

## Metadata preservation and stewardship for genomic data is possible, but must happen now

Eric D. Crandall<sup>1</sup>, Rachel H. Toczydlowski<sup>2,3</sup>, Libby Liggins<sup>4</sup>, Ann E. Holmes<sup>5</sup>, Maryam Ghoojaei<sup>6</sup>, Michelle R. Gaither<sup>6</sup>, Briana E. Wham<sup>7</sup>, Andrea L. Pritt<sup>8</sup>, Cory Noble<sup>4</sup>, Tanner J. Anderson<sup>9</sup>, Randi L. Barton<sup>10,11</sup>, Justin T. Berg<sup>12,13</sup>, Sofia G. Beskid<sup>14</sup>, Alonso Delgado<sup>15</sup>, Emily Farrell<sup>6</sup>, Nan Himmelsbach<sup>16</sup>, Samantha R. Queeno<sup>9</sup>, Thienthanh Trinh<sup>6</sup>, Courtney Weyand<sup>17</sup>, Andrew Bentley<sup>18</sup>, John Deck<sup>19</sup>, Cynthia Riginos<sup>20</sup>, Gideon S. Bradburd<sup>2,20</sup>, Robert J. Toonen<sup>22</sup>

1. Department of Biology, Pennsylvania State University, University Park, Pennsylvania, PA 16801
2. Department of Integrative Biology - Ecology, Evolution, and Behavior Program, Michigan State University, East Lansing, MI 48824
3. Current Affiliation: Institute for Applied Ecosystem Studies, Northern Research Station, United States Forest Service, Rhinelander, WI, 54501, U.S.
4. School of Natural Sciences, Massey University, Auckland, Aotearoa New Zealand
5. Department of Animal Science, University of California-Davis, Davis, CA 95616
6. Department of Biology, University of Central Florida, Orlando, FL 32816
7. Department of Research Informatics and Publishing, The Pennsylvania State University Libraries, University Park, Pennsylvania, PA 16801
8. Madlyn L. Hanes Library, The Pennsylvania State University Libraries, Pennsylvania State University, Harrisburg Campus, 351 Olmsted Drive, Middletown, PA 17057
9. Department of Anthropology, University of Oregon, Eugene, OR 97403
10. California State University Monterey Bay, 100 Campus Circle, Seaside, CA, 93955
11. Moss Landing Marine Laboratories, Moss Landing, CA 95039
12. UOG Marine Laboratory, University of Guam, Mangilao, Guam 96910
13. Current Affiliation: Department of Oceanography, University of Hawai'i at Mānoa, Kane'ohe, HI 96744
14. Department of Integrative Biology, University of Texas at Austin, Austin, TX 78712
15. Department of Evolution, Ecology, and Organismal Biology, The Ohio State University; Columbus, OH 43210
16. Department of Natural Science, Hawai'i Pacific University, Honolulu, HI 96813
17. Department of Biological Sciences, Auburn University, Auburn, AL 36849
18. Biodiversity Institute, University of Kansas, Lawrence, KS, 66045
19. Berkeley Natural History Museums, University of California, Berkeley, California
20. School of Biological Sciences, The University of Queensland, Brisbane, Australia
21. Current affiliation: Department of Ecology and Evolutionary Biology, University of Michigan 1105 N. University Ave, Ann Arbor, MI, 48109
22. Hawai'i Institute of Marine Biology, University of Hawai'i at Mānoa, Kane'ohe, HI 96744

## 1 **Abstract**

2  
3 Genetic diversity within species represents a fundamental yet underappreciated level of  
4 biodiversity. Because genetic diversity can indicate species and population resilience to changing  
5 climate, its measurement is relevant to many national and global conservation policy targets.  
6 Many studies of evolutionary biology, molecular ecology and conservation genetics produce  
7 large amounts of genome-scale genetic diversity data for wild populations. While open data  
8 policies have ensured an abundance of freely available genomic data stored in the databases of  
9 the International Nucleotide Sequence Database Collaboration (INSDC), only about 13% of  
10 current accessions have the associated spatial and temporal metadata in INSDC necessary to be  
11 reused in monitoring programs, macrogenetic studies, or for acknowledging the sovereignty of  
12 nations or Indigenous Peoples. We undertook a “distributed datathon” to quantify the availability  
13 of these missing metadata in sources external to the INSDC and to test the hypothesis that these  
14 metadata decay with time. We also worked to remediate these missing metadata by extracting  
15 them, when present, from associated published papers, online repositories, and/or from direct  
16 communication with authors. Starting with 848 programmatically identified candidate datasets  
17 (INSDC BioProjects), we manually determined that 561 contained samples from wild  
18 populations. We successfully restored spatiotemporal metadata (locality name and/or geospatial  
19 coordinates and collection year) for 78% of these 561 datasets (N = 440 BioProjects comprising  
20 45,105 individuals or BioSamples from 762 species in 17 phyla). We also quantified the  
21 availability of 33 additional categories of metadata in sources external to the INSDC.  
22 Information about associated publications and the type of habitat from which the samples were  
23 taken was the most easily found; information about sampling permits was the most challenging  
24 to locate. Looking at papers and online repositories was much more fruitful than contacting  
25 authors, who only replied to our email requests 45% of the time. Overall, 23% of our email  
26 queries to authors discovered useful metadata. Importantly, we found that the probability of  
27 retrieving spatiotemporal metadata declines significantly with the age of the dataset, with a  
28 13.5% yearly decrease for metadata located in published papers or online repositories and up to a  
29 22% yearly decrease for metadata that were only available from authors. This observable  
30 metadata decay, mirrored in studies of other types of biological data, should motivate swift  
31 updates to data sharing policies and researcher practices to ensure that the valuable context  
32 provided by metadata is not lost forever.

33  
34  
35

36

## 37 **Introduction**

38

39 Genetic diversity is the foundational layer of biodiversity. Just as ecosystem health and  
40 resilience depends on the diversity of its component species, the health and resilience of each  
41 species depends on its genomic diversity (Clark, 2010; Reusch, Ehlers, Hämmerli, & Worm,  
42 2005). Without genetic diversity in the form of standing allelic variation, populations and species  
43 cannot adapt to a rapidly changing climate and other anthropogenically-induced stresses  
44 (Blanchet et al., 2020; Raffard, Santoul, Cucherousset, & Blanchet, 2019). Local or global  
45 extinctions of species in turn threaten the ecosystems upon which the quality of human lives  
46 depend (Brauman et al., 2020; Des Roches, Pendleton, Shapiro, & Palkovacs, 2021).  
47 Concerningly, genetic diversity, like all levels of biodiversity, is declining rapidly during the  
48 Anthropocene across the tree of life (Exposito-Alonso et al., 2022; Leigh, Hendry, Vázquez-  
49 Domínguez, & Friesen, 2019; Miraldo et al., 2016; Pinsky & Palumbi, 2014).

50

51 Recognizing the vital importance of biodiversity to human well-being and the future of our  
52 planet, several international agreements strongly encourage the monitoring and conservation of  
53 genetic diversity in both wild and domesticated species. Foremost among these are the United  
54 Nations Sustainable Development Goal 2.5 and the international Convention on Biological  
55 Diversity (CBD) treaty, which explicitly acknowledge the importance of monitoring and  
56 conserving any component of biological diversity (including genetic diversity) that may have  
57 “actual or potential use or value for humanity.” Moreover, the CBD’s article 15 and attendant  
58 Nagoya Protocol codify procedures to ensure the sharing of benefits arising from genetic  
59 resources (such as digital sequence information; DSI) discovered or accessed within a nation’s  
60 sovereign borders. The subsequent Strategic Plan for Biodiversity 2011-2020 laid out the 20  
61 Aichi Biodiversity Targets, including target 13, which aims to maintain the “genetic diversity of  
62 cultivated plants and farmed and domesticated animals and of wild relatives, including other  
63 socio-economically as well as culturally valuable species.” Now, even as we are facing shortfalls  
64 on all 20 of the Aichi Biodiversity Targets (CBD, 2020; Hoban et al., 2021), there have been  
65 calls to broaden genetic diversity targets to include *all* extant species in the New Post-2020  
66 Global Biodiversity Framework, in what has been described as a “moonshot for biology” (Hoban  
67 et al., 2020; Laikre et al., 2020; Lewin et al., 2018).

68

69 Over the last decade, advances in DNA sequencing technology have enabled the generation  
70 of genome-scale datasets of ever larger numbers of individuals, drawn from a growing variety of  
71 species (Allendorf, 2017; Hendricks et al., 2018). Researchers are now able to genotype  
72 thousands of genomic loci or sequence whole genomes from non-model species for which they  
73 have no prior genetic resources (Lou, Jacobs, Wilder, & Therkildsen, 2021; Willette et al., 2014).  
74 The shift from genetic to genomic-scale datasets is catalyzing novel conservation insights  
75 including: the detection of inbreeding depression (e.g. Kardos, Taylor, Ellegren, Luikart, &  
76 Allendorf, 2016), the discovery of subtle, previously undetectable population structure (e.g.  
77 Cheng, Gold, Rodriguez, & Barber, 2021; Gaither et al., 2018), reconstruction of demographic  
78 histories (Prada et al., 2016), the precise identification of distant pedigree relationships (e.g.  
79 Baetscher et al., 2019), uncovering cryptic species (e.g. Quattrini et al., 2019), clues about the  
80 genomic basis of local adaptation (e.g. Wilder, Palumbi, Conover, & Therkildsen, 2020) and  
81 important traits such as nutritional components (e.g. Kumar et al., 2021). Accordingly, the DSI

82 derived in these studies is highly valued as a resource equivalent to biobanks, providing essential  
83 information for conservation (Hoban et al., 2022) as well as ensuring future food security  
84 (Castañeda-Álvarez et al., 2016; Halewood et al., 2018).

85  
86 Genomic datasets record the genetic diversity of a species at a particular time and  
87 location, providing a benchmark for how populations are responding to human drivers of  
88 changing environmental conditions, cultivation, and land and sea use, as well as measuring  
89 indicators of progress toward conservation targets and goals (Hoban et al., 2022, 2020) and the  
90 genetic resources available for future cultivation or domestication (Halewood et al., 2018).  
91 However, genomic datasets can only be useful for monitoring global genetic biodiversity and the  
92 sustainable human use of genetic diversity (including benefit-sharing, Cowell et al., 2022) when  
93 archived publicly with accompanying metadata about the spatiotemporal, environmental and  
94 methodological context of the sequenced sample (Riginos et al., 2020; Scholz et al., 2022;  
95 Schriml et al., 2020).

