

1 **Title:** FAIR enough? A perspective on the status of nucleotide sequence data and
2 metadata on public archives

3

4 **Authors:** Christiane Hassenrück^{1,2}, Tobias Poprick³, Véronique Helfer³, Massimiliano Molari⁴,
5 Raissa Meyer⁵, Ivaylo Kostadinov⁶

6

7 **Affiliations:**

8 ¹ Leibniz Institute for Baltic Sea Research Warnemünde (IOW), Seestrasse 15, 18119 Rostock-
9 Warnemünde, Germany

10 ² MARUM - Center for Marine Environmental Sciences, University of Bremen, Leobener Straße 8, 28359
11 Bremen, Germany

12 ³ Leibniz Centre for Tropical Marine Research (ZMT), Fahrenheitstrasse 6, 28359 Bremen, Germany

13 ⁴ Max Planck Institute for Marine Microbiology, Celsiusstr. 1, 28359 Bremen, Germany

14 ⁵ Alfred Wegener Institute for Polar and Marine Research, Am Handelshafen 12, 27570 Bremerhaven,
15 Germany

16 ⁶ GFBio - Gesellschaft für Biologische Daten e.V., c/o Research II, Campus Ring 1, 28759 Bremen,
17 Germany

18

19 **Correspondence:** christiane.hassenrueck@io-warnemuende.de

20

21 **Abstract:**

22 Knowledge derived from nucleotide sequence data is increasing in importance in the life sciences,
23 as well as decision making (mainly in biodiversity policy). Metadata standards have been
24 established to facilitate sustainable sequence data management according to the FAIR principles
25 (Findability, Accessibility, Interoperability, Reusability). Here, we review the status of metadata
26 available for raw read Illumina amplicon and whole genome shotgun sequencing data derived
27 from ecological metagenomic material that are accessible at the European Nucleotide Archive
28 (ENA), as well as the compliance of the primary sequence data (fastq files) with data submission
29 requirements. While overall basic metadata, such as geographic coordinates, were retrievable in
30 98% of the cases for this type of sequence data, interoperability was not always ensured and
31 other (mainly conditionally) mandatory parameters were often not provided at all. Metadata
32 standards, such as the 'Minimum Information about any(x) Sequence (MIxS)', were only
33 infrequently used despite a demonstrated positive impact on metadata quality. Furthermore, the
34 sequence data itself did not meet the prescribed requirements in 31 out of 39 studies that were
35 manually inspected. To tackle the most immediate needs to improve FAIR sequence data

36 management, we provide a list of minimal suggestions to researchers, research institutions,
37 funding agencies, reviewers, publishers, and databases, that we believe might have a potentially
38 large positive impact on sequence data and metadata FAIRness, which is crucial for further
39 research and its derived applications.

40

41 **Keywords:** sequence data management, metadata standards, interoperability, reusability,
42 ontology, digital sequence information (DSI), biodiversity, next generation sequencing

43

44 INTRODUCTION

45

46 Next generation sequencing has gained increasing popularity and is now firmly established as a
47 routine tool in multiple fields of the life sciences, such as ecology (foremost microbial ecology),
48 biodiversity research, and conservation biology. Furthermore, knowledge derived from nucleotide
49 sequence data, also referred to as digital sequence information (DSI), is becoming increasingly
50 relevant for decision-making in natural resource management (e.g. [Sustainable Development](#)
51 [Goals](#)) and as part of international agreements (e.g. [Convention on Biological Diversity](#)). The
52 amount of nucleotide sequence data has been and still is growing exponentially (Harrison et al.,
53 2021). However, a string of nucleotides (ACTG) on its own does not contain much information -
54 metadata and contextual parameters (Box 1) are required to describe sample origin and sequence
55 generation for the data to be meaningful within and beyond the scope of the study, for which the
56 sequence was obtained. Capturing and communicating not only the primary (sequence) data, but
57 also its metadata and contextual data, is a crucial part of good data management. To promote
58 sustainable data management and usage, the FAIR principles have been introduced. They offer
59 guidance on how to make data Findable, Accessible, Interoperable, and Reusable (Box 1;
60 Fillinger, de la Garza, Peltzer, Kohlbacher, & Nahnsen, 2019; Wilkinson et al., 2016), to prepare
61 for a future of more automated analyses with the aim that the value of the data will not be restricted
62 to a single study, but will extend to the reuse and integration across multiple studies over time.
63 Recently, the trend of an increasing data volume, which is being more and more sustainably
64 managed, has resulted in nucleotide sequence data being used more frequently in the emerging
65 field of data science, answering new scientific questions with existing data, as such constituting
66 a public good for the scientific community (Box 1). One prime example is the TARA Oceans data
67 set, which has so far resulted in hundreds of publications making secondary use of the data – a
68 number that is constantly increasing¹.

69

70 One key aspect of FAIR data is the implementation of standards for metadata (Box 1; Wilkinson
71 et al., 2016). To this end, the Genomic Standards Consortium (GSC; Field et al., 2011) has
72 established the 'Minimum Information about any (x) Sequence ([MIxS](#))' family of standards, to
73 describe sample collection and sequence generation in a consistent manner (Yilmaz et al., 2011;
74 Fig 1). MIxS consists of several customized checklists tailored for a diverse set of sequencing
75 applications and investigated environments, defining mandatory (always to be provided, core

¹ https://oceans.taraexpeditions.org/wp-content/uploads/2020/10/TARA_RA_EN_.pdf

76 parameters), conditionally mandatory (mandatory for a specific sequencing application),
77 environment-specific, and optional parameters. In addition to standardizing metadata parameter
78 names, MIxS suggests a consistent format for units and syntax for data values, thereby promoting
79 the interoperability of the data provided in compliance with this standard. For instance, MIxS
80 makes use of ontologies, such as the Environmental Ontology ([ENVO](#); Buttigieg, Morrison, Smith,
81 Mungall, & Lewis, 2013; Buttigieg et al., 2016) or the Experimental Factor Ontology ([EFO](#); Malone
82 et al., 2010), to describe the sampled environment or experimental conditions using a controlled
83 vocabulary and to make the metadata more machine-actionable.

84

85 Box 1: Glossary

FAIR (paraphrased after Wilkinson et al. (2016):

- Findable: Data and metadata are linked and findable via a unique and persistent identifier (e.g. accession number). Metadata is further searchable.
- Accessible: Data and metadata are retrievable (by humans and machines) via their identifier. Metadata remains accessible even if associated data is not available anymore.
- Interoperable: Data and metadata use a common language for knowledge representation understandable by humans and machines.
- Reusable: Data are described by rich metadata to provide the context required for reuse.

Ontology: Ontologies impose a (machine-readable) hierarchical structure of relationships for the components of a given system, using a controlled and clearly defined vocabulary. In the case of the Environmental Ontology (ENVO), this increases the interoperability of environmental descriptions, helping (meta)data records achieve demonstrable FAIRness.

Metadata: Collection of parameters that describe the primary data, in this example nucleotide sequencing data. Most metadata parameters are intrinsic to the sampling or experimental design and the laboratory or analytical procedures. As such they are often known *a priori*, i.e. before the primary data collection. For instance, sampling location, experimental treatments, sequencing platform.

Contextual data: While often grouped together with metadata, contextual data is referring to parameters which are recorded alongside the primary data. For instance, temperature, salinity, inorganic nutrient concentrations. They are often primary data for other research fields.

Metadata standards: Metadata standards provide a structured framework for metadata documentation in compliance with the FAIR principles. They provide checklists of well-defined parameters to be reported and determine the vocabulary and units to be used to ensure consistent data across studies.

Data mining: In the context of this publication, data mining is referring to the retrieval and reuse of data sets that have been archived in open access data repositories.

86

87

88 The main hubs for long-term nucleotide sequence data storage are the three resources making
89 up the International Nucleotide Sequence Database Collaboration ([INSDC](#); Fig 1): the European
90 Nucleotide Archive ([ENA](#)), the National Center for Biotechnology Information ([NCBI](#)), and the

91 DNA Data Bank of Japan ([DDBJ](#)). These databases are mirrored so that the same data is
92 available in all three. The establishment of this infrastructure has been instrumental in propagating
93 global standards for sequence data and metadata (e.g. fasta, and fastq data formats) and offers
94 services far beyond the provision of access to such data (see e.g. Cook, Bergman, Cochrane,
95 Apweiler, & Birney, 2018; Fukuda, Kodama, Mashima, Fujisawa, & Ogasawara, 2021; NCBI
96 Resource Coordinators, 2018). Furthermore, as part of the GSC, the members of the INSDC have
97 contributed to designing the MixS standards, pioneering the implementation of such metadata
98 standards in nucleotide sequence databases with their official adoption in 2011 (Yilmaz et al.,
99 2011).

100

101 ENA offers a [minimal metadata standard](#) for sequence submissions, although the use of more
102 extensive [checklists](#), such as those based on MixS, is strongly recommended². Metadata on ENA
103 is organized on several levels³ (Fig. 1): Sequencing **runs** are associated with specific
104 **experiments**, which are referring to individual nucleic acid extractions and/or library preparations.
105 Experiments are collected into **studies**, which usually use a common methodological approach.
106 Several studies can be summarized by an umbrella project that may correspond to the larger
107 scientific project, for which the data was generated. **Samples**, referring to the biological material,
108 can be associated with multiple studies through different experiments. This flexible metadata
109 model allows representing complex experimental set-ups correctly, but can be hard for
110 inexperienced submitters to navigate properly and provide the necessary information for checklist
111 compliance.

112

113 When accessing sequence data as data consumer (Fig. 1), all provided metadata for each level
114 (run, experiment, study, sample) can be retrieved as XML. To simplify metadata access, the [ENA](#)
115 [advanced search](#) offers a collection of indexed parameters that use standardized names and are
116 searchable (i.e. usable to restrict the search) and returnable (i.e. downloadable in a user-friendly
117 TSV format; [ENA Portal API](#); Fig. 1). On the run level, these indexed parameters are also inherited
118 from sample and experiment metadata. At the moment, the implementation of indexed
119 parameters is limited to metadata parameters that are mandatory for most checklists and/or most
120 frequently provided. As such, many conditionally mandatory, environment-specific, and optional
121 MixS parameters, mainly due to a lack of consistent and widespread use, are not indexed and
122 only accessible in XML format, where no standardized nomenclature, controlled vocabulary or

² <https://ena-docs.readthedocs.io/en/latest/submit/samples.html#checklists>

³ <https://ena-docs.readthedocs.io/en/latest/submit/reads/programmatic.html#object-relationships>

123 specific data value format are enforced. Therefore, some of the value of MIxS is intrinsically lost,
124 making non-indexed parameters not interoperable.

125

126 In addition to metadata requirements, ENA also standardizes the format of the submitted
127 sequence data depending on the sequencing approach. For instance, paired-end Illumina raw
128 reads have to be submitted as demultiplexed R1 and R2 files (fastq) without artificial sequences
129 (e.g. adapters, linkers, barcodes/tags, primers) and prior to any quality trimming⁴. To provide
130 sequencing data in such a format, initial sequence processing steps are necessary, starting from
131 the multiplexed sequencer output. As bioinformatic sequence analysis pipelines vary, it is
132 important to not deviate from the sequence format requirements by adjusting analysis workflows
133 accordingly.

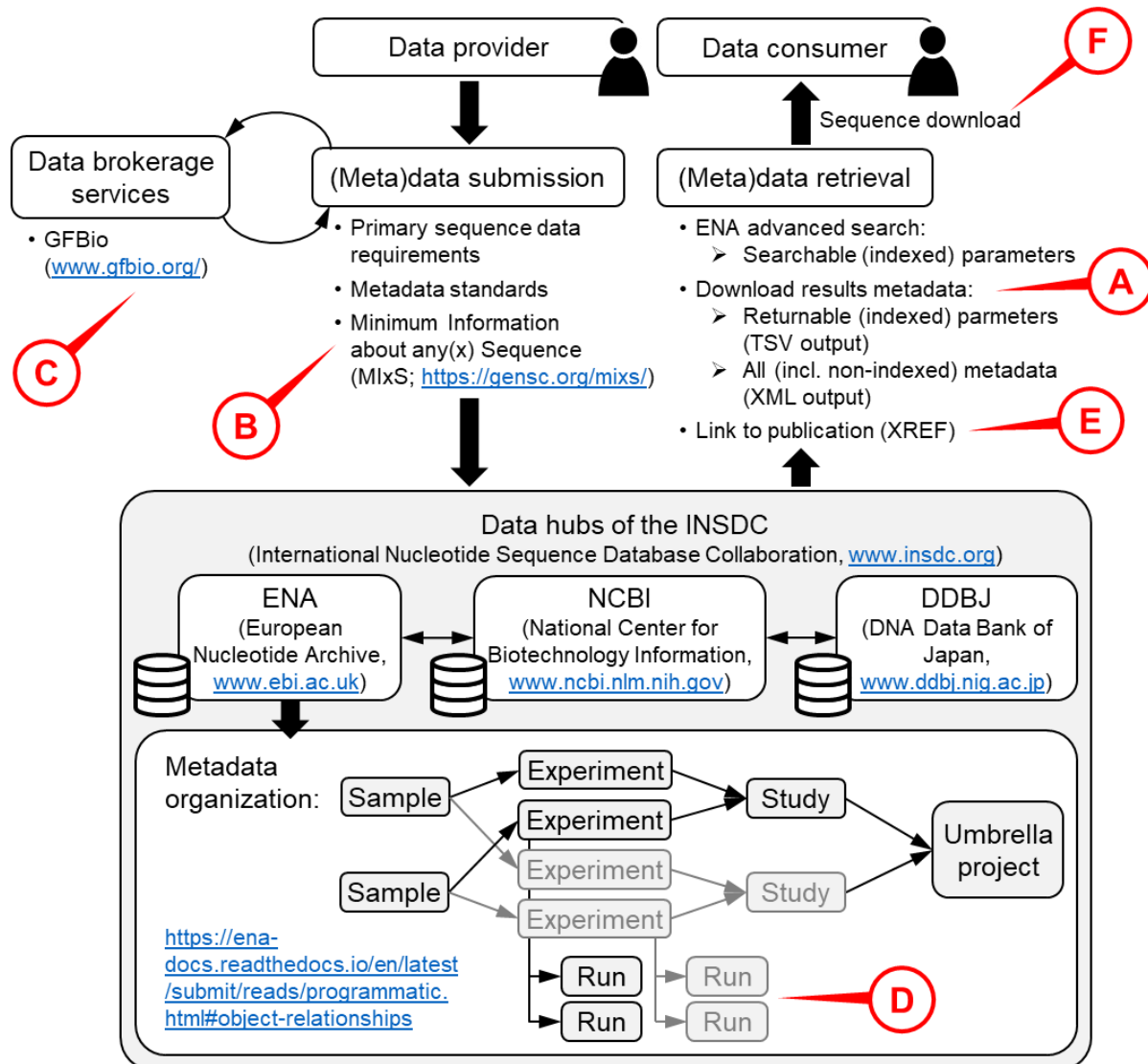
134

135 To support sustainable data management, data brokerage services, such as the German
136 Federation for Biological Data ([GFBio](#)), have been established (Diepenbroek et al., 2014).
137 Brokerage services offer a central entry point for data submissions, providing personal guidance
138 (helpdesk) on FAIR data, supporting and often simplifying the data submission process, and
139 ensuring data deposition on the most appropriate archive. As an additional checkpoint, brokerage
140 services therefore constitute a valuable resource for each individual researcher to improve the
141 FAIRness of their data, which is now becoming a strict requirement from many funding agencies.

142

143

⁴ <https://ena-docs.readthedocs.io/en/latest/submit/fileprep/reads.html?highlight=read%20format#fastq-format>



- A: Which metadata is available and how accessible is it?
Investigated parameters include geographic coordinates, environment description according to ENVO, target gene, nominal length, status as environmental sample.
- B: How does use of MIxS checklists according to environmental packages affect (A)?
- C: How does use of a brokerage service (GFBio) affect (A)?
- D: Are sequences deposited in the correct format?
- E: Can additional information be easily retrieved by linking the publication?
- F: The analysis is sequence-centric, focusing on the metadata available per run.

144

145 Figure 1: Summary of the submission and retrieval process for nucleotide sequence data using the
146 databases of the INSDC. The letters highlight the specific aspects in that process that were investigated
147 here, focusing on how various factors (B, C) affect metadata quality (A) in a sequence-centric data mining
148 approach (F), the compliance of the primary sequence data with data archiving requirements (D), and the
149 retrieval of further information via associated publications (E).

150

151 To facilitate data reuse in data mining endeavors, access to and retrieval of the raw read
152 sequencing data and metadata on run level, together with the inherited metadata parameters

153 describing the sample and sequencing experiment, are most crucial. However, despite the
154 available framework for FAIR data archiving, data interoperability and reusability are still limited
155 and often complicated by insufficient metadata (Eckert et al., 2020; Hoopen et al., 2016; Jurburg,
156 Konzack, Eisenhauer, & Heintz-Buschart, 2020; personal observation). Therefore, we decided to
157 conduct a review of the status of nucleotide sequence data and associated metadata accessible
158 through ENA to (i) identify deficits in metadata quality and (ii) provide suggestions for improving
159 FAIR data management (Fig. 1).

160

161 We restricted our analysis to a very popular example for biodiversity assessment in ecology:
162 paired-end amplicon (metabarcoding) raw read data generated from ecological metagenomes
163 (NCBI taxid: 410657) as source material on the Illumina platform. We focus on metadata
164 parameters, which are mandatory and/or crucial for the reuse of this kind of data, evaluating the
165 impact that the use of MIxS checklists (Fig. 1B) and the brokerage service GFBio (Fig. 1C) had
166 on metadata quality. We searched ENA on 13.12.2020 for raw read data using the following
167 search query: `tax_tree(410657) AND library_selection = "PCR" AND`
168 `library_strategy = "AMPLICON" AND library_layout = "PAIRED" AND`
169 `instrument_platform = "ILLUMINA" AND library_source = "METAGENOMIC".` For
170 the resulting 413 849 search results on run level (Fig. 1F; hereafter referred to as cases), we
171 downloaded all available metadata parameters in TSV format as well as the sample and
172 experiment XML. Specifically, we checked for the following parameters if data was provided, ease
173 of access, correctness (if applicable), and compliance with standards (Fig. 1A): (i) geographic
174 coordinates, i.e. **latitude** and **longitude**, (ii) information about the **target gene** or subfragment or
175 primers, (iii) **nominal length**, i.e. insert size, and (iv) the parameters **environment_biome**,
176 **environment_material**, and **environment_feature**, which make use of ENVO to ensure a
177 standardized description of the sampled environment using a controlled vocabulary. As
178 comparison to the amplicon sequencing example, we repeated this assessment also for shotgun
179 metagenomic paired-end Illumina reads using the same query apart from the following changes:
180 `library_selection = "RANDOM" AND library_strategy = "WGS"` (Whole Genome
181 Sequencing; date accessed 08.01.2021; SI figures 1 + 2). The scripts to search ENA, download
182 the metadata, and calculate the summaries presented in this study are available on:
183 <https://github.com/chassenr/ENA-metadata>. The summaries of cases per year visualized in the
184 subsequent figures are further available in SI table 1.

185

186 Apart from metadata quality, the reusability of nucleotide sequence data - especially the
187 automation thereof - strongly depends of the format of the submitted raw read data itself.
188 Therefore, using a recent data mining effort focused on amplicon studies of the V3-V4
189 hypervariable region of the bacterial 16S gene (Molari et al. in prep.) as a **case study**, we
190 evaluated if the raw read data (fastq files) had been archived according to the ENA guidelines
191 (Fig. 1D). Furthermore, we checked the correct declaration as **environmental sample** and the
192 availability of the associated **manuscript publication** comparing a manual search to the [ENA](#)
193 [XREF API](#) (Fig. 1E), as these may provide further relevant information about the provenance of
194 the data.

195 RESULTS

196

197 **General trends:** The number of cases (runs) retrieved by the above-mentioned query has been
198 increasing continuously over the last decade, with more than 120 000 runs submitted in 2020
199 alone. In total, only 6.5% of runs (26 903 of 413 849) were linked to samples compliant with a
200 MlxS checklist and environmental package, with 21% in 2015 and (with the exception of 2018)
201 decreasing proportions since. On average 1.22 runs were submitted per sample, with a highly
202 skewed distribution where only 8% of the samples were associated with more than one run. As
203 such, most cases corresponded to a single sample, from which several of the investigated
204 metadata parameters were inherited.

205

206 **Geographic coordinates** (Fig. 2 top): Latitude and longitude are mandatory MlxS parameters
207 and essential for the reuse of sequencing data from environmental samples, especially in
208 molecular ecology. However, these parameters are not necessarily enforced by ENA for
209 submissions that do not use MlxS. Nevertheless, in our example the majority of all cases (80%)
210 were archived with latitude and longitude available as decimal degrees in the TSV search output.
211 For an additional 18% of the cases, latitude and longitude values were retrievable from the sample
212 XML as part of non-indexed parameters. There, this data was stored under [23 different parameter](#)
213 [names](#), and was therefore not easily accessible or interoperable. Considering only cases
214 submitted according to a MlxS checklist and specifying a MlxS environmental package, latitude
215 and longitude were always provided in some form and the proportion of cases with latitude and
216 longitude available in the TSV output was slightly higher across all years (86%), although it has
217 been declining from more than 99% to 61% since 2017.

218

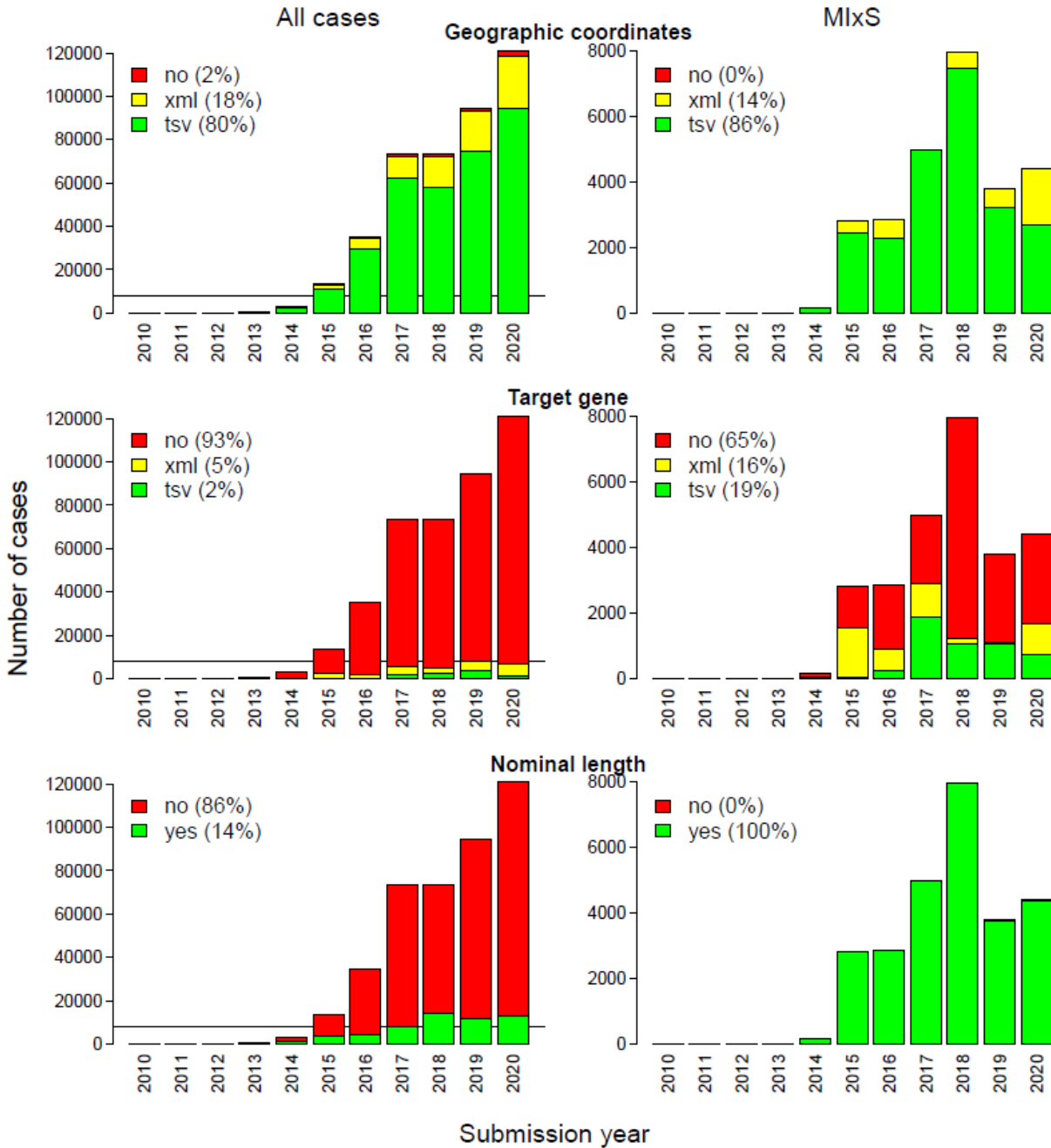
219 **Target gene** (Fig. 2 middle): Missing information about which DNA region was targeted by the
220 amplicon sequencing approach is one of the main obstacles for data interpretation and reuse.
221 The ENA search output in TSV format includes the indexed parameter `target_gene`, which can be
222 used to specify the amplified gene or locus name using free text. In MlxS, `target_gene` is further
223 specified as a mandatory parameter for amplicon sequencing studies (MIMARKS survey),
224 although its use is not enforced by ENA for data submissions. Additionally, the non-indexed
225 parameters `target_subfragment` and `pcr_primers` are listed as conditionally mandatory
226 parameters to supply additional metadata about the amplified gene region, and as such, should
227 be supplied for all amplicon sequencing experiments. Among all cases investigated here, only 2%
228 provided the target gene, with an additional 5% where some information about the amplified

229 region could be retrieved from non-indexed parameters in the sample XML (stored under [58](#)
230 [different parameter names](#)) and from the library construction protocol included in the experiment
231 XML. However, such entries were extremely inconsistent, ranging from gene and gene region
232 names (or a combination of both) to primer names, primer sequences, and references for the
233 applied PCR protocol. To reuse this data, each entry would have to be inspected manually,
234 making available target gene information not only difficult to access, but also not interoperable.
235 The proportion of cases with target gene information available among those submitted according
236 to a MIxS checklist and environmental package was considerably higher with 35%, although still
237 far from optimal bearing in mind that this is a mandatory parameter. The correct identification of
238 the amplified region without any respective metadata is cumbersome and computationally
239 expensive. Complete and correct metadata entries, using a standardized format or even
240 controlled vocabulary preferably in accordance with existing ontologies, would drastically reduce
241 computational and man-power requirements for post-deposition data curation and data reuse.

242

243 **Nominal length** (Fig. 2 bottom): Nominal length specifies the insert size, i.e. the length of the
244 amplified fragment between the sequencing adapters (i.e. including primers) in the library. It is
245 mandatory for all paired-end sequencing runs according to ENA, NCBI, DDBJ submission
246 tutorials, and should have been enforced since 2014. However, in 86% of all cases this parameter
247 was not provided. The use of a MIxS checklist and environmental package during the submission
248 increased the percentage of cases providing nominal length to almost 100%, with only 55 cases
249 in total in 2019 and 2020 lacking these values, suggesting that submitters who made the effort to
250 be MIxS compliant were also more likely to provide metadata parameters outside of MIxS.
251 Interestingly, peaks in the distribution of supplied nominal length values occurred at 250bp,
252 300bp, 500bp, and 600bp (data not shown), which correspond to the length of individual reads or
253 the combined length of paired reads in the popular 2x250bp and 2x300bp sequencing
254 approaches, and may therefore represent read length rather than insert size. This suggests that
255 misconceptions exist about the definition of the parameter nominal length among sequence data
256 submitters.

257



258

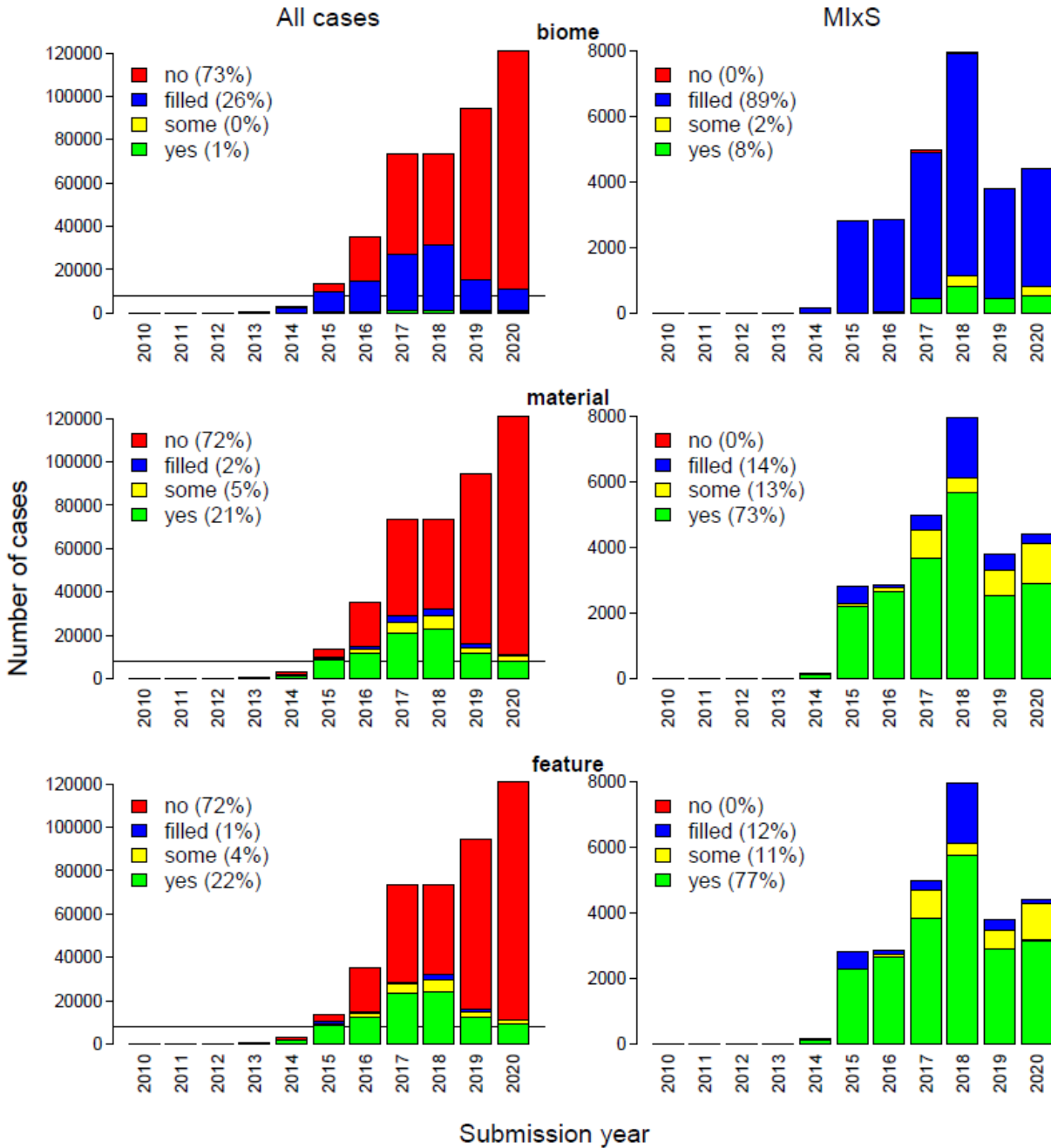
259 **Figure 2:** Number of cases (amplicon example) with and without metadata available for geographic
 260 coordinates (latitude, longitude), target gene (or related information, such as subfragment or pcr primers),
 261 and nominal length. For geographic coordinates and target gene, 'tsv' is referring to the information
 262 provided for indexed metadata parameters in the TSV search output, while 'xml' is referring to non-indexed
 263 metadata only accessible in the XML view of the ENA sample or experiment. The percentages summarize
 264 cases over all years and are rounded to integers. The plots on the right only show the cases with samples
 265 submitted according to a MixS checklist and environmental package. The horizontal line in the plots on the
 266 left indicates the y-axes range of the plots on the right.

267

268 **Environment description using ENVO** (Fig. 3): The parameters `environment_biome`,
269 `environment_material`, and `environment_feature` use ENVO terms to characterize various
270 characteristics of the sample and the environment it originated from. They are mandatory
271 parameters for any MlxS checklist. Across all cases 72-73% did not provide data for any of these
272 three parameters. Conversely, values were supplied for all cases, which specified a MlxS
273 checklist and environmental package, with the exception of 125 cases for the parameter
274 `environment_biome`. However, ENVO terminology was inconsistently used despite the format
275 required by MlxS. Especially for the parameter `environment_biome`, the majority of the provided
276 values did not resemble any existing or even obsolete ENVO terms. Exact matches to ENVO term
277 IDs, which are mandatory to be included according to the MlxS documentation, were only found
278 for 8% of the cases, with an additional 2% with character string matches to ENVO term names.
279 For `environment_material` and `environment_feature`, these proportions were considerably higher
280 with 73% and 77% matches to ENVO term IDs and an additional 13% and 11% matches to ENVO
281 term names, respectively. This demonstrates that the use of MlxS checklists drastically improved
282 the availability of an environment description via the parameters `environment_biome`,
283 `environment_material`, and `environment_feature`, but also that, if provided, the interoperability of
284 this metadata is severely impaired by non-ENVO entries.

285
286 Interestingly, a considerable number of cases, none of which specified a MlxS checklist upon
287 submission, provided metadata for the parameters `env_broad_scale`, `env_local_scale`,
288 `env_medium`. Those parameters are included in a more recent version of MlxS (version 5.0)
289 instead of `environment_biome`, `environment_material`, and `environment_feature` (version 4.0). As
290 the new parameters are not (yet) indexed on ENA, we retrieved the values from the sample XML.
291 Since 2018, 47 302 cases consistently provided values for all three parameters, corresponding to
292 16% of all cases submitted during that time period, and following an increasing trend with 27% of
293 all cases in 2020. This trend counteracted the decreasing use of `environment_biome`,
294 `environment_material`, and `environment_feature` over the investigated time period, resulting in
295 stable proportions of approximately 35% of the cases per year archived, regardless of MlxS
296 checklist usage, with either of these two sets of environmental descriptors since 2016. Lack of
297 compliance with ENVO terminology was also observed here, especially for the parameter
298 `env_broad_scale` (data not shown).

299
300



301
302
303
304
305
306
307
308
309
310

Figure 3: Number of cases (amplicon example) with values for environment_biome, environment_material, and environment_feature according to ENVO. Yes: exact match to ENVO term ID; some: character string match to ENVO term name (including matches after extensive character string manipulation); filled: a value is provided, but not ENVO-compatible; no: no entry. The percentages summarize cases over all years and are rounded to integers. The plots on the right only show the cases with samples submitted according to a MixS checklist and environmental package. The horizontal line in the plots on the left indicates the y-axes range of the plots on the right.

311 **Use of brokerage services:** Since 2016, 1475 cases have been submitted via GFBio, making
312 this brokerage service the most frequently used (as assessed by the number of cases) throughout
313 the whole investigated time period, closely followed by CNSA (China Nucleotide Sequence
314 Archive). Here, we briefly highlight the effect the use of this brokerage service had on the
315 metadata supplied for the above parameters. We found that the quality of the metadata,
316 specifically accessibility and interoperability, were considerably improved in submissions via
317 GFBio (SI table 1): all cases (with the exception of 95 cases from one study) provided latitude
318 and longitude data retrievable from the TSV search output of ENA and also used the correct
319 format and terminology for the parameters `environment_biome`, `environment_material`, and
320 `environment_feature`. Furthermore, all cases (without exception) contained nominal length data,
321 although the suspicious peaks in the data distribution at 300bp and 500bp were still present. If
322 information about the amplified gene was provided (26% of the cases), this was mainly accessible
323 in the TSV search output of ENA (24%), and restricted to two values: 16S rRNA and 18S rRNA.

324

325 **Metagenomic sequencing data:** Until the end of 2020, 29 253 cases had been submitted with
326 `library_selection` "RANDOM" and `library_strategy` "WGS" in an analogous example to the
327 amplicon data sets that we explored above. The number of cases submitted each year has been
328 consistently increasing, with a spike in 2019 caused by one study with an exceptionally high
329 number of associated runs. Nonetheless, the proportion of WGS of the total number of amplicon
330 and WGS cases in our particular example has been declining over the last decade, reaching a
331 mostly stable value at approximately 5% since 2017 (SI table 1). Regarding the investigated
332 metadata parameters, similar trends were observed in the WGS data compared to the amplicon
333 example, although overall metadata quality was slightly higher (SI figures 1 and 2). Specifically,
334 this improvement was related to the more frequent usage of MIxS checklists for in total 14% of all
335 WGS cases (SI table 1). These cases associated with samples submitted in compliance with MIxS
336 displayed an almost perfect track record for the parameters latitude and longitude, nominal length,
337 `environment_feature` and `environment_material` in terms of metadata interoperability and
338 consequently reusability.

339

340 **Case study:** In a recent study (Molari et al. in prep.), the raw reads archived on ENA from 39
341 studies using paired-end Illumina amplicon sequencing of the V3-V4 hypervariable region of the
342 bacterial 16S rRNA gene were downloaded and bioinformatically processed to generate quality-
343 trimmed merged fasta files to be used for oligotyping (Eren et al., 2014). This analysis required
344 that the sequences were generated from the exact same gene region to be comparable.

345 Information about the primers used in the sequencing library preparation were obtained from
346 associated publications or the submitters directly. After download, the raw read data was checked
347 for compliance with ENA submission requirements, i.e. that paired-end Illumina reads were
348 archived as demultiplexed, unmerged forward and reverse reads, without artificial sequences and
349 prior to any quality trimming. Of the 39 inspected studies, only eight were submitted as required.
350 The majority (28 studies) did not remove the primer sequences, eight studies contained already
351 merged sequences, and one study even provided only the sequencer output prior to
352 demultiplexing (sample barcode information had to be obtained from the author). This data mining
353 experience showed that even if metadata that enables the findability and accessibility of the data
354 is provided, the raw read data itself may often not be submitted as required, making manual
355 checks mandatory and limiting the interoperability and reusability of the data.

356
357 As we investigated each of these studies in detail, comparing the information provided in the
358 sample and study description as well as the associated publication (if available), we also checked
359 the values for the logical (boolean) sample metadata parameter `environmental_sample`. This
360 parameter “identifies sequences derived by direct molecular isolation from an environmental DNA
361 sample” ([ENA Portal API](#)). This description applied to the samples from all 39 studies, however
362 all samples were declared as non-environmental. Based on this observation, we revisited the
363 metadata inspected in the current study: Of the 413 849 cases of amplicon data less than 2%
364 were marked as originating from environmental samples. This seemed unlikely, although we were
365 not able to check all submissions manually for the correctness of this parameter. Specifically, we
366 found it paradoxical that none of the cases submitted according to a MIxS environmental package
367 were actually declared as originating from environmental samples. If our suspicions about the
368 incorrect use of this parameter were confirmed, it would make this metadata parameter unsuited
369 for selecting data for reuse in data mining endeavors.

370
371 Lastly, we assessed the availability of a PubMed record retrievable via the ENA XREF API for the
372 ENA study accessions in the case study, as such a publication may provide further information
373 about a data set than available in the metadata on ENA. Associated publications were retrievable
374 for 20 of the 39 investigated studies. Of the remaining 19, publications were found manually for
375 14. This shows the limitations of the automated approach using XREF. To improve metadata
376 completeness, publications would have to be linked manually to the ENA study upon request by
377 the author.

378

379 DISCUSSION

380

381 All investigated cases met the criteria for metadata findability and accessibility since those were
382 a prerequisite of conducting this study. However, interoperability and therefore reusability
383 remained a challenge. Overall, our results revealed a high variability in data and metadata
384 interoperability and reusability on ENA for the particular examples relevant for molecular ecology.
385 Laudably, with few exceptions, geographic coordinates (latitude and longitude) were always
386 provided, and mostly available as indexed (searchable and returnable) parameters. Beyond the
387 scientific relevance, geographic information about the sample and, by extension, sequence origin
388 is essential for equitable [Access and Benefit Sharing](#) and the key parameter to linking sequence-
389 based biodiversity observations to the Ocean Biodiversity Information System ([OBIS](#)), the Global
390 Biodiversity Information Facility ([GBIF](#)), and other platforms for biodiversity assessment (Bax et
391 al., 2019; Buttigieg et al., 2018; Canonico et al., 2019). However, other mandatory metadata
392 parameters (nominal length) or those often crucial for the reuse of the data (target gene,
393 classification as environmental sample, environment description) showed a very low
394 interoperability, if values were provided at all.

395

396 While the MIxS metadata checklists have been established a decade ago (Yilmaz et al., 2011),
397 they were only infrequently used despite their evident positive impact on metadata quality and
398 consequently the reusability of the data. Furthermore, despite an increasing number of calls to
399 action (Reiser, Harper, Freeling, Han, & Luan, 2018; Ryan et al., 2020; Stevens et al., 2020), the
400 use of this community standard has been declining over the last years, especially for amplicon
401 sequencing data in the selected example. The number of data sets being submitted each year is
402 expected to continue to rise considering decreasing sequencing costs and the undiminished
403 popularity of amplicon sequencing for multiple applications in molecular and microbial ecology
404 and biodiversity research (e.g. environmental DNA studies). Therefore, it is even more worrisome
405 that the use of standards in data submissions has not been following this same upward trend. We
406 further noticed that even when data was submitted according to MIxS, this metadata standard
407 was often not used as intended or to its full potential. In part, this situation may have arisen from
408 inconsistencies in the documentation about metadata requirements provided by separate
409 resources. For instance, the MIxS checklists are not implemented in their entirety by INSDC due
410 to a lack of demand to archive such parameters. Especially conditionally mandatory parameters,
411 such as `target_gene`, `target_subfragment`, and `pcr_primers` in the case of amplicon data (i.e.
412 MIMARKS standard), are listed as optional for the MIxS checklists on ENA, the latter two also not

413 being indexed as searchable or returnable parameters. Furthermore, the description of the
414 parameters `environment_biome`, `environment_material`, `environment_feature` in the ENA
415 documentation specifies free text, whereas MIxS specifies the use of a controlled vocabulary and
416 data syntax. In such cases, the more stringent standard should be communicated and adhered
417 to in compliance with the original standard description.

418
419 Other issues, which we did not explore in more detail here, included duplicated data submissions,
420 contradictory primer references, and conflicting metadata entries for the same run. The latter is
421 especially difficult to track and often only detectable after a manual check of the data, metadata,
422 and publication. In some such instances, we discovered contradictory entries for the sequencing
423 method, instrument model, library selection, sample and library names, NCBI taxonomy ID
424 (`tax_id`), and MIxS environmental package and checklist, often hidden only among the non-
425 indexed parameters. Additional to deficient metadata, based on our case study (Molari et al. in
426 prep.) we further suspect that a large proportion of raw read amplicon data (fastq files) has not
427 been archived correctly according to the ENA submission guidelines. Our study adds another
428 facet to the increasing body of work illustrating the deficits in nucleotide sequence databases
429 (Eckert et al., 2020; Hoopen et al., 2016; Jurburg et al., 2020) with drastic consequences:
430 Terabytes to Petabytes of data may not be readily interoperable and reusable, severely limiting
431 their added value, long-term impact, and future relevance.

432
433 While ongoing development of standards and their integration across disciplines⁵ is an essential
434 endeavor to increase the added-value by standard-compliant (meta)data, we think that it is crucial
435 to avoid further delays in improving (meta)data quality and FAIRness by making better use of
436 existing standards. This task is up to each individual researcher, to voluntarily use more stringent
437 checklists and provide optional parameters. Brokerage services, such as GFBio, fundamentally
438 improved metadata quality and therefore data reusability, but are too personnel-intensive to solve
439 all challenges described here. Luckily, the number of tools, platforms, and tutorials to inform and
440 facilitate sustainable data management has increased rapidly over the last two years (Olsson &
441 Hartley, 2019; Quiñones et al., 2020; Riginos et al., 2020; Sansone et al., 2019). However, to
442 encourage such initiatives, primarily a shift in the recognition and scientific value system is
443 required to provide incentives for proper data archival and publication, including standardized

⁵ <https://www.tdwg.org/community/gbwg/MIxS/>

444 metadata, to enable long-term reuse of the data (Riginos et al., 2020; Sansone et al., 2019;
445 Westoby, Falster, & Schrader, 2021).

446

447 In the following we provide a (non-exhaustive) list of suggestions to address some of the most
448 acute deficits of data and metadata FAIRness for nucleotide sequence data from the perspective
449 of molecular ecologists. We hope that they are easy to implement and can have a potentially large
450 positive impact on the research fields relying on such data as well as derived applications in
451 biodiversity and conservation policy and management strategies.

452

453 Suggestions @researchers:

- 454 ● Make use of existing checklists and data brokerage services, and use checklists beyond
455 mandatory parameters. For instance, the MIxS parameters `target_gene`,
456 `target_subfragment`, and `pcr_primers` should be supplied for all sequencing read data that
457 was generated with `library_selection="PCR"` AND `library_strategy="AMPLICON"`.
- 458 ● Enter data diligently and according to the specified format to facilitate interoperability. It is
459 not only important **that** (meta)data is archived, but also **how**.
- 460 ● Whenever possible, use ontologies and actively contribute to the improvement of
461 ontologies by suggesting so far missing terms to the ontology developers. Many
462 ontologies, like ENVO, have a low-threshold for suggesting a new term, e.g. opening an
463 issue on GitHub, and provide guidance on using ontology terms in the context of MIxS⁶.
- 464 ● Update data submissions if additional information (manuscript DOI, accession numbers of
465 related data sets) becomes available.

466

467 Suggestions @research institutions

- 468 ● Invest in capacity development and training of early career researchers to avoid incorrect
469 data (and associated metadata) submissions due to inexperience, and raise awareness
470 for the persistently high value of FAIR data. Data archiving is not a trivial task, and it is at
471 least as important as manuscript publications.
- 472 ● Incentivization of good data management through recognition of data publications towards
473 career progression metrics.

474

475

⁶ <https://github.com/EnvironmentOntology/envo/wiki/Using-ENVO-with-MIxS>

476 Suggestions @funding agencies

- 477 • Funding as an incentive for good data management to promote a stronger inclusion of
478 data and information stewardship in scientific projects (e.g. data management plan as
479 prerequisite as already implemented by several funding agencies)
- 480 • Allocate additional funding for data technicians and data managers for the successful
481 implementation of data management plans.

482

483 Suggestions @reviewers:

- 484 • Review the submitted data and metadata for a scientific manuscript as thoroughly as the
485 manuscript text. FAIR data archiving should be as important a criterion for manuscript
486 publication as scientific soundness.

487

488 Suggestions @publishers

- 489 • If feasible, make data availability statements (accession numbers) accessible outside
490 access restrictions, so that publications can be more easily and automatically linked to the
491 data sets.

492

493 Suggestions @databases:

- 494 • Implement automated checkpoints for data consistency, specifically related to the use of
495 controlled vocabulary (ontologies) and empty entries for mandatory parameters (e.g.
496 nominal length).
- 497 • Harmonize the documentation about the parameters that use a controlled vocabulary
498 (ontologies) across the different resources (MIxS, ENA) by choosing the more stringent
499 standard.
- 500 • Expand the list of indexed metadata (MIxS) parameters in a concerted effort with the
501 scientific community to promote stronger adherence to more extensive checklists and
502 standards.
- 503 • Facilitate easy access to embargoed data sets for reviewing purposes of data and
504 metadata.

505

506 **Acknowledgements**

507

508 We would like to thank Pier Luigi Buttigieg for inspiring discussions about data management

509 and data mining and the FAIRness of it all.

510

511 References

- 512
- 513 Bax, N. J., Miloslavich, P., Muller-Karger, F. E., Allain, V., Appeltans, W., Batten, S. D., ... Tyack, P. L.
514 (2019). A response to scientific and societal needs for marine biological observations. *Frontiers in*
515 *Marine Science*, 6(JUL), 1–22. doi: 10.3389/fmars.2019.00395
- 516 Buttigieg, P. L., Fadeev, E., Bienhold, C., Hehemann, L., Offre, P., & Boetius, A. (2018). Marine microbes
517 in 4D — using time series observation to assess the dynamics of the ocean microbiome and its links
518 to ocean health. *Current Opinion in Microbiology*, 43, 169–185. doi: 10.1016/j.mib.2018.01.015
- 519 Buttigieg, P. L., Morrison, N., Smith, B., Mungall, C. J., & Lewis, S. E. (2013). The environment ontology:
520 Contextualising biological and biomedical entities. *Journal of Biomedical Semantics*, 4(1), 1–9. doi:
521 10.1186/2041-1480-4-43
- 522 Buttigieg, P. L., Pafilis, E., Lewis, S. E., Schildhauer, M. P., Walls, R. L., & Mungall, C. J. (2016). The
523 environment ontology in 2016: Bridging domains with increased scope, semantic density, and
524 interoperability. *Journal of Biomedical Semantics*, 7(1), 1–12. doi: 10.1186/s13326-016-0097-6
- 525 Canonico, G., Buttigieg, P. L., Montes, E., Muller-Karger, F. E., Stepien, C., Wright, D., ... Murton, B. (2019).
526 Global observational needs and resources for marine biodiversity. *Frontiers in Marine Science*,
527 6(JUL), 1–20. doi: 10.3389/fmars.2019.00367
- 528 Cook, C. E., Bergman, M. T., Cochrane, G., Apweiler, R., & Birney, E. (2018). The European Bioinformatics
529 Institute in 2017: Data coordination and integration. *Nucleic Acids Research*, 46(D1), D21–D29. doi:
530 10.1093/nar/gkx1154
- 531 Diepenbroek, M., Glöckner, F. O., Grobe, P., Güntsch, A., Huber, R., König-Ries, B., ... Triebel, D. (2014).
532 Towards an Integrated Biodiversity and Ecological Research Data Management and Archiving
533 Platform : The German Federation for the Curation of Biological Data (GFBio). *Informatik 2014 – Big*
534 *Data Komplexität Meistern. GI-Edition: Lecture Notes in Informatics (LNI) - Proceedings*, 1711–1724.
- 535 Eckert, E. M., Di Cesare, A., Fontaneto, D., Berendonk, T. U., Bürgmann, H., Cytryn, E., ... Corno, G.
536 (2020). Every fifth published metagenome is not available to science. *PLoS Biology*, 18(4), 1–7. doi:
537 10.1371/journal.pbio.3000698
- 538 Eren, A. M., Morrison, H. G., Lescault, P. J., Reveillaud, J., Vineis, J. H., & Sogin, M. L. (2014). Minimum
539 entropy decomposition : Unsupervised oligotyping for sensitive partitioning of high- throughput marker
540 gene sequences. *The ISME Journal*, 9(4), 968–979. doi: 10.1038/ismej.2014.195
- 541 Field, D., Amaral-Zettler, L., Cochrane, G., Cole, J. R., Dawyndt, P., Garrity, G. M., ... Wooley, J. (2011).
542 The Genomic Standards Consortium. *PLoS Biology*, 9(6), 8–10. doi: 10.1371/journal.pbio.1001088
- 543 Fillinger, S., de la Garza, L., Peltzer, A., Kohlbacher, O., & Nahnsen, S. (2019). Challenges of big data
544 integration in the life sciences. *Analytical and Bioanalytical Chemistry*, 411(26), 6791–6800. doi:
545 10.1007/s00216-019-02074-9
- 546 Fukuda, A., Kodama, Y., Mashima, J., Fujisawa, T., & Ogasawara, O. (2021). DDBJ update: Streamlining
547 submission and access of human data. *Nucleic Acids Research*, 49(D1), D71–D75. doi:
548 10.1093/nar/gkaa982
- 549 Harrison, P. W., Ahamed, A., Aslam, R., Alako, B. T. F., Burgin, J., Buso, N., ... Cochrane, G. (2021). The
550 European Nucleotide Archive in 2020. *Nucleic Acids Research*, 49(D1), D82–D85. doi:
551 10.1093/nar/gkaa1028
- 552 Hoopen, P. Ten, Amid, C., Buttigieg, P. L., Pafilis, E., Bravakos, P., O-Tárraga, A. M. C., ... Cochrane, G.

- 553 (2016). Value, but high costs in post-deposition data Curation. *Database*, 2016, 1–10. doi:
554 10.1093/database/bav126
- 555 Jurburg, S. D., Konzack, M., Eisenhauer, N., & Heintz-Buschart, A. (2020). The archives are half-empty:
556 an assessment of the availability of microbial community sequencing data. *Communications Biology*,
557 3(1). doi: 10.1038/s42003-020-01204-9
- 558 Malone, J., Holloway, E., Adamusiak, T., Kapushesky, M., Zheng, J., Kolesnikov, N., ... Parkinson, H.
559 (2010). Modeling sample variables with an Experimental Factor Ontology. *Bioinformatics*, 26(8),
560 1112–1118. doi: 10.1093/bioinformatics/btq099
- 561 NCBI Resource Coordinators. (2018). Database resources of the National Center for Biotechnology
562 Information. *Nucleic Acids Research*, 46(D1), D8–D13. doi: 10.1093/nar/gkx1095
- 563 Olsson, T. S. G., & Hartley, M. (2019). Lightweight data management with dtool. *PeerJ*, 2019(3). doi:
564 10.7717/peerj.6562
- 565 Quiñones, M., Liou, D. T., Shyu, C., Kim, W., Vujkovic-Cvijin, I., Belkaid, Y., & Hurt, D. E. (2020).
566 “mETAGENOTE: A simplified web platform for metadata annotation of genomic samples and
567 streamlined submission to NCBI’s sequence read archive.” *BMC Bioinformatics*, 21(1), 1–12. doi:
568 10.1186/s12859-020-03694-0
- 569 Reiser, L., Harper, L., Freeling, M., Han, B., & Luan, S. (2018). FAIR: A Call to Make Published Data More
570 Findable, Accessible, Interoperable, and Reusable. *Molecular Plant*, 11(9), 1105–1108. doi:
571 10.1016/j.molp.2018.07.005
- 572 Riginos, C., Crandall, E. D., Liggins, L., Gaither, M. R., Ewing, R. B., Meyer, C., ... Deck, J. (2020). Building
573 a global genomics observatory: Using GEOME (the Genomic Observatories Metadatabase) to
574 expedite and improve deposition and retrieval of genetic data and metadata for biodiversity research.
575 *Molecular Ecology Resources*, 20(6), 1458–1469. doi: 10.1111/1755-0998.13269
- 576 Ryan, M., Schloter, M., Berg, G., Kinkel, L. L., Eversole, K., Macklin, J. A., ... Sessitsch, A. (2020). Towards
577 a unified data infrastructure to support European and global microbiome research- A call to action.
578 *Environmental Microbiology*, 00, 1–4. doi: 10.1111/1462-2920.15323
- 579 Sansone, S. A., McQuilton, P., Rocca-Serra, P., Gonzalez-Beltran, A., Izzo, M., Lister, A. L., & Thurston,
580 M. (2019). FAIRsharing as a community approach to standards, repositories and policies. *Nature*
581 *Biotechnology*, 37(4), 358–367. doi: 10.1038/s41587-019-0080-8
- 582 Stevens, I., Mukarram, A. K., Hörtenhuber, M., Meehan, T. F., Rung, J., & Daub, C. O. (2020). Ten simple
583 rules for annotating sequencing experiments. *PLoS Computational Biology*, 16(10), 1–7. doi:
584 10.1371/journal.pcbi.1008260
- 585 Westoby, M., Falster, D. S., & Schrader, J. (2021). Motivating data contributions via a distinct career
586 currency. *Proceedings of the Royal Society B: Biological Sciences*, 288(1946). doi:
587 10.1098/rspb.2020.2830
- 588 Wilkinson, M. D., Dumontier, M., Aalbersberg, Ij. J., Appleton, G., Axton, M., Baak, A., ... Mons, B. (2016).
589 The FAIR Guiding Principles for scientific data management and stewardship. *Scientific Data*, 3, 1–9.
590 doi: 10.1038/sdata.2016.18
- 591 Yilmaz, P., Kottmann, R., Field, D., Knight, R., Cole, J. R., Amaral-Zettler, L., ... Glöckner, F. O. (2011).
592 Minimum information about a marker gene sequence (MIMARKS) and minimum information about
593 any (x) sequence (MIxS) specifications. *Nature Biotechnology*, 29(5), 415–420. doi: 10.1038/nbt.1823
- 594

595 **Supplement**

596

597 **SI table 1:** Number of amplicon and WGS cases (runs) submitted per year, separated by the
598 availability of various metadata parameters, use of MlxS environmental package and GFBio as
599 brokerage service. Metadata availability (Category) as explained in figures 1 and 2.

600

601 **SI figure 1:** Number of cases (WGS example) with and without metadata available for geographic
602 coordinates (latitude, longitude) and nominal length. For geographic coordinates, 'tsv' is referring
603 to the information provided for indexed metadata parameters in the TSV search output, while 'xml'
604 is referring to non-indexed metadata only accessible in the XML view of the ENA sample or
605 experiment. The percentages summarize cases over all years and are rounded to integers. The
606 plots on the right only show the cases with samples submitted according to a MlxS checklist and
607 environmental package. The horizontal line in the plots on the left indicates the y-axis range of
608 the plots on the right.

609

610 **SI figure 2:** Number of cases (WGS example) with values for environment_biome,
611 environment_material, and environment_feature according to ENVO. Yes: exact match to ENVO
612 term ID; some: character string match to ENVO term name (including matches after extensive
613 character string manipulation); filled: a value is provided, but not ENVO-compatible; no: no entry.
614 The percentages summarize cases over all years and are rounded to integers. The plots on the
615 right only show the cases with samples submitted according to a MlxS checklist and
616 environmental package. The horizontal line in the plots on the left indicates the y-axis range of
617 the plots on the right. The category 'filled' is overrepresented in the year of 2019 for 'biome' and
618 'material', mainly due to runs from a single study with more than 2000 runs.