronic under-
564         exploration of the deep pelagic ocean. PLoS ONE **5**, e10223.
565         (doi:10.1371/journal.pone.0010223)

566  60.  Ramirez-Llodra, E. et al. 2010 Deep, diverse and definitely different: unique
567        attributes of the world's largest ecosystem. Biogeosciences **7**, 2851–2899.
568        (doi:10.5194/bg-7-2851-2010)

569  61.  Danovaro, R., Snelgrove, P. V. R. & Tyler, P. 2014 Challenging the paradigms
570        of deep-sea ecology. Trends Ecol Evol **29**, 465–475.
571        (doi:10.1016/j.tree.2014.06.002)

572  62.  Higgs, N. D. & Attrill, M. 2015 Biases in biodiversity: wide-ranging species are
573        discovered first in the deep sea. Front. Mar. Sci. **2**, 717.
574        (doi:10.3389/fmars.2015.00061)

575  63.  Jouffray, J. B., Blasiak, R., Norström, A. V., Österblom, H. & Nyström, M. 2020
576        The Blue Acceleration: The Trajectory of Human Expansion into the Ocean. One
577        Earth **2**, 43–54. (doi:10.1016/j.oneear.2019.12.016)

578  64.  Jones, D. O. B., Amon, D. J. & Chapman, A. S. A. 2018 Mining Deep-Ocean
579        Mineral Deposits: What are the Ecological Risks? Elements **14**, 325–330.
580        (doi:10.2138/gselements.14.5.325)

581  65.  Hidalgo, M. & Browman, H. I. 2019 Developing the knowledge base needed to
582        sustainably manage mesopelagic resources. ICES Journal of Marine Science
583        **76**, 609–615. (doi:10.1093/icesjms/fsz067)

584  66.  Jetz, W. & Freckleton, R. P. 2015 Towards a general framework for predicting
585        threat status of data-deficient species from phylogenetic, spatial and
586        environmental information. Philos T Roy Soc B **370**, 20140016–20140016.
587        (doi:10.1098/rstb.2014.0016)

588  67.  González-del-Pliego, P., Freckleton, R. P., Edwards, D. P., Koo, M. S.,
589        Scheffers, B. R., Pyron, R. A. & Jetz, W. 2019 Phylogenetic and Trait-Based
590        Prediction of Extinction Risk for Data-Deficient Amphibians. Current Biology **29**,
591        1557–1563.e3. (doi:10.1016/j.cub.2019.04.005)

592  68.  Faulwetter, S. et al. 2016 EMODnet Workshop on mechanisms and guidelines
593        to mobilise historical data into biogeographic databases. RIO **2**, e9774–28.
594        (doi:10.3897/rio.2.e9774)

595  69.  Mieszkowska, N., Sugden, H., Firth, L. B. & Hawkins, S. J. 2014 The role of
596        sustained observations in tracking impacts of environmental change on marine
597        biodiversity and ecosystems. Philos T R Soc A **372**, 20130339.
598        (doi:10.1098/rsta.2013.0339)

599  70.  Goodwin, K. D., Thompson, L. R., Duarte, B., Kahlke, T., Thompson, A. R.,
600        Marques, J. C. & Caçador, I. 2017 DNA Sequencing as a Tool to Monitor Marine
601        Ecological Status. Front. Mar. Sci. **4**, e1002358.
602        (doi:10.3389/fmars.2017.00107)

603  71.  Tahsin, T., Weissenbacher, D., Rivera, R., Beard, R., Firago, M., Wallstrom, G.,
604        Scotch, M. & Gonzalez, G. 2016 A high-precision rule-based extraction system
605        for expanding geospatial metadata in GenBank records. J Am Med Inform
606        Assoc **23**, 934–941. (doi:10.1093/jamia/ocv172)

607  72.  Tahsin, T., Weissenbacher, D., O'Connor, K., Magge, A., Scotch, M. &
608        Gonzalez-Hernandez, G. 2017 GeoBoost: accelerating research involving the

609    geospatial metadata of virus GenBank records. Bioinformatics **34**, 1606–1608.
610    (doi:10.1093/bioinformatics/btx799)

611    73.    Chamberlain, S. In press. bold: Interface to Bold Systems API. R package
612    version 1.1.0. httpsCRAN.R-project.orgpackagebold.

613    74.    Webb, T. J., Lines, A. & Howarth, L. M. 2020 Occupancy-derived thermal
614    affinities reflect known physiological thermal limits of marine species. Ecology
615    and Evolution **75**, 209. (doi:10.1002/ece3.6407)

616    75.    Pinsky, M. L., Selden, R. L. & Kitchel, Z. J. 2020 Climate-Driven Shifts in Marine
617    Species Ranges: Scaling from Organisms to Communities. Annual Review of
618    Marine Science **12**, 153–179. (doi:10.1146/annurev-marine-010419-010916)

619    76.    Costello, M. J., Claus, S., Dekeyzer, S., Vandepitte, L., Tuama, É. Ó., Lear, D. &
620    Tyler-Walters, H. 2015 Biological and ecological traits of marine species. PeerJ
621    **3**, e1201. (doi:10.7717/peerj.1201)

622    77.    Beauchard, O., Veríssimo, H., Queirós, A. M. & Herman, P. M. J. 2017 The use
623    of multiple biological traits in marine community ecology and its potential in
624    ecological indicator development. Ecol Indic **76**, 81–96.
625    (doi:10.1016/j.ecolind.2017.01.011)

626    78.    Kindsvater, H. K., Mangel, M., Reynolds, J. D. & Dulvy, N. K. 2016 Ten
627    principles from evolutionary ecology essential for effective marine conservation.
628    Ecology and Evolution **6**, 2125–2138. (doi:10.1002/ece3.2012)

629    79.    Hiddink, J. G. et al. 2019 Assessing bottom trawling impacts based on the
630    longevity of benthic invertebrates. J Appl Ecol **56**, 1075–1084.
631    (doi:10.1111/1365-2664.13278)

632    80.    Álvarez-Noriega, M., Burgess, S. C., Byers, J. E., Pringle, J. M., Wares, J. P. &
633    Marshall, D. J. 2020 Global biogeography of marine dispersal potential. Nature
634    Ecology & Evolution 2017 1:9 **4**, 1–8. (doi:10.1038/s41559-020-1238-y)

635    81.    Menegotto, A. & Rangel, T. F. 2018 Mapping knowledge gaps in marine
636    diversity reveals a latitudinal gradient of missing species richness. Nature
637    Communications **9**, 4713. (doi:10.1038/s41467-018-07217-7)

638    82.    Ratnasingham, S. & Hebert, P. 2007 BOLD: The Barcode of Life Data System
639    (http://www. barcodinglife. org). Mol. Ecol. Notes (doi:doi: 10.1111/j.1471-
640    8286.2006.01678.x)

641

642

643

644  **Table 1.** Data sources used to link different dimensions of diversity across all marine
645  animals.

| Dimension of diversity | Data source | Data type | Reference |
|---|---|---|---|
| Taxonomy | WoRMS | Authoritative classification and catalogue of marine taxonomic names | [37] |
| Functional Groups | WoRMS | Classification of marine species into broad ecological groups | [37] |
| Biogeography | OBIS | Global database of marine species occurrence records | [22] |
| Genetics | GenBank | The NIH genetic sequence database, an annotated collection of all publicly available DNA sequences | [3] |
| Molecular taxonomy | BOLD | Barcode of Life Data System for DNA barcodes | [82] |
| Conservation status | IUCN Red List | The IUCN Red List of threatened species | [6] |

646

**Table 2.** Breakdown of 206,849 marine animal species by number of global occurrence records in OBIS, and numbers of nucleotide sequences in GenBank.

| Number of OBIS records | Number of GenBank nucleotides | | | | | | | | | Totals | In OBIS |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | 1 | 2-10 | 11-100 | 101-1,000 | 1,001-10,000 | 10,001-100,000 | 100,001-1,000,000 | >1,000,000 | | |
| 0 | 93,519 | 1,312 | 4,484 | 1,164 | 116 | 13 | 15 | 13 | 0 | *100,636* | |
| 1 | 16,905 | 356 | 1,253 | 314 | 33 | 3 | 3 | 2 | 0 | *18,869* | |
| 2-10 | 35,613 | 1,086 | 3,714 | 990 | 122 | 17 | 11 | 8 | 0 | *41,561* | |
| 11-100 | 19,998 | 1,392 | 5,931 | 2,733 | 351 | 32 | 30 | 26 | 2 | *30,495* | |
| 101-1,000 | 4,274 | 594 | 3,334 | 2,917 | 512 | 51 | 35 | 37 | 1 | *11,755* | |
| 1,001-10,000 | 402 | 86 | 630 | 1,113 | 315 | 33 | 53 | 33 | 4 | *2,669* | *106,213* |
| 10,001-100,000 | 42 | 4 | 107 | 406 | 167 | 31 | 22 | 31 | 1 | *811* | |
| 100,001-1,000,000 | 2 | 0 | 0 | 14 | 20 | 5 | 3 | 8 | 0 | *52* | |
| >1,000,000 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | *1* | |
| *Totals* | *170,755* | *4,830* | *19,453* | *9,651* | *1,636* | *185* | *173* | *158* | *8* | **206,849** | |
| *In GenBank* | | *36,094* | | | | | | | | | |

**Table 3.** Species with high numbers of GenBank nucleotide records but few OBIS occurrences, or species with large numbers of OBIS occurrences but few GenBank nucleotides.

| Species | Phylum | Class | Functional Group | GenBank Nucleotides | OBIS Records |
|---|---|---|---|---|---|
| Olavius algarvensis | Annelida | Clitellata | benthos | 173,609 | 0 |
| Capitella teleta | Annelida | Polychaeta | benthos | 208,794 | 1 |
| Platynothrus peltifer | Arthropoda | Arachnida | other/unknown | 106,099 | 0 |
| Caligus rogercresseyi | Arthropoda | Hexanauplia | other/unknown | 628,843 | 0 |
| Proasellus racovitzai | Arthropoda | Malacostraca | benthos | 127,716 | 0 |
| Proasellus ibericus | Arthropoda | Malacostraca | benthos | 150,798 | 0 |
| Bragasellus molinai | Arthropoda | Malacostraca | benthos | 209,419 | 0 |
| Proasellus beticus | Arthropoda | Malacostraca | benthos | 228,033 | 0 |
| Seriola quinqueradiata | Chordata | Actinopterygii | fish | 105,911 | 6 |
| Theragra finnmarchica | Chordata | Actinopterygii | fish | 130,916 | 0 |
| Takifugu flavidus | Chordata | Actinopterygii | fish | 138,301 | 0 |
| Takifugu rubripes | Chordata | Actinopterygii | fish | 466,790 | 5 |
| Molgula tectiformis | Chordata | Ascidiacea | benthos | 106,904 | 0 |
| Halocynthia roretzi | Chordata | Ascidiacea | benthos | 116,123 | 4 |
| Pelecanus crispus | Chordata | Aves | birds | 231,775 | 0 |
| Balaenoptera acutorostrata scammoni | Chordata | Mammalia | mammals | 238,976 | 0 |
| Emydocephalus ijimae | Chordata | Reptilia | reptiles | 157,876 | 0 |
| Hemicentrotus pulcherrimus | Echinodermata | Echinoidea | benthos | 153,541 | 3 |
| Apostichopus parvimensis | Echinodermata | Holothuroidea | benthos | 166,764 | 1 |
| Apostichopus japonicus | Echinodermata | Holothuroidea | benthos | 401,310 | 4 |
| Cumia reticulata | Mollusca | Gastropoda | benthos | 144,517 | 2 |
| Amphimedon queenslandica | Porifera | Demospongiae | benthos | 142,554 | 9 |
| Thunnus alalunga | Chordata | Actinopterygii | fish | 0 | 114,485 |
| Chrysophrys auratus | Chordata | Actinopterygii | fish | 0 | 104,066 |

**Table 4A** Breakdown of marine animal species by functional group and IUCN Assessment status. Listed for each IUCN assessment status are the total number of species per functional group, the number of these species with occurrences in OBIS, and the associated percentage.

| Functional Group | IUCN Assessment Status | | | | | | | | | | | |
| | Not assessed | | | Data Deficient | | | Threatened | | | Non-threatened | | |
| | N(species) | N(species in OBIS) | % species in OBIS | N(species) | N(species in OBIS) | % species in OBIS | N(species) | N(species in OBIS) | % species in OBIS | N(species) | N(species in OBIS) | % species in OBIS |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Benthos | 144,097 | 73,610 | 51% | 749 | 530 | 71% | 305 | 258 | 85% | 1,400 | 1,206 | 86% |
| Zooplankton | 5,742 | 3,027 | 53% | 0 | 0 | - | 4 | 2 | 50% | 2 | 2 | 100% |
| Nekton | 3,076 | 1,878 | 61% | 160 | 127 | 79% | 7 | 2 | 29% | 156 | 151 | 97% |
| Fish | 8,599 | 6,161 | 72% | 1,780 | 1,350 | 76% | 523 | 457 | 87% | 7,359 | 6,997 | 95% |
| Mammals | 66 | 26 | 39% | 20 | 17 | 85% | 36 | 29 | 81% | 70 | 68 | 97% |
| Birds | 179 | 71 | 40% | 1 | 0 | 0% | 125 | 92 | 74% | 382 | 340 | 89% |
| Reptiles | 20 | 9 | 45% | 21 | 14 | 67% | 11 | 11 | 100% | 44 | 37 | 84% |
| Other / Unknown | 31,891 | 9,725 | 31% | 3 | 3 | 100% | 10 | 4 | 40% | 11 | 9 | 82% |
| *Totals* | *193,670* | *94,507* | *49%* | *2,734* | *2,041* | *75%* | *1,021* | *855* | *84%* | *9,424* | *8,810* | *94%* |

**Table 4B** Breakdown of marine animal species by functional group and presence in the BOLD DNA Barcode database. Listed for species absente from or present in BOLD are the total number of species per functional group, the number of these species with occurrences in OBIS, and the associated percentage.

| Functional Group | In Barcode of Life Database? | | | | | |
| | No | | | Yes | | |
| | N(species) | N(species in OBIS) | % species in OBIS | N(species) | N(species in OBIS) | % species in OBIS |
|---|---|---|---|---|---|---|
| Benthos | 131,390 | 62,316 | 47% | 15,161 | 13,288 | 88% |
| Zooplankton | 4,768 | 2,117 | 44% | 980 | 914 | 93% |
| Nekton | 2,506 | 1,355 | 54% | 893 | 803 | 90% |
| Fish | 8,683 | 5,842 | 67% | 9,578 | 9,123 | 95% |
| Mammals | 85 | 37 | 44% | 107 | 103 | 96% |
| Birds | 238 | 108 | 45% | 449 | 395 | 88% |

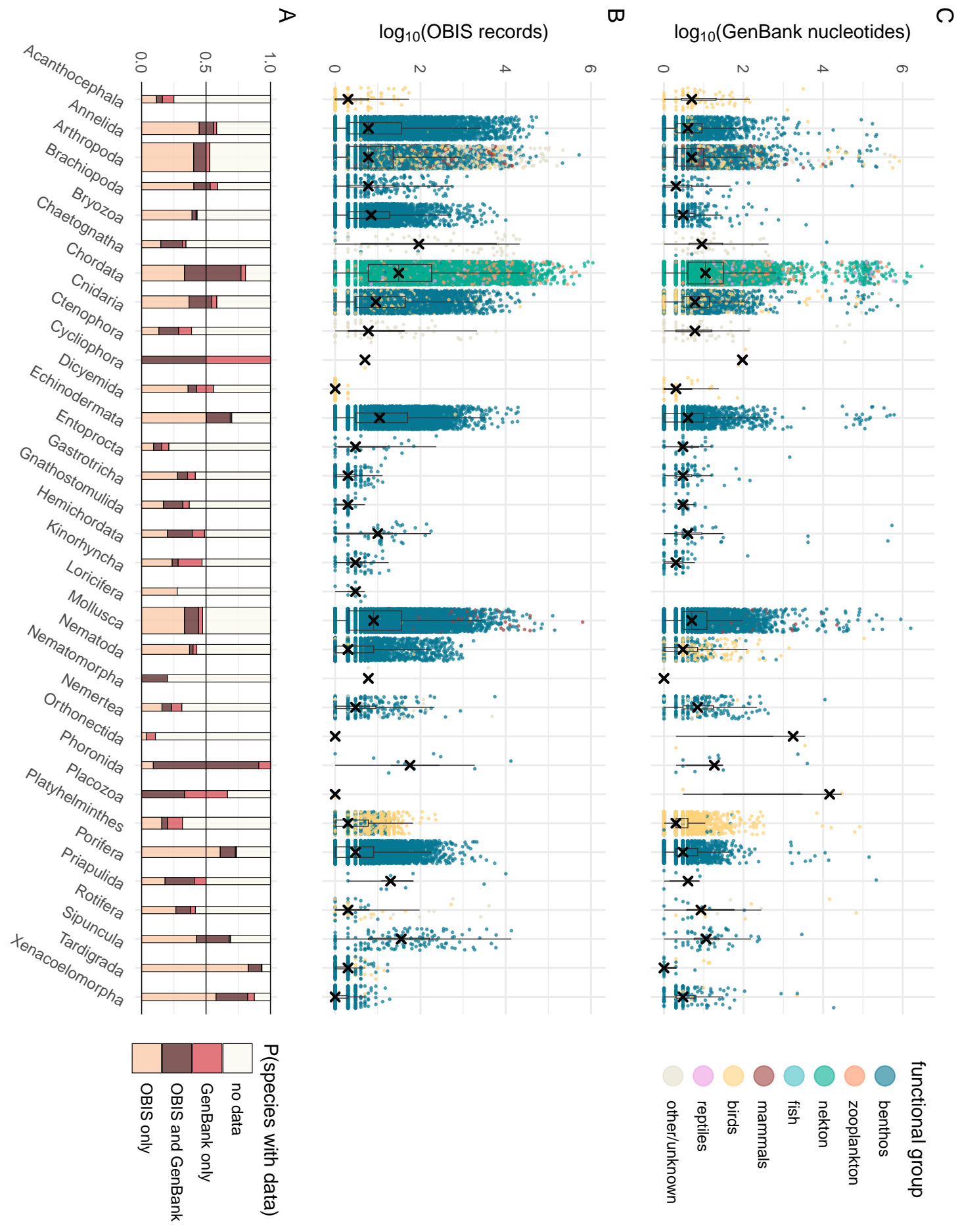| | | | | | | |
|---|---|---|---|---|---|---|
| Reptiles | 80 | 59 | 70% | 16 | 15 | 94% |
| Other / Unknown | 30,292 | 8,692 | 29% | 1,623 | 1,049 | 65% |
| *Totals* | *178,042* | *80,523* | *45%* | *28,807* | *25,690* | *89%* |

Figure Legends

**Figure 1.** Availability of biogeographic (>45M OBIS occurrence records) and genetic (>56M GenBank nucleotides) data across 206,849 marine animal species, summarised by phylum and by broad functional group. **(A)** Proportion of species in each phylum with data in either database, both databases, or neither. Bar width is proportional to the number of species in each phylum. Number of **(B)** OBIS occurrence records and **(C)** Genbank nucleotide sequences are shown for species that occur in the respective database. Each point represents a species, coloured by functional group. Box plots are superimposed with X marking the median number of records within each phylum. Phylum size varies from 2 species (Cycliophora) to 57,336 species (Arthropoda).

**Figure 2.** Coefficients from the hurdle models of data availability across functional groups, first modelling presence in a database with a binomial model, and then non-zero counts of records in a database as a negative binomial model. Species presence in OBIS **(A)** or GenBank nucleotide database **(C)** across functional groups is indicated with binomial coefficients (with 95% confidence intervals) on the response scale, representing the ratio of the probabilities of species within a group having records in the database versus not having records in the database. For the subset of species present in **(B)** OBIS or **(D)** GenBank, the empirical mean number of records per species is plotted together with bootstrapped 95% confidence intervals. For each group, the predicted non-zero count from the hurdle model is indicated with an X. Point size is scaled to the total number of species in each functional group (A, C, ranging from 96 reptiles to 146,551 benthos) and to the number of species in each group with records in OBIS (B, 71 reptiles to 75,604 benthos) or GenBank (D, 78 reptiles to 19,235 benthos).

**Figure 3.** Mosaic plot showing the joint distribution of species between categories of OBIS records and GenBank nucleotides. (A) shows all species, and is dominated by species with no records in either database. (B) zooms in on species with high numbers (>100) of OBIS records, and (C) reverses the axes and zooms in on species with high numbers (>100) of GenBank nucleotides. Axis labels indicate the number of records at the right-hand bound of each category.
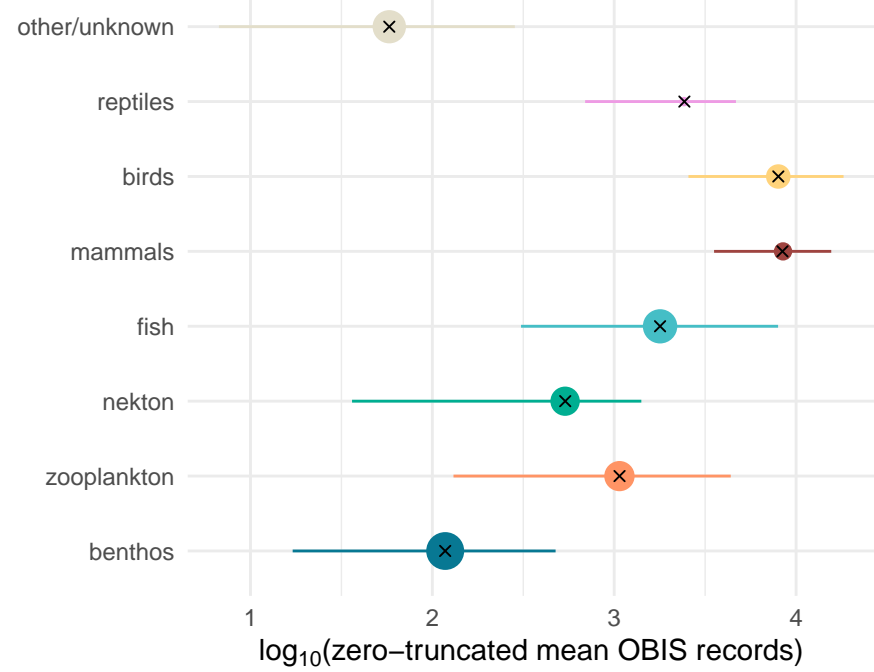
**Figure 4.** Distribution of occurrence records across 106,213 marine animal species present in OBIS by functional group and by **(A)** IUCN assessment status and **(B)** presence in the Barcode of Life Data System. Each point represents a species.
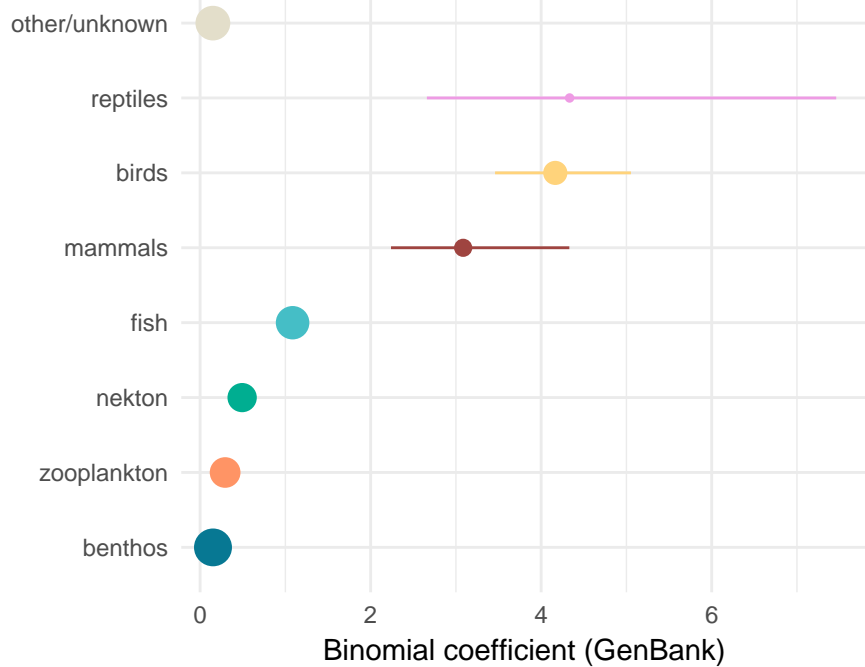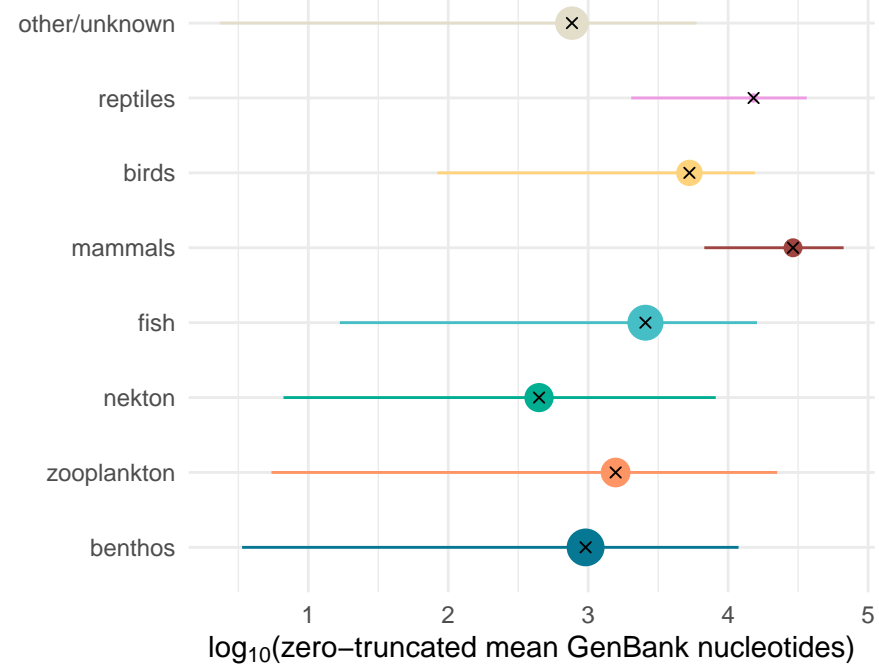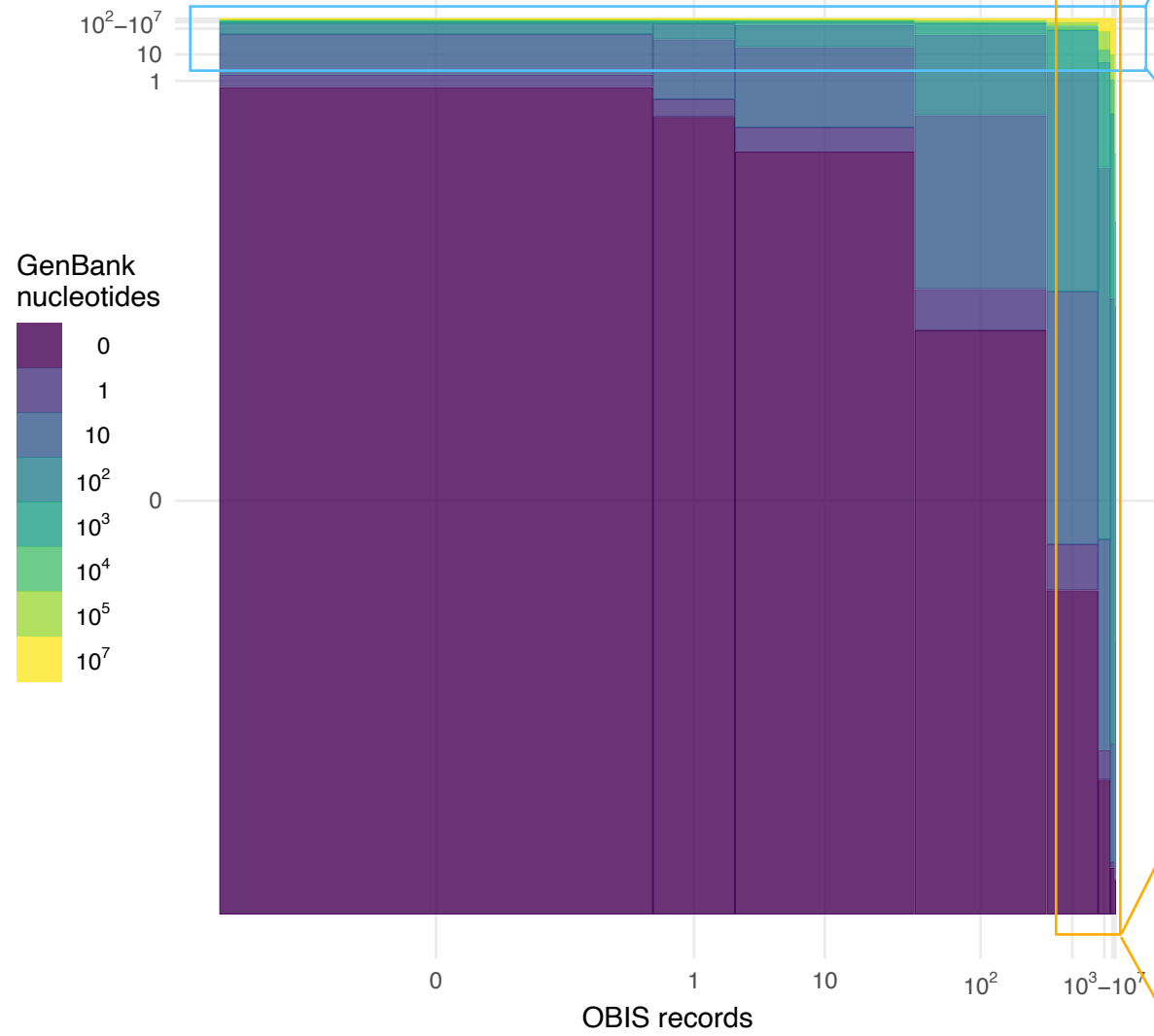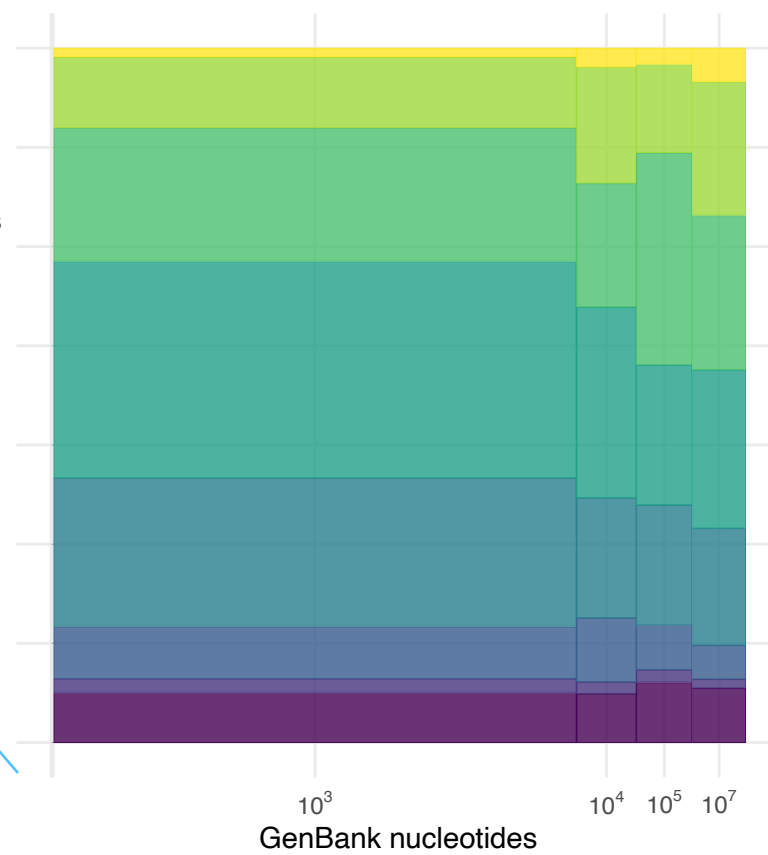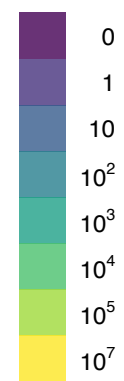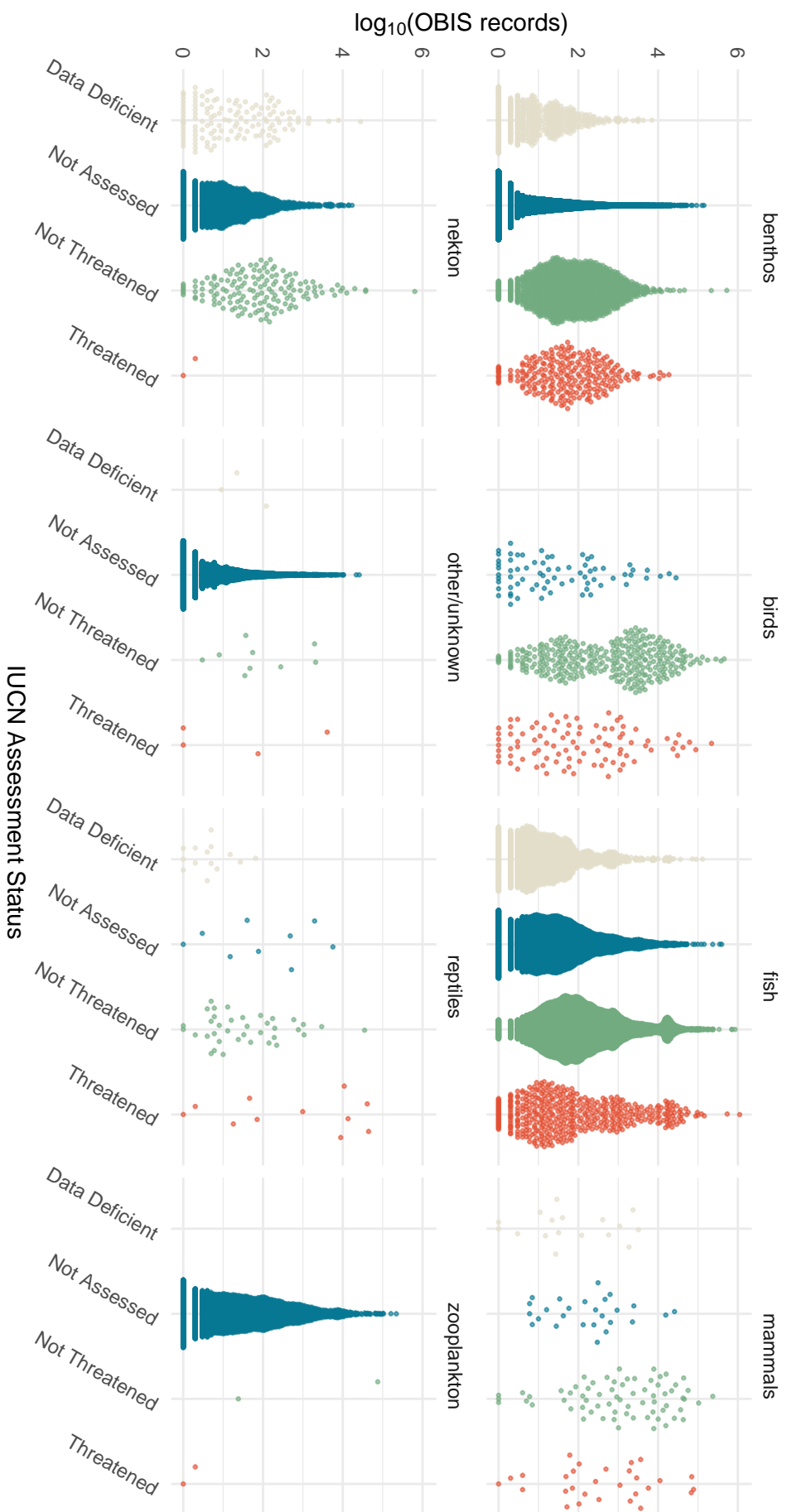
**A**

P(species with data)

0.0    0.5    1.0

**B**

$\log_{10}$(OBIS records)

0    2    4    6

**C**

$\log_{10}$(GenBank nucleotides)

0    2    4    6

Acanthocephala
Annelida
Arthropoda
Brachiopoda
Bryozoa
Chaetognatha
Chordata
Cnidaria
Ctenophora
Cycliophora
Dicyemida
Echinodermata
Entoprocta
Gastrotricha
Gnathostomulida
Hemichordata
Kinorhyncha
Loricifera
Mollusca
Nematoda
Nematomorpha
Nemertea
Orthonectida
Phoronida
Placozoa
Platyhelminthes
Porifera
Priapulida
Rotifera
Sipuncula
Tardigrada
Xenacoelomorpha

P(species with data)

- OBIS only
- OBIS and GenBank
- GenBank only
- no data

functional group

- other/unknown
- reptiles
- birds
- mammals
- fish
- nekton
- zooplankton
- benthos