96  
97 The genetics community has long championed open data publication with the  
98 foundational databases of the International Nucleotide Sequence Database Collaboration  
99 (INSDC; Cochrane, Karsch-Mizrachi, Takagi, & INSDC, 2016) forming in the early 1980's. In  
100 2009, the INSDC launched the Sequence Read Archive as a repository dedicated to second-  
101 generation sequence data. It has since grown exponentially to include over 600 terabytes of  
102 freely-available DNA sequence data from over 16,700 wild and domesticated eukaryotic species  
103 as of 2021 (Toczydlowski et al., 2021). Around the same time, the MIxS metadata standards  
104 (Field et al., 2008; Yilmaz et al., 2011) were defined to inform the minimum information about  
105 *what* (detailed taxonomy), *where* (GPS coordinates and habitat), *when* (collection date), *how*  
106 (sampling and sequencing protocols) and *by whom* a genetic sample was collected. Enabled by  
107 the INSDC infrastructure and encouraged by the Joint Data Archiving Policy (JDAP;  
108 <http://datadryad.org/pages/jdap>) implemented by top journals in 2011, the proportion of papers  
109 providing open access to their genetic data increased dramatically (Pope, Liggins, Keyse,  
110 Carvalho, & Riginos, 2015). However, the inclusion of accompanying metadata crucial for the  
111 reuse of these data for genetic diversity monitoring, macrogenetic studies, or identifying their  
112 provenance within national boundaries or the lands and waters of Indigenous Peoples, has lagged  
113 behind (Pope et al., 2015; Toczydlowski et al., 2021). As of 2021, out of over 300,000 SRA  
114 BioSamples that are potentially relevant to global genetic biodiversity, only ~13% had metadata  
115 indicating both the time and precise location from which they were sampled (Toczydlowski et  
116 al., 2021).

117  
118 In a timely and welcome update to their policy, INSDC now intends to extend their minimum  
119 metadata requirements to include collection date and country of origin  
120 (<https://www.insdc.org/spatio-temporal-annotation-policy-18-11-2021>). Although 'country' is  
121 legislatively aligned with the Nagoya Protocol, it is not spatially aligned with the lands and  
122 waters of Indigenous Peoples (e.g. <https://native-land.ca/>) and does not provide adequate spatial  
123 resolution for conservation monitoring. Moreover, this policy and infrastructure change will take  
124 time to implement (anticipated to be end of 2022), meaning that much of the genomic data  
125 collated over the last ~12 years for past and present populations, of immeasurable value to  
126 understanding and monitoring the biodiversity crisis, are not Findable, Accessible, Interoperable  
127 or Reusable (FAIR; Wilkinson et al., 2016). This absence of appropriate spatio-temporal

128 metadata represents the effective loss of tens to hundreds of millions of dollars of research effort  
129 for most future purposes (Schriml et al., 2020; Toczydlowski et al., 2021), rendering associated  
130 genetic data invisible to government ministries and non-governmental organizations tasked with  
131 protecting the world’s natural environment (Laikre, 2010; Laikre et al., 2020). Moreover,  
132 without spatiotemporal provenance of genomic data enabling connection to the lands and waters  
133 of Indigenous Peoples, these peoples will potentially lose out on benefits (e.g. capacity  
134 development, food security, biomedical advances) arising from genomic information originating  
135 within their territories (Liggins, Hudson, & Anderson, 2021; Marden et al., 2021; McCartney et  
136 al., 2022; Scholz et al., 2022). There is urgency in addressing this metadata gap: previous studies  
137 of morphological (Vines et al., 2014) and genetic (Pope et al., 2015) data suggest that the  
138 probability of existing metadata ever being linked to the genomic data significantly decreases  
139 over time.

140  
141 In the summer of 2020, we convened a distributed remote datathon to (1) assess the  
142 availability of metadata outside of the INSDC, (2) recover and curate metadata missing in  
143 INSDC from external sources (i.e. published research papers, other online repositories, or the  
144 authors themselves), and (3) extend our initial report on the metadata gap (Toczydlowski et al.,  
145 2021) to investigate how the recoverability of these metadata is affected by dataset age and to  
146 document shortfalls and costs of our remedial efforts. In our datathon, 14 graduate students from  
147 the USA and 12 professional academics (authors on this paper) across 3 countries worked  
148 together via Zoom, Slack and Google Sheets as “metadata curators” to establish and execute  
149 curation protocols and infill missing metadata. Collectively, we searched for metadata external to  
150 the INSDC (e.g. associated scientific publications, Dryad, museum collections) for 848 genomic  
151 datasets (INSDC BioProjects) representing 94,416 individual samples (BioSamples). The  
152 BioSamples and associated genetic sequence data in these projects were selected as they were  
153 missing at least latitude and longitude metadata in the INSDC. Our findings underscore the  
154 importance of appropriate and immediate metadata archival going forward. We provide guidance  
155 based on our collective experience gained over the datathon on practices to retain crucial  
156 metadata.

## 157 158 **Methods**

### 159 160 *Identifying BioProjects for metadata curation*

161  
162 Our datathon followed a workflow illustrated in Figure 1. We first searched the INSDC to  
163 identify datasets that were potentially relevant to monitoring genetic diversity of wild  
164 populations (including wild relatives of domesticated species) but were lacking critical metadata  
165 about the latitude and longitude of the sampling location (as described in Toczydlowski et al.,  
166 2021). On November 7, 2019, we searched the INSDC BioProject database using the *rentrez* R  
167 package (Winter, 2017) and custom R scripts to locate datasets that contained genomic-level  
168 DNA sequence data from eukaryotic species, excluding human and metagenomic datasets. We  
169 then downloaded all metadata associated with these BioProjects from the INSDC PubMed,  
170 Taxonomy, BioSample and SRA databases. We further filtered the BioProjects to remove  
171 BioSamples (sequenced individuals) from species whose population dynamics and evolution are  
172 largely governed by humans: pathogens and their vectors, model organisms, and domesticated  
173 species. We used custom lists for each category of non-wild organisms (Supplementary

174 Materials S1). These filters helped to target our efforts at recovering metadata for the datasets  
175 most likely to be of relevance to genetic diversity monitoring of wild populations. Additional  
176 details about the search and filtering steps can be found in Supplemental Materials of  
177 Toczydlowski et al. 2021. These filtering steps yielded 848 BioProjects that entered our datathon  
178 curation pipeline. These BioProjects contained at least 5 individuals that lacked geospatial  
179 coordinates in the INSDC and were potentially sampled from wild populations.

180

181 We built a blank template to enter metadata into that we located external to the INSDC using  
182 the Genomic Observatories Metadatabase (GEOME; Deck et al., 2017; Riginos et al., 2020)  
183 This template included terms that described collection date and location, habitat etc. Table 1  
184 gives definitions of 16 required and recommended metadata terms, while Supplemental Materials  
185 S2 provides an example template with all 36 metadata terms and their definitions and controlled  
186 vocabularies. We then programmatically filled this template for each of the 848 BioProjects in  
187 our pipeline with the metadata already present in the INSDC using custom R scripts.

188

189 *Recovering and curating metadata external to the INSDC*

190

191 Metadata curators were each randomly assigned a set of BioProjects. Curators followed a  
192 standard protocol (Figure 1; Supplemental Materials S3) to locate and enter associated metadata  
193 for samples in each BioProject that were missing in the INSDC but reported in external sources  
194 (e.g. associated published scientific papers). Briefly, we first searched for associated publications  
195 by googling the BioProject PRJ accession number and/or key words associated with the  
196 BioProject (e.g. author, affiliation, funding information, species, and locality names). While the  
197 INSDC includes terms to capture associated publications, these were only populated for 7% of  
198 the 848 BioProjects that we examined. We then skimmed each paper to determine if the  
199 BioProject contained sequence data from wild individuals. BioProjects comprising over 50%  
200 BioSamples from non-wild individuals (brood stock, laboratory stock, domesticated species, and  
201 some seed banks) were marked as “not relevant” and removed from further curation. Non-wild  
202 BioSamples that were a minority in their BioProject were designated as such under the  
203 *establishmentMeans* metadata term, and further metadata recovery was not prioritized for these  
204 samples. BioProjects that pooled their genetic samples without barcodes (precluding estimates of  
205 genetic diversity at the level of individual) were also marked as “not relevant” and removed from  
206 the curation pipeline. Once relevance to our efforts was determined, curators looked in figures,  
207 tables, and supplemental information of associated publications and/or linked online repositories  
208 (e.g. Data Dryad, Github, Zenodo, and databases of biocollections such as museums and  
209 herbaria) for sample-level metadata from 36 metadata terms (most defined as Darwin Core  
210 terms; Wieczorek et al., 2012). When necessary, curators converted metadata to standard formats  
211 (e.g. degree minutes seconds data were converted to decimal degrees), and then added metadata  
212 into the pre-generated GEOME template for that BioProject. After performing quality control,  
213 these metadata could then be easily uploaded to GEOME and potentially cross-walked into the  
214 appropriate INSDC databases.

215

216 After adding all metadata that could be gleaned from the associated paper(s) into the  
217 GEOME templates, curators made a structured comment on a master spreadsheet (Supplemental  
218 Materials S4) indicating whether metadata for each of the required and recommended terms were  
219 absent for all BioSamples (“none”), present for less than 50% of BioSamples (“some”), present

220 for greater than 50% of BioSamples (“most”), or iv) present for all BioSamples (“all”). If the  
221 paper was missing information from one of six or seven “required” terms (georeference-able  
222 *locality* OR [*decimalLatitude* AND *decimalLongitude*], *coordinateUncertaintyInMeters*,  
223 *georeferenceProtocol*, *habitat*, *environmentalMedium*, *yearCollected*), the curator flagged the  
224 BioProject to initiate author contact. An additional nine metadata terms were considered  
225 “recommended”: missing metadata in these fields alone did not trigger an author contact but

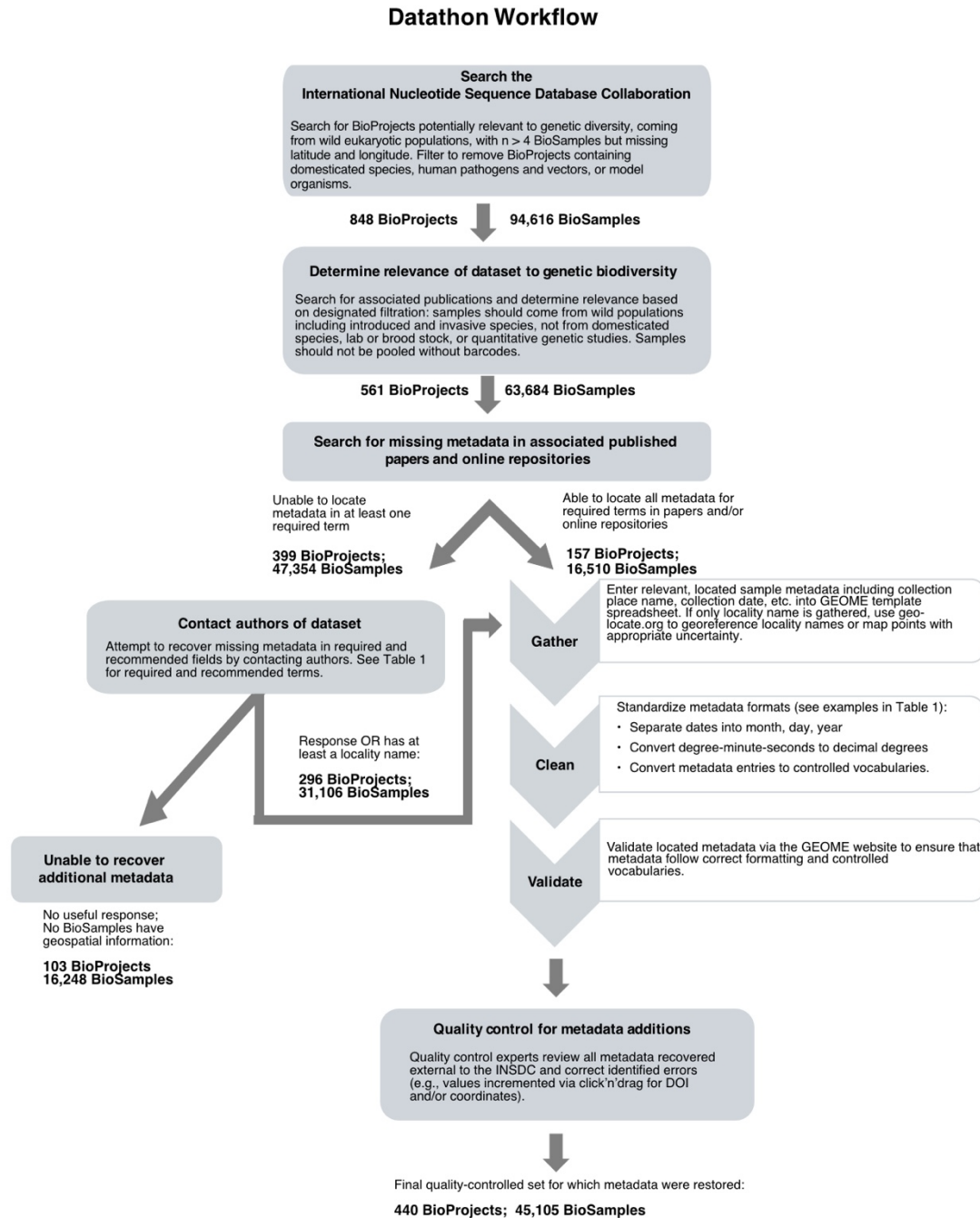


Figure 1. Datathon workflow. The number of BioProjects and BioSamples remaining after each step are given below the step.

226 curators and authors were asked to populate these fields as completely as possible. These  
227 recommended terms included *country*, *establishmentMeans*, *permitInformation*,  
228 *associatedReferences*, *preservative*, and 4 terms that tracked genetic data derived from the raw  
229 reads such as SNP genotypes or sequence alignments. Progress and notes at each curation step  
230 were tracked as meta-metadata on the master spreadsheet.

231  
232 After a quality-control step to ensure that author names and email addresses found in papers  
233 were input correctly, corresponding authors of the paper were contacted by email (see  
234 Supplemental Material S5 for email text) using the Yet Another Mail Merge add-on for Google  
235 Sheets (yamm.com). Using YAMM allowed us to track email receipts and responses. If an email  
236 was undeliverable, we used our best efforts to locate an alternate email. We were able to  
237 successfully deliver email queries for all 351 of 492 relevant BioProjects that met the criteria for  
238 author contact. About two weeks after sending the initial email, curators sent reminder emails to  
239 unresponsive authors at least once and at most twice. This process emulated the efforts of a  
240 reasonably persistent researcher to obtain metadata important to their research.

241  
242 Once curators had gathered and entered all available metadata into the GEOME template for  
243 a BioProject (from online sources external to the INSDC and/or directly from authors), curators  
244 validated the metadata at the GEOME website (geome-db.org). This programmatic GEOME  
245 validation ensured that metadata in each field were correctly formatted and within controlled  
246 vocabularies, e.g. for *country* and *habitat* (see lists tab of Supplemental Materials S3 for all  
247 controlled vocabularies). Following validation, MRG, BEW, AP and EDC performed a final  
248 quality-control check (QC; steps described in Supplemental Material S1). Filled and QCed  
249 GEOME templates for each BioProject will be uploaded to the GEOME database.

250  
251 *Investigating metadata decay*

252  
253 We investigated the effect of BioProject age on the probability that we were able to recover  
254 metadata information for 11 metadata categories. Previous investigations of metadata have  
255 indicated rapid decay when data are not publicly archived (Roche et al., 2014; Vines et al., 2014,  
256 2013). We used Bayesian logistic regression to fit four distinct models to investigate the  
257 relationships between BioProject age (number of days between publication in the INSDC and  
258 November 7, 2019) and: A) the probability that metadata could be retrieved from INSDC,  
259 associated published papers, and/or repositories, B) the probability that we received an author  
260 response for the 351 BioProjects that triggered an author contact via email, C) the probability  
261 that authors provided any metadata, given that they responded and D) the probability that authors  
262 provided metadata for a majority of samples, given that they responded.

263  
264 Information about the collection date and location of a sample are the most critical pieces of  
265 metadata required to identify the Indigenous provenance of genomic data and make genomic  
266 sequence data repurposable for monitoring and fundamental data synthesis research, so we  
267 focused our investigations on these two categories; we refer to the aggregate as spatiotemporal  
268 metadata. We defined a BioProject as having spatiotemporal data if collection dates, and  
269 latitudes and longitudes and/or locality were present for at least 50% of the BioSamples that it  
270 contained. In model C, we counted a gain in collection year, or place name, or latitude/longitude  
271 for any number of BioSamples as recovery of metadata. In model D, we only counted increases

272 in metadata where BioProjects had incomplete spatiotemporal metadata for > 50% of its  
273 BioSamples and then had spatiotemporal metadata present for > 50% of BioSamples after  
274 contacting authors. That is, model C assessed the probability of recovering any metadata external  
275 to the INSDC, and model D assessed the probability of recovering metadata for the majority of  
276 samples. In a supplemental analysis, we investigated how metadata within individual  
277 spatiotemporal terms and other important metadata terms (i.e. required and recommended terms,  
278 Table 1) decayed, including: *decimalLongitude* and *decimalLatitude*, *collectionYear*, *country*,  
279 *locality*, *habitat*, *environmentalMedium*, *materialSampleID*, *permitInformation*,  
280 *associatedReferences* (publication DOI), *preservative*, and derived genetic data terms  
281 (Supplemental Materials S6).

282  
283 We conducted all statistical analyses at the level of BioProject (as opposed to BioSamples or  
284 genomic sequences), because presence/absence of metadata for BioSamples within a given  
285 BioProject was highly correlated (Toczydlowski et al., 2021). Curators also tracked metadata  
286 recovery efforts at the level of BioProject for convenience, and we contacted authors about entire  
287 BioProjects rather than individual BioSamples. In each set of models, we removed BioProjects  
288 that already had complete metadata in the category of interest, and therefore could not gain any  
289 more.

290  
291 We analyzed the effect of BioProject age on our response variables (model A: successful  
292 metadata retrieval; model B: author response; model C: author provided any metadata  
293 conditional on a response to our query; model D: author provided metadata for the majority of  
294 samples conditional on a response to our query) using generalized linear models. In each  
295 analysis, we modeled our response variable as a Bernoulli-distributed variable with a probability  
296 of success that was a linear function of our predictor variable: BioProject age. In each analysis,  
297 the parameters of our model were a global mean probability of success and an effect size of  
298 BioProject age on probability of success for that response variable. Analyses used the canonical  
299 inverse-logit inverse link function. In mathematical notation, our model was:

300

$$Y_i \sim \text{Bernoulli}(p_i)$$

$$p_i = \frac{1}{1 + e^{-\theta_i}}$$

$$\theta_i = \mu + \beta \times X_i$$

301  
302 Where  $Y_i$  is the  $i$ th outcome (response variable),  $p_i$  is the probability of successfully observing that  
303 outcome,  $\mu$  is the global mean probability of success, and  $\beta$  is the effect of BioProject age on the  
304 transformed probability of success for that outcome ( $\theta_i$ ). We had no strong prior beliefs about the  
305 effect of BioProject age on success in each of the four analyses we ran; to reflect these beliefs,  
306 the priors we placed on our parameters were:  $\beta \sim N(0,10)$ ;  $\mu \sim N(0,10)$ . All statistical analyses  
307 were performed using Rstan version 2.21.2 [50] running 4 independent chains for 2,000  
308 iterations, thinning to sample only every 4th iteration to reduce autocorrelation, and discarding  
309 the first 1,000 iterations as burn-in. To assess the significance of the effect of BioProject age on

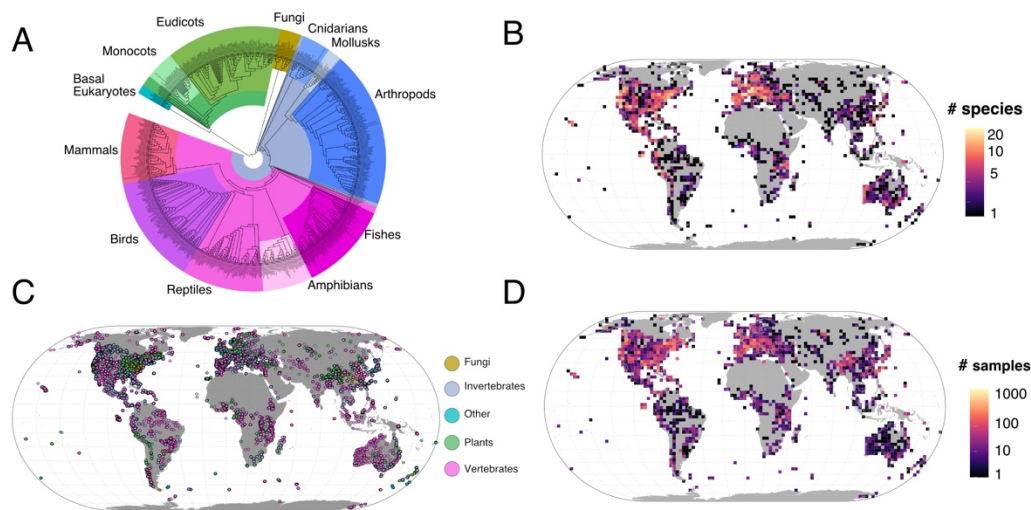


Figure 2. The taxonomic and geographic scope of the datathon. A) A cladogram of 719 of the 762 species from BioProjects that passed through the final quality control step. This is a subtree of the Open Tree of Life (Hinchcliff et al. 2015) generated with the rotl package for R (Michonneau et al. 2016) and visualized with iTOL software (itol.embl.de; Letunic and Bork 2021). B) Heatmap of species for BioSamples for which spatial coordinates were recovered by the datathon. C) Map showing the geographic distribution of broad taxonomic categories of these BioSamples. D) Heatmap of these BioSamples.

310 success of each outcome, we determined whether the 95% equal-tailed credible interval of the  
311 marginal distribution on  $\beta$  contained 0; if it did, the effect of BioProject age was deemed not  
312 significant.

313

## 314 Results

315

316 We identified 848 INSDC BioProjects, representing 94,416 BioSamples from individual  
317 eukaryotic organisms that lacked geospatial coordinates and had at least five putatively wild  
318 individuals as determined by our filters. Curators were able to locate associated published  
319 scientific papers for 741 of these 848 BioProjects. Reading these papers revealed 561  
320 BioProjects with a majority of relevant, truly wild individuals, comprising 63,684 individuals  
321 from 873 species. After scouring associated published papers for metadata and contacting  
322 authors, a total of 440 BioProjects with 45,105 BioSamples from 762 species in 17 eukaryotic  
323 phyla (Figure 2A) had geospatial data (either coordinates or a locality name) and were passed  
324 through quality control for eventual upload to GEOME. BioSamples that passed through the  
325 datathon came from all continents and all major oceans (Figure 2B-D).

326

327 For the subset of BioProjects that we focused on (those that were missing latitude and  
328 longitude), datathon curators were able to recover metadata for a majority of BioSamples in a  
329 BioProject as follows (depicted in Figure 3). For geospatial coordinates, nearly 60% could be

330 found in an associated publication or online repository. While nearly 30% of these BioProjects  
331 did already contain information about collection year in the INSDC, curators were only able to  
332 recover an additional 21% from papers or online repositories. Datathon curators recovered  
333 metadata regarding habitat, environmental medium (the media displaced by the sampled  
334 organism) and publication DOI for over 80% of BioProjects from published papers and their  
335 supplemental information. Additional large gains in BioProjects were made from online sources  
336 external to the INSDC for locality (48.8%), and country name (39.8%). Notably, permit  
337 information was the least available of any of the metadata categories that we explored. There is  
338 no permit metadata term in the INSDC and curators found permit information in papers for only  
339 21% of BioProjects.

340  
341 Contacting authors yielded comparatively less metadata than our search of papers and  
342 supplemental information, although it should be noted that this step was secondary to looking in  
343 papers and online. Out of 351 author contact attempts, we received 158 responses (45% response  
344 rate). Of the 158 responses, 80 (51%) provided at least some missing metadata, yielding an  
345 overall “useful author response rate” of 23%. Through contacting authors, we were able to  
346 recover collection year metadata for an additional 9% of BioProjects, and geospatial coordinates  
347 for an additional 8.5% percent of BioProjects. Gains in other metadata categories were all less  
348 than 5%, with permit information showing only a 1.2% increase from authors.

349  
350

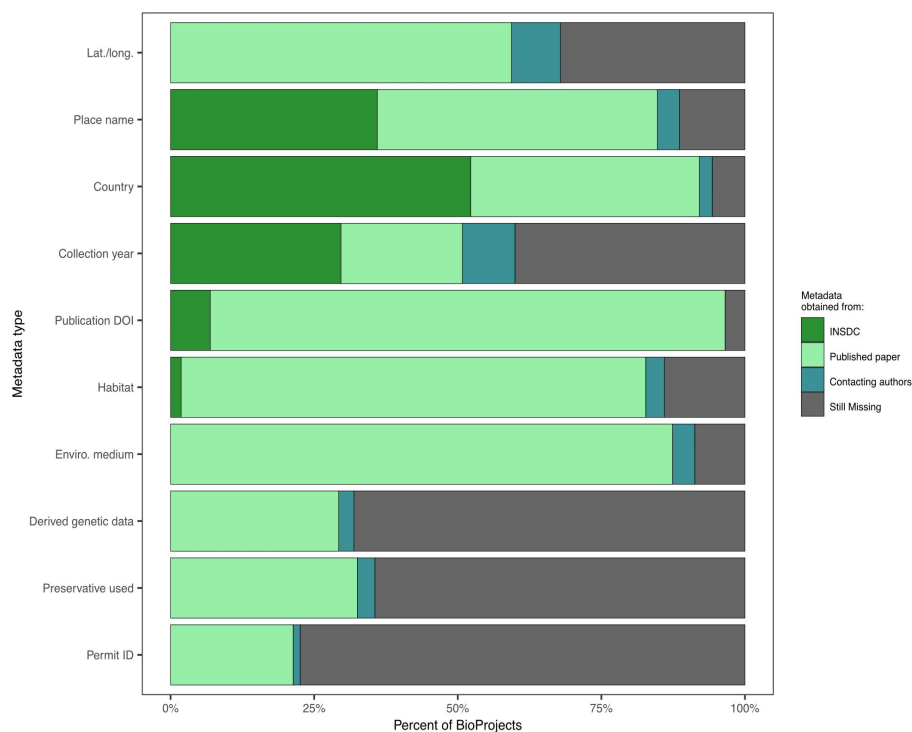


Figure 3. Stacked bars showing the percent of BioProjects for which metadata were found from each of three sources, across 10 priority metadata categories.

351

352 The age (time since deposition into the INSDC) of the BioProject had a strong effect on  
 353 whether metadata could be recovered. After searching for metadata within the INSDC and within  
 354 published papers, we found that spatiotemporal metadata (defined as year AND geospatial  
 355 coordinates OR locality), had a mean odds ratio of 0.865 (95% highest posterior density credible  
 356 interval [HPD CI]: 0.775 - 0.964; Figure 4A). This indicates that for every year after a  
 357 BioProject is published to the SRA, there is about a 13.5% decrease (HPD CI: 3.6% - 22.5%) in  
 358 the probability that its metadata can be found in the SRA, in papers or elsewhere online. On the  
 359 other hand, there was a strong positive effect of BioProject age on whether an attempt to contact  
 360 the authors was answered, with a 25.5% increase in the probability of a reply of any kind for  
 361 every year after SRA publication (mean odds ratio of 1.255; 95% HPD CI: 1.120 - 1.412; Figure  
 362 4B). In other words, we were more likely to get an email response for older datasets. However,  
 363 given a response, the probability that authors furnished any amount of metadata for year OR  
 364 coordinates OR locality decreased with BioProject age by 21% per year (odds ratio 0.810; 95%  
 365 HPD CI: 0.680 - 0.949; Figure 4C). Similarly, the probability that the authors provided metadata  
 366 for year AND coordinates OR locality for a majority of BioSamples decreased by 22% per year  
 367 (odds ratio: 0.819; 95% HPD CI: 0.671 - 0.994; Figure 4D).  
 368

369 Figures for Bayesian logistic regressions of BioProject age on other metadata categories can  
 370 be found in Supplemental Materials S6 (figures) and S7 (tables of  $\beta$  values). In accordance with  
 371 the results for spatiotemporal metadata, supplementary analyses indicated that metadata for  
 372 collection year (posterior mean slope = -0.133, 95% Credible Interval: -0.233 - -0.034; Table S1,  
 373 Figure S6-11A) and preservative used (posterior mean slope = -0.111, 95% HPD: -0.218 - -  
 374 0.009; Figure S6-9A) were significantly less likely to be recovered from INSDC, publications,  
 375 and/or online repositories with increasing age of a BioProject. Furthermore, and as with  
 376 spatiotemporal metadata, the probability that responding authors provided additional metadata

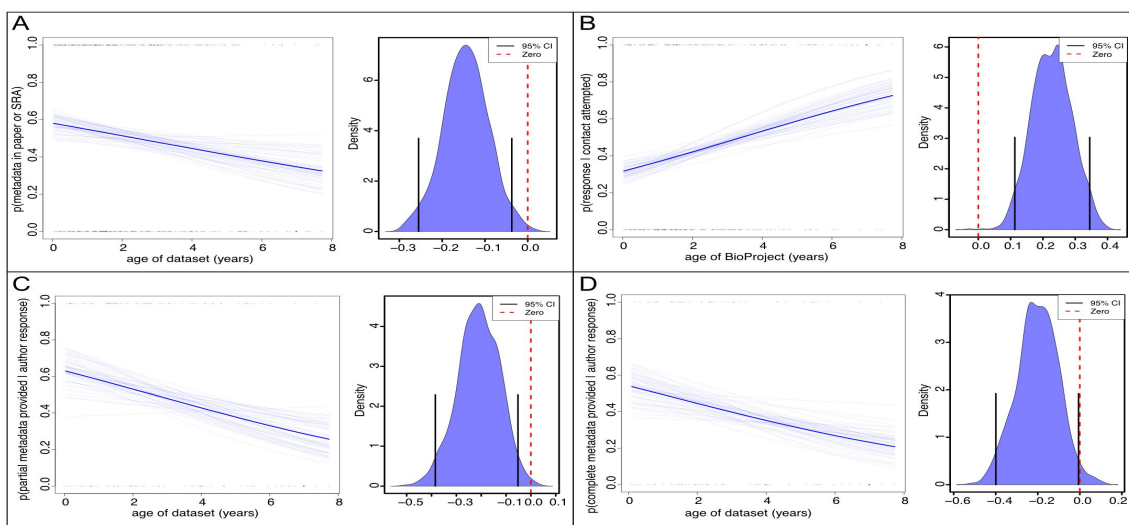


Figure 4. The effect of BioProject age on the probability of recovering spatiotemporal metadata. Density plots depict posterior distribution for log(Odds Ratio), with black lines showing 95% highest posterior density (HPD). (A) Probability that metadata were found in the INSDC or in the INSDC or associated papers and repositories. (B) Probability of receiving a reply from BioProject authors to our contact email. (C) Probability of receiving any amount of additional metadata for year OR coordinates OR locality. (D) Probability of receiving metadata for year AND coordinates OR locality for a majority of BioSamples. 95% HPD intervals exclude 0 in all panels.

377 for georeferences (decimalLatitude and decimalLongitude; posterior mean slope = -0.151, 95%  
378 Credible Interval: -0.386 – -0.05; Figure S6-2C), collection year (posterior mean slope = -0.174,  
379 95% Credible Interval: -0.363 – 0.000; Figure S6-11C), and preservative used (posterior mean  
380 slope = -0.438, 95% Credible Interval: -0.873 – -0.081; Figure S6-9C) was significantly greater  
381 for younger BioProjects. The provisioning of permit information followed this same trend  
382 (although marginally insignificant, posterior mean slope = -0.555, 95% Credible Interval: -1.31 -  
383 0.003; Figure S6-5C), suggesting these metadata are relatively available within the personal data  
384 management system of authors.

385  
386 Concerningly and counter to our result for spatiotemporal metadata, supplementary analyses  
387 indicated that metadata for habitat (Table S1, posterior mean slope = 0.141, 95% Credible  
388 Interval: 0.006 - 0.285; Figure S6-6C) and environmental medium (posterior mean slope = 0.176,  
389 95% Credible Interval: 0.016 - 0.355; Figure S6-5C) were less frequently recovered from  
390 INSDC, publications and/or repositories for younger BioProjects. The reasons for these results  
391 are unclear, but may indicate a decline in author attention and appreciation of organismal natural  
392 history. Retrieval of these metadata through author contact had no relationship with BioProject  
393 age.

## 394 **Discussion**

395  
396  
397 With this distributed datathon, we have demonstrated that crucial metadata can be restored  
398 for many genomic investigations of wild organisms. However, our analyses show that metadata  
399 are more difficult to recover the longer we wait, and many are locked in non-standard formats.  
400 Because the great majority of publicly available genomic datasets lack important metadata  
401 (Toczydlowski et al., 2021), they are not findable, accessible, interoperable nor reusable (FAIR;  
402 Wilkinson et al., 2016). Only genomic data that are FAIR will allow systematic monitoring of  
403 the fundamental layer of biodiversity (Hoban et al., 2021), and enable assertions regarding  
404 provenance for informing CBD Nagoya Protocol obligations. Our results illustrate that: (1)  
405 metadata availability is dependent on type: location, publication and habitat metadata are much  
406 more available or inferable than metadata about permits and preservatives; (2) with considerable  
407 time and paid effort, it is possible to recover some of these important metadata from the non-  
408 standardized and non-machine-readable formats in which they are currently being stored; and (3)  
409 while metadata archival practices may be improving incrementally, genomic metadata are  
410 subject to the same decay processes demonstrated for other types of scientific data (Pope et al.,  
411 2015; Vines et al., 2014).

412  
413 There are likely multiple factors underlying the observed metadata decay. First, it is not  
414 surprising that older metadata are less likely to have been archived. Metadata archival practices  
415 are gradually improving, with more metadata being recorded into the INSDC, research papers,  
416 and online repositories such as Data Dryad (Figure 4A). This is consistent with increasing  
417 acknowledgement that these metadata are relevant and important to future research. However,  
418 the rate of metadata archival is apparently not keeping up with the rapid growth of genomic  
419 datasets (see Figure 1 of Toczydlowski et al., 2021) and it is certainly not closing the gap.  
420 Second, we found that authors of recent SRA datasets were significantly less likely to reply to  
421 our queries than those of older datasets (Figure 4B), although the overall response rate of 45%  
422 was comparable to previous studies (Vines et al., 2014, 2013). This result may indicate that

423 recent SRA depositions are part of ongoing research projects for which authors are unwilling to  
424 share metadata for fear of getting “scooped” by others working on similar research questions. It  
425 is also true that younger authors are more likely to leave science than older authors (Reithmeier  
426 et al., 2019) and thus may no longer support their publications. Similarly, there may be a cohort  
427 effect in which authors of older studies are more established in their careers and have more time,  
428 and/or are more aware of increasing expectations around FAIR data, and thus more willing to  
429 communicate and share. Finally, of the authors that did reply, there was a significant decrease  
430 with the age of the BioProject in whether partial or complete spatiotemporal metadata were  
431 provided (Figure 4C,D), suggesting that if metadata are not properly archived to public  
432 repositories, they are subject to being lost over time, as previously highlighted for morphological  
433 data (Vines et al., 2014).

434  
435 Taken together, our results support assertions that the current research system overly weights  
436 publications and citations, while underweighting scientific openness and transparency (S. W.  
437 Davies et al., 2021; McNutt et al., 2016; Nosek et al., 2015). Changing the system will likely  
438 require a combination of carrots and sticks (Whitlock, 2011). Carrots can take the form of citable  
439 data publications (Dimitrova et al., 2021), recognition of open data practices by hiring,  
440 promotion and tenure committees, or commendations from professional societies or departments  
441 (Roche, Kruuk, Lanfear, & Binning, 2015; Roche et al., 2014). Sticks in the form of open  
442 metadata mandates, must come from journals (Sibbett, Rieseberg, & Narum, 2020, Gareth  
443 Jenkins, pers. comm.), funding agencies, and data repositories, which all have a responsibility to  
444 respond to the needs of the research community (Lin et al., 2020). While we applaud the  
445 INSDC’s new spatiotemporal metadata annotation policy requiring country of origin metadata  
446 and their adoption of the MIxS metadata standards, we call for greater mandated spatial  
447 resolution to include at least a locality name or spatial coordinates (Table 1) with appropriate  
448 uncertainty or additional terms (such as Darwin Core’s coordinateUncertaintyInMeters,  
449 informationWithheld; Wicczorek et al., 2012) to protect endangered species or sovereignty of  
450 Indigenous Peoples (Hudson et al., 2020; McCartney et al., 2022).

451  
452 Our datathon provided an unparalleled opportunity to train graduate students in the  
453 importance of proper data curation, and to raise awareness that almost every dataset has a  
454 potential for re-use. We suggest that training in data curation and metadata usage should be part  
455 of reproducible research training in every science graduate program, with emphasis on avoiding  
456 some of the metadata practices that hinder metadata recovery described in Box 1. Additionally,  
457 “datathons” such as that undertaken here could help to close the metadata gap in the short term,  
458 as they are very cost effective. If we assume a mean cost of sequencing of USD 50 per  
459 BioSample (and ignore the much higher, additional cost of sample collection and processing),  
460 this datathon rescued over USD 2.1 million worth of genomic sequence data for future research  
461 purposes. Co-authors of this paper spent about 2,300 hours on this metadata retrieval effort,  
462 which, if valued at an average wage of USD 19 per hour, yields a return on investment of nearly  
463 4,700%, with average costs of remediating a BioSample or BioProject at USD 1.05 and USD 110  
464 respectively. But ultimately, datathons are a stopgap solution.

465  
466 Going forward, the entire biodiversity genomics research community should give the same  
467 priority to sharing metadata that they have given to sharing primary data, because it is only the  
468 metadata that make primary data FAIR. From a process standpoint, the collection of metadata

469 should begin at the time of sampling, with the assignment of a globally unique identifier (GUID)  
470 to the actual material sample. This identifier, which should be assigned as early as possible after  
471 collection, serves as the root to which all subsequent derived products could be linked in an  
472 extended specimen cloud to establish clear provenance and thereby prevent duplication of data or  
473 effort (N. Davies et al., 2021; Lendemer et al., 2020). Through the use of GUIDs, both physical  
474 and digital products of the sample (digital sequence information, but also DNA or RNA  
475 extractions, subsamples, images, video, audio, CT scans, measurements of morphology, traits,  
476 gut contents, parasites, and other related data and associated metadata) will be linked to their  
477 material sample GUID to provide an extensive, holistic metadata cloud that can be used to better  
478 inform current research endeavors as well as create additional data-intensive research pathways.  
479 GEOME (Deck et al., 2017; Riginos et al., 2020) is an example of an easy-to-use “metadata  
480 broker” platform that can provide spreadsheet templates with definitions that can be filled in  
481 offline when the sample is collected. It can then mint a GUID for any sample that is added to it,  
482 and then harvest the INSDC accession numbers for genomic reads that are submitted to the SRA  
483 through GEOME, thereby maintaining permanent links between the sample metadata and  
484 genomic data.

485  
486 The challenge then is to integrate these metadata downstream into databases (such as  
487 INSDC) which describe data derived from the sample. INSDC enables such linkages to other  
488 metadata platforms through the use of both Structured Voucher  
489 (<https://www.ncbi.nlm.nih.gov/biollections/docs/faq/>) and Linkout  
490 (<https://www.ncbi.nlm.nih.gov/projects/linkout/>) facilities for both Nucleotide and SRA (through  
491 their corresponding BioSample record) datasets respectively (e.g.  
492 <https://www.ncbi.nlm.nih.gov/nuccore/KC825472>). Through these linkages, metadata  
493 corresponding to the original material sample can be tied to the resulting sequence(s) to both  
494 validate the metadata associated with the sequence record as well as provide updated information  
495 should specimens be reidentified or georeferenced after the lodging of the sequence with INSDC.  
496 Using the INSDC as a long-term repository for metadata about the sample may not make sense,  
497 in part because researchers who submit the sequences to INSDC have sole editing rights to the  
498 sequence record and it is currently quite difficult for others (such as the collections who hold the  
499 vouchers) to keep the INSDC metadata up to date or add additional information. Thus, the  
500 integration of these metadata from an upstream source somewhat negates the necessity for this  
501 information to be duplicated by the sequence depositor and ensures that the metadata are  
502 constantly up to date. This not only supports open, reproducible science (Buckner, Sanders,  
503 Faircloth, & Chakrabarty, 2021) but also exemplifies the Findable and Accessible principles of  
504 FAIR data (Wilkinson et al., 2016).

505  
506 What this piecemeal system currently lacks, however, is support for data Interoperability and  
507 Reusability. This is because of the siloed nature of the data and our inability to compile it into a  
508 single resource for machine readability, data manipulation or downstream use. This shortcoming  
509 is being addressed through various initiatives such as the Extended Specimen Network (ESN;  
510 Lendemer et al., 2020; Thiers et al., 2021), the Digital Extended Specimen  
511 (<https://dissco.tech/2020/03/31/what-is-a-digital-specimen/>), the Distributed System for  
512 Scientific Collections (DiSCCo; <https://www.dissco.eu/>), iSamples (N. Davies et al., 2021) and  
513 others. Such a system would require all actors in the data landscape (researchers, collections,  
514 data aggregators, publishers, etc.) to utilize and publish resolvable GUIDs on all specimens,

515 datasets and products of research to make these linkages possible, and thereby create an  
516 extensive online network of knowledge, and increase the potential for scientific research  
517 questions to be answered.

518  
519 We join others in calling for ambitious policy goals that safeguard genetic diversity and  
520 scientific practices to enable this (Des Roches et al., 2021; Díaz et al., 2020; Laikre et al., 2020).  
521 Swift collective action is required to protect all levels of global biodiversity, and the first step  
522 towards protecting the evolutionary health of eukaryotic species worldwide is to close the  
523 metadata gap highlighted here. Simultaneously, conservation geneticists, molecular ecologists  
524 and evolutionary biologists must engage with global biodiversity assessment programs to ensure  
525 genomic data can be collected, interpreted and archived appropriately (Brodersen & Seehausen,  
526 2014). Several exemplary international networks (e.g. GEOBON Genetic Composition Working  
527 Group, IUCN Conservation Genetics Specialist Group, and EU COST Action Genetic  
528 Biodiversity Knowledge for Ecosystem resilience [GBiKE]) have already made a case for  
529 protecting the genetic diversity of all species (Laikre et al., 2020), proposed indicators to gauge  
530 progress toward goals (Hoban et al., 2020; Laikre et al., 2020). These groups have asserted their  
531 rationale for these changes to stakeholders in policy documents, providing essential clarity in the  
532 use of genetic data, and reporting against targets (Hoban et al., 2021). These actions and  
533 advances encourage the uptake of genetic diversity monitoring by national authorities and  
534 international bodies. The vision for many of these biodiversity monitoring networks is to develop  
535 agile pipelines that intake raw biodiversity data and produce outputs that can directly inform  
536 conservation policies and decisions (Hoban et al., 2021). Yet, without appropriate archival of  
537 genomic data that includes the spatiotemporal metadata, we will be unable to deliver on the  
538 promise of genetic diversity monitoring.

539  
540 The GEOME datathon enabled 13 students and 7 academics from 15 institutions and 4  
541 countries to take account of the growing metadata gap for genomics data and begin to remediate  
542 it. The serendipity of being able to run a remote, distributed datathon due to travel restrictions  
543 and funding reallocation forced by COVID-19, in a time when Indigenous rights, biodiversity  
544 conservation and the value of genetic diversity have been front of mind, has not been lost on the  
545 participants. While our efforts have just begun to address the growing metadata gap, it is our  
546 hope that most researchers will start to ensure the FAIRness of their genomic data and metadata  
547 before or upon publication, thereby honoring the work that went into creating it and providing  
548 limitless opportunities for reuse of their data to help answer the important scientific questions of  
549 the future.

550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560

## 561 Acknowledgements

562 This effort arose from an Evolution in Changing Seas Research Coordination Network  
563 (RCN) working group (NSF-OCE-1764316, Katie Lotterhos) and was funded by the Diversity of  
564 the Indo-Pacific Network RCN (NSF-DEB-1457848, to R.J.T.). We gratefully thank all of the  
565 authors who took the time to provide helpful responses to our metadata inquiries. We also thank  
566 Neil Davies, Chris Meyer, Beth Davis and Kiersey Nielsen for their input.

## 567 Data Availability Statement

568 The metadata that were recovered by the datathon are available from the Genomic Observatories  
569 Metadatabase (GEOME; [geome-db.org](https://geome-db.org)) <https://geome-db.org/workbench/project-overview?projectId=305>. Relevant code for analyses in this paper may be found at  
570 [https://github.com/ericcrandall/geome\\_metadatathon1](https://github.com/ericcrandall/geome_metadatathon1). The meta-metadata for BioProjects that  
571 were determined to be relevant to the datathon are in Supplemental Materials S4.  
572

573

## 574 Literature Cited

575

- 576 Allendorf, F. W. (2017). Genetics and the conservation of natural populations: Allozymes to genomes.  
577 *Molecular Ecology*, 26(2), 420–430. doi: 10.1111/mec.13948
- 578 Baetscher, D. S., Anderson, E. C., Gilbert-Horvath, E. A., Malone, D. P., Saarman, E. T., Carr, M. H., &  
579 Garza, J. C. (2019). Dispersal of a nearshore marine fish connects marine reserves and adjacent  
580 fished areas along an open coast. *Molecular Ecology*, 1(2), 0148–13. doi: 10/gf8rkk
- 581 Blanchet, S., Prunier, J. G., Paz-Vinas, I., Saint-Pé, K., Rey, O., Raffard, A., ... Dubut, V. (2020). A river  
582 runs through it: The causes, consequences, and management of intraspecific diversity in river  
583 networks. *Evolutionary Applications*, 13(6), 1195–1213. doi: 10.1111/eva.12941
- 584 Brauman, K. A., Garibaldi, L. A., Polasky, S., Aumeeruddy-Thomas, Y., Brancalion, P. H. S., DeClerck,  
585 F., ... Verma, M. (2020). Global trends in nature's contributions to people. *Proceedings of the*  
586 *National Academy of Sciences*, 117(51), 32799–32805. doi: 10.1073/pnas.2010473117
- 587 Brodersen, J., & Seehausen, O. (2014). Why evolutionary biologists should get seriously involved in  
588 ecological monitoring and applied biodiversity assessment programs. *Evolutionary Applications*,  
589 7(9), 968–983. doi: 10.1111/eva.12215
- 590 Buckner, J. C., Sanders, R. C., Faircloth, B. C., & Chakrabarty, P. (2021). The critical importance of  
591 vouchers in genomics. *ELife*, 10, e68264. doi: 10.7554/eLife.68264
- 592 Buttigieg, P. L., Morrison, N., Smith, B., Mungall, C. J., Lewis, S. E., & the ENVO Consortium. (2013).  
593 The environment ontology: Contextualising biological and biomedical entities. *Journal of*  
594 *Biomedical Semantics*, 4(1), 43. doi: 10.1186/2041-1480-4-43
- 595 Buttigieg, P. L., Pafilis, E., Lewis, S. E., Schildhauer, M. P., Walls, R. L., & Mungall, C. J. (2016). The  
596 environment ontology in 2016: Bridging domains with increased scope, semantic density, and  
597 interoperability. *Journal of Biomedical Semantics*, 7(1), 57. doi: 10.1186/s13326-016-0097-6
- 598 Castañeda-Álvarez, N. P., Khoury, C. K., Achicanoy, H. A., Bernau, V., Dempewolf, H., Eastwood, R. J.,  
599 ... Toll, J. (2016). Global conservation priorities for crop wild relatives. *Nature Plants*, 2(4), 1–6.  
600 doi: 10.1038/nplants.2016.22
- 601 CBD. (2020). *Global Biodiversity Outlook 5*. Montreal: United Nations. Retrieved from United Nations  
602 website: <https://www.cbd.int/gbo/gbo5/publication/gbo-5-en.pdf>
- 603 Cheng, S. H., Gold, M., Rodriguez, N., & Barber, P. H. (2021). Genome-wide SNPs reveal complex fine  
604 scale population structure in the California market squid fishery (*Doryteuthis opalescens*).  
605 *Conservation Genetics*, 22(1), 97–110. doi: 10.1007/s10592-020-01321-2

- 606 Clark, J. S. (2010). Individuals and the Variation Needed for High Species Diversity in Forest Trees.  
607 *Science*, 327(5969), 1129–1132. doi: 10.1126/science.1183506
- 608 Cochrane, G., Karsch-Mizrachi, I., Takagi, T., & INSDC. (2016). The International Nucleotide Sequence  
609 Database Collaboration. *Nucleic Acids Research*, 44(D1), D48–D50. doi: 10.1093/nar/gkv1323
- 610 Cowell, C., Paton, A., Borrell, J. S., Williams, C., Wilkin, P., Antonelli, A., ... Kersey, P. J. (2022). Uses  
611 and benefits of digital sequence information from plant genetic resources: Lessons learnt from  
612 botanical collections. *PLANTS, PEOPLE, PLANET*, 4(1), 33–43. doi: 10.1002/ppp3.10216
- 613 Davies, N., Deck, J., Kansa, E. C., Kansa, S. W., Kunze, J., Meyer, C., ... Lehnert, K. (2021). Internet of  
614 Samples (iSamples): Toward an interdisciplinary cyberinfrastructure for material samples.  
615 *GigaScience*, 10(5). doi: 10.1093/gigascience/giab028
- 616 Davies, S. W., Putnam, H. M., Ainsworth, T., Baum, J. K., Bove, C. B., Crosby, S. C., ... Bates, A. E.  
617 (2021). Promoting inclusive metrics of success and impact to dismantle a discriminatory reward  
618 system in science. *PLOS Biology*, 19(6), e3001282. doi: 10.1371/journal.pbio.3001282
- 619 Deck, J., Gaither, M. R., Ewing, R., Bird, C. E., Davies, N., Meyer, C., ... Crandall, E. D. (2017). The  
620 Genomic Observatories Metadatabase (GeOMe): A new repository for field and sampling event  
621 metadata associated with genetic samples. *PLoS Biology*, 15(8), e2002925. doi: 10/gbqt2p
- 622 Des Roches, S., Pendleton, L. H., Shapiro, B., & Palkovacs, E. P. (2021). Conserving intraspecific  
623 variation for nature’s contributions to people. *Nature Ecology & Evolution*, 5(5), 574–582. doi:  
624 10.1038/s41559-021-01403-5
- 625 Díaz, S., Zafra-Calvo, N., Purvis, A., Verburg, P. H., Obura, D., Leadley, P., ... Zanne, A. E. (2020). Set  
626 ambitious goals for biodiversity and sustainability. *Science*, 370(6515), 411–413. doi:  
627 10.1126/science.abe1530
- 628 Dimitrova, M., Meyer, R., Buttigieg, P. L., Georgiev, T., Zhelezov, G., Demirov, S., ... Penev, L. (2021).  
629 A streamlined workflow for conversion, peer review, and publication of genomics metadata as  
630 omics data papers. *GigaScience*, 10(5). doi: 10.1093/gigascience/giab034
- 631 Exposito-Alonso, M., Booker, T. R., Czech, L., Gillespie, L., Hateley, S., Kyriazis, C. C., ... Zess, E.  
632 (2022). *Genetic diversity loss in the Anthropocene*. 6. doi: 10.1126/science.abn5642
- 633 Field, D., Garrity, G., Gray, T., Morrison, N., Selengut, J., Sterk, P., ... Wipat, A. (2008). The minimum  
634 information about a genome sequence (MIGS) specification. *Nature Biotechnology*, 26(5), 541–  
635 547. doi: 10/dms
- 636 Gaither, M. R., Gkafas, G. A., de Jong, M., Sarigol, F., Neat, F., Regnier, T., ... Hoelzel, A. R. (2018).  
637 Genomics of habitat choice and adaptive evolution in a deep-sea fish. *Nature Ecology &*  
638 *Evolution*, 2(4), 680–687. doi: 10.1038/s41559-018-0482-x
- 639 Halewood, M., Lopez Noriega, I., Ellis, D., Roa, C., Rouard, M., & Sackville Hamilton, R. (2018). Using  
640 Genomic Sequence Information to Increase Conservation and Sustainable Use of Crop Diversity  
641 and Benefit-Sharing. *Biopreservation and Biobanking*, 16(5), 368–376. doi:  
642 10.1089/bio.2018.0043
- 643 Hendricks, S., Anderson, E. C., Antao, T., Bernatchez, L., Forester, B. R., Garner, B., ... Luikart, G.  
644 (2018). Recent advances in conservation and population genomics data analysis. *Evolutionary*  
645 *Applications*, 11(8), 1197–1211. doi: 10/gd5nc2
- 646 Hinchliff, C. E., Smith, S. A., Allman, J. F., Burleigh, J. G., Chaudhary, R., Coghill, L. M., ... Cranston,  
647 K. A. (2015). Synthesis of phylogeny and taxonomy into a comprehensive tree of life.  
648 *Proceedings of the National Academy of Sciences*, 112(41), 12764–12769. doi: 10/f7vmqg
- 649 Hoban, S., Archer, F. I., Bertola, L. D., Bragg, J. G., Breed, M. F., Bruford, M. W., ... Hunter, M. E.  
650 (2022). Global genetic diversity status and trends: Towards a suite of Essential Biodiversity  
651 Variables (EBVs) for genetic composition. *Biological Reviews, Early Online*. doi:  
652 10.1111/brv.12852
- 653 Hoban, S., Bruford, M., D’Urban Jackson, J., Lopes-Fernandes, M., Heuertz, M., Hohenlohe, P. A., ...  
654 Laikre, L. (2020). Genetic diversity targets and indicators in the CBD post-2020 Global  
655 Biodiversity Framework must be improved. *Biological Conservation*, 248, 108654. doi:  
656 10/gg32bm

- 657 Hoban, S., Bruford, M. W., Funk, W. C., Galbusera, P., Griffith, M. P., Grueber, C. E., ... Vernesi, C.  
658 (2021). Global Commitments to Conserving and Monitoring Genetic Diversity Are Now  
659 Necessary and Feasible. *BioScience*, (biab054). doi: 10.1093/biosci/biab054
- 660 Hudson, M., Garrison, N. A., Sterling, R., Caron, N. R., Fox, K., Yracheta, J., ... Carroll, S. R. (2020).  
661 Rights, interests and expectations: Indigenous perspectives on unrestricted access to genomic  
662 data. *Nature Reviews Genetics*, 21(6), 377–384. doi: 10.1038/s41576-020-0228-x
- 663 Kardos, M., Taylor, H. R., Ellegren, H., Luikart, G., & Allendorf, F. W. (2016). Genomics advances the  
664 study of inbreeding depression in the wild. *Evolutionary Applications*, 9(10), 1205–1218. doi:  
665 10.1111/eva.12414
- 666 Kumar, A., Anju, T., Kumar, S., Chhapekar, S. S., Sreedharan, S., Singh, S., ... Lim, Y. P. (2021).  
667 Integrating Omics and Gene Editing Tools for Rapid Improvement of Traditional Food Plants for  
668 Diversified and Sustainable Food Security. *International Journal of Molecular Sciences*, 22(15),  
669 8093. doi: 10.3390/ijms22158093
- 670 Laikre, L. (2010). Genetic diversity is overlooked in international conservation policy implementation.  
671 *Conservation Genetics*, 11(2), 349–354. doi: 10.1007/s10592-009-0037-4
- 672 Laikre, L., Hoban, S., Bruford, M. W., Segelbacher, G., Allendorf, F. W., Gajardo, G., ... Vernesi, C.  
673 (2020). Post-2020 goals overlook genetic diversity. *Science*, 367(6482), 1083. doi: 10/ggnjfx
- 674 Leigh, D. M., Hendry, A. P., Vázquez-Domínguez, E., & Friesen, V. L. (2019). Estimated six per cent  
675 loss of genetic variation in wild populations since the industrial revolution. *Evolutionary  
676 Applications*, 12(8), 1505–1512. doi: 10.1111/eva.12810
- 677 Lendemer, J., Thiers, B., Monfils, A. K., Zaspel, J., Ellwood, E. R., Bentley, A., ... Aime, M. C. (2020).  
678 The Extended Specimen Network: A Strategy to Enhance US Biodiversity Collections, Promote  
679 Research and Education. *BioScience*, 70(1), 23–30. doi: 10.1093/biosci/biz140
- 680 Letunic, I., & Bork, P. (2021). Interactive Tree Of Life (iTOL) v5: An online tool for phylogenetic tree  
681 display and annotation. *Nucleic Acids Research*, 49(W1), W293–W296. doi:  
682 10.1093/nar/gkab301
- 683 Lewin, H. A., Robinson, G. E., Kress, W. J., Baker, W. J., Coddington, J., Crandall, K. A., ... Zhang, G.  
684 (2018). Earth BioGenome Project: Sequencing life for the future of life. *Proceedings of the  
685 National Academy of Sciences*, 115(17), 4325–4333. doi: 10/gdh5vz
- 686 Liggins, L., Hudson, M., & Anderson, J. (2021). Creating space for Indigenous perspectives on access  
687 and benefit-sharing: Encouraging researcher use of the Local Contexts Notices. *Molecular  
688 Ecology*, 30(11), 2477–2482. doi: 10.1111/mec.15918
- 689 Lin, D., Crabtree, J., Dillo, I., Downs, R. R., Edmunds, R., Giarretta, D., ... Westbrook, J. (2020). The  
690 TRUST Principles for digital repositories. *Scientific Data*, 7(1), 144. doi: 10/ggwrtj
- 691 Lou, R. N., Jacobs, A., Wilder, A. P., & Therkildsen, N. O. (2021). A beginner’s guide to low-coverage  
692 whole genome sequencing for population genomics. *Molecular Ecology*, 30(23), 5966–5993. doi:  
693 10.1111/mec.16077
- 694 Marden, E., Abbott, R. J., Austerlitz, F., Ortiz-Barrientos, D., Baucom, R. S., Bongaerts, P., ... Rieseberg,  
695 L. H. (2021). Sharing and reporting benefits from biodiversity research. *Molecular Ecology*,  
696 30(5), 1103–1107. doi: 10.1111/mec.15702
- 697 McCartney, A. M., Anderson, J., Liggins, L., Hudson, M. L., Anderson, M. Z., TeAika, B., ... Phillippy,  
698 A. M. (2022). Balancing openness with Indigenous data sovereignty: An opportunity to leave no  
699 one behind in the journey to sequence all of life. *Proceedings of the National Academy of  
700 Sciences*, 119(4). doi: 10.1073/pnas.2115860119
- 701 McNutt, M., Lehnert, K., Hanson, B., Nosek, B. A., Ellison, A. M., & King, J. L. (2016). Liberating field  
702 science samples and data. *Science*, 351(6277), 1024–1026. doi: 10.1126/science.aad7048
- 703 Michonneau, F., Brown, J. W., & Winter, D. J. (2016). rotl: An R package to interact with the Open Tree  
704 of Life data. *Methods in Ecology and Evolution*, 7(12), 1476–1481. doi: 10.1111/2041-  
705 210X.12593

- 706 Miraldo, A., LI, S., Borregaard, M. K., Flórez-Rodríguez, A., Gopalakrishnan, S., Rizvanovic, M., ...  
707 Nogués-Bravo, D. (2016). An Anthropocene map of genetic diversity. *Science*, 353(6307), 1532–  
708 1535. doi: 10/f9ccgr
- 709 Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D., Breckler, S. J., ... Yarkoni, T.  
710 (2015). Promoting an open research culture. *Science*, 348(6242), 1422–1425. doi: 10/gcpzwn
- 711 Pinsky, M. L., & Palumbi, S. R. (2014). Meta-analysis reveals lower genetic diversity in overfished  
712 populations. *Molecular Ecology*, 23(1), 29–39. doi: 10.1111/mec.12509
- 713 Pope, L. C., Liggins, L., Keyse, J., Carvalho, S. B., & Riginos, C. (2015). Not the time or the place: The  
714 missing spatio-temporal link in publicly available genetic data. *Molecular Ecology*, 24(15), 3802–  
715 3809. doi: 10.1111/mec.13254
- 716 Prada, C., Hanna, B., Budd, A. F., Woodley, C. M., Schmutz, J., Grimwood, J., ... Medina, M. (2016).  
717 Empty Niches after Extinctions Increase Population Sizes of Modern Corals. *Current Biology*,  
718 26(23), 3190–3194. doi: 10.1016/j.cub.2016.09.039
- 719 Quattrini, A. M., Wu, T., Soong, K., Jeng, M.-S., Benayahu, Y., & McFadden, C. S. (2019). A next  
720 generation approach to species delimitation reveals the role of hybridization in a cryptic species  
721 complex of corals. *BMC Evolutionary Biology*, 19(1), 116. doi: 10.1186/s12862-019-1427-y
- 722 Raffard, A., Santoul, F., Cucherousset, J., & Blanchet, S. (2019). The community and ecosystem  
723 consequences of intraspecific diversity: A meta-analysis. *Biological Reviews*, 94(2), 648–661.  
724 doi: 10.1111/brv.12472
- 725 Reithmeier, R., O’Leary, L., Zhu, X., Dales, C., Abdulkarim, A., Aquil, A., ... Zou, C. (2019). The  
726 10,000 PhDs project at the University of Toronto: Using employment outcome data to inform  
727 graduate education. *PLOS ONE*, 14(1), e0209898. doi: 10/ggfnvv
- 728 Reusch, T. B. H., Ehlers, A., Hämmerli, A., & Worm, B. (2005). Ecosystem recovery after climatic  
729 extremes enhanced by genotypic diversity. *Proceedings of the National Academy of Sciences*,  
730 102(8), 2826–2831. doi: 10.1073/pnas.0500008102
- 731 Riginos, C., Crandall, E. D., Liggins, L., Gaither, M. R., Ewing, R. B., Meyer, C., ... Deck, J. (2020).  
732 Building a global genomics observatory: Using GEOME (the Genomic Observatories  
733 Metadatabase) to expedite and improve deposition and retrieval of genetic data and metadata for  
734 biodiversity research. *Molecular Ecology Resources*, 20(6), 1458–1469. doi: 10/ghf4wp
- 735 Roche, D. G., Kruuk, L. E. B., Lanfear, R., & Binning, S. A. (2015). Public Data Archiving in Ecology  
736 and Evolution: How Well Are We Doing? *PLoS Biology*, 13(11), e1002295-12. doi:  
737 10.1371/journal.pbio.1002295
- 738 Roche, D. G., Lanfear, R., Binning, S. A., Haff, T. M., Schwanz, L. E., Cain, K. E., ... Kruuk, L. E. B.  
739 (2014). Troubleshooting Public Data Archiving: Suggestions to Increase Participation. *PLoS*  
740 *Biology*, 12(1), e1001779. doi: 10/ggvr85
- 741 Scholz, A. H., Freitag, J., Lyal, C. H. C., Sara, R., Cepeda, M. L., Cancio, I., ... Overmann, J. (2022).  
742 Multilateral benefit-sharing from digital sequence information will support both science and  
743 biodiversity conservation. *Nature Communications*, 13(1), 1086. doi: 10.1038/s41467-022-  
744 28594-0
- 745 Schriml, L. M., Chuvochina, M., Davies, N., Eloe-Fadrosch, E. A., Finn, R. D., Hugenholtz, P., ... Walls,  
746 R. (2020). COVID-19 pandemic reveals the peril of ignoring metadata standards. *Scientific Data*,  
747 7(1), 188. doi: 10.1038/s41597-020-0524-5
- 748 Sibbett, B., Rieseberg, L. H., & Narum, S. (2020). The Genomic Observatories Metadatabase. *Molecular*  
749 *Ecology Resources*, 20(6), 1453–1454. doi: 10/ghf4wm
- 750 Thiers, B., Bates, J., Bentley, A. C., Ford, L. S., Jennings, D., Monfils, A. K., ... Pandey, J. L. (2021).  
751 Implementing a Community Vision for the Future of Biodiversity Collections. *BioScience*, 71(6),  
752 561–563. doi: 10.1093/biosci/biab036
- 753 Toczydlowski, R. H., Liggins, L., Gaither, M. R., Anderson, T. J., Barton, R. L., Berg, J. T., ... Crandall,  
754 E. D. (2021). Poor data stewardship will hinder global genetic diversity surveillance. *Proceedings*  
755 *of the National Academy of Sciences*, 118(34), e2107934118. doi: 10.1073/pnas.2107934118

756 Vines, T. H., Albert, A. Y. K., Andrew, R. L., Débarre, F., Bock, D. G., Franklin, M. T., ... Rennison, D.  
757 J. (2014). The Availability of Research Data Declines Rapidly with Article Age. *Current Biology*,  
758 24(1), 94–97. doi: 10/qpm  
759 Vines, T. H., Andrew, R. L., Bock, D. G., Franklin, M. T., Gilbert, K. J., Kane, N. C., ... Yeaman, S.  
760 (2013). Mandated data archiving greatly improves access to research data. *The FASEB Journal*,  
761 27(4), 1304–1308. doi: 10.1096/fj.12-218164  
762 Whitlock, M. C. (2011). Data archiving in ecology and evolution: Best practices. *Trends in Ecology*  
763 *& Evolution*, 26(2), 61–65. doi: 10.1016/j.tree.2010.11.006  
764 Wieczorek, J., Bloom, D., Guralnick, R., Blum, S., Döring, M., Giovanni, R., ... Vieglais, D. (2012).  
765 Darwin Core: An evolving community-developed biodiversity data standard. *PLoS ONE*, 7(1),  
766 e29715. doi: 10.1371/journal.pone.0029715  
767 Wilder, A. P., Palumbi, S. R., Conover, D. O., & Therkildsen, N. O. (2020). Footprints of local adaptation  
768 span hundreds of linked genes in the Atlantic silverside genome. *Evolution Letters*, 4(5), 430–  
769 443. doi: 10.1002/evl3.189  
770 Wilkinson, M. D., Dumontier, M., Aalbersberg, Ij. J., Appleton, G., Axton, M., Baak, A., ... Mons, B.  
771 (2016). Comment: The FAIR Guiding Principles for scientific data management and stewardship.  
772 *Scientific Data*, 3(1), 1–9. doi: 10/bdd4  
773 Willette, D. A., Allendorf, F. W., Barber, P. H., Barshis, D. J., Carpenter, K. E., Crandall, E. D., ... Seeb,  
774 J. E. (2014). So, you want to use next-generation sequencing in marine systems? Insight from the  
775 Pan-Pacific Advanced Studies Institute. *Bulletin Of Marine Science*, 90(1), 79–122. doi:  
776 10.5343/bms.2013.1008  
777 Winter, D. J. (2017). rentrez: An R package for the NCBI eUtils API. *The R Journal*, 9(2), 520–526. doi:  
778 10.32614/rj-2017-058  
779 Yilmaz, P., Kottmann, R., Field, D., Knight, R., Cole, J. R., Amaral-Zettler, L., ... Glöckner, F. O.  
780 (2011). Minimum information about a marker gene sequence (MIMARKS) and minimum  
781 information about any (x) sequence (MIxS) specifications. *Nature Biotechnology*, 29(5), 415–  
782 420. doi: 10.1038/nbt.1823  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806

807  
808  
809  
810  
811  
812

**Box 1.** Summary of metadata practices encountered during the datathon that hindered metadata recovery and recommended practices to improve future usability of samples and genetic sequence data. In general, we recommend that authors use metadata software (GEOME: [geome-db.org](http://geome-db.org), COPO: [copo-project.org](http://copo-project.org), or museum database software such as Specify) to organize and archive their sample metadata.

Practice	Challenge	Solution
INSDC materialSampleID (SRA) does not match any sample identifiers in associated scientific publication(s)	Metadata external to the INSDC cannot be assigned to genetic samples or metadata are associated with the wrong sample	Use consistent, persistent, globally unique sample identifiers (e.g. Darwin Core <i>materialSampleID</i> ) across data repositories and publications; if sample identifiers are not consistent, provide an explicit cross-reference table in all associated publications and data repositories
Large amounts of metadata are only available in associated publications in PDF format & lack consistent formatting	Metadata cannot be programmatically converted to standard table formats (e.g. entries formatted column by row, rather than row by column) & require time-consuming manual extraction	Provide metadata in comma- or tab-delimited files (.csv or .txt) using standard column headers (i.e. terms suggested by Darwin Core, MiXs, or GEOME or COPO) & associated vocabulary, where possible
Specimens, BioSamples or metadata are deposited in a biodiversity collection (e.g., museum, herbarium, biobank, zoo), but biodiversity collection accession numbers are not provided in the associated publications or INSDC	Biodiversity collection record searches can be time-consuming & may not yield enough information to link samples in collection databases back to INSDC databases	Use consistent sample identifiers across all databases & publications; or provide a cross-reference table in associated publication(s) that links biodiversity collection accessions to INSDC materialSampleIDs and identifies the biocollection repository by name. The Darwin Core standard accommodate multiple terms
Associated publication references previous publication(s) for details about the sample metadata	Time consuming & challenging to track citation trail back to metadata in original/earlier publication. Sample metadata or identifiers may be absent or in inconsistent formats across associated publications	Compile & include relevant metadata from previous publications in supplementary materials for the publication linked to the INSDC BioProject; if needed, include a column flagging whether data are new to the present study or originated from another source (& identify that source)
Sample collection geospatial coordinates or location name withheld in order to protect endangered species, sensitive habitat, or Indigenous sovereignty	Sample lacks spatial metadata	Provide imprecise geospatial coordinates & use large, defined <i>coordinateUncertainty</i> to maintain local anonymity of sensitive collection sites. Provide additional comments using Darwin Core <i>informationWithheld</i>
Codes used to abbreviate sample collection locations are inconsistent or hard to find throughout publication & related materials	Sample collection locations cannot be determined or require time-consuming manual curation	Include site codes in the sample identifier ( <i>materialSampleID</i> ). Use consistent site codes throughout associated publications & provide a key with codes & geospatial coordinates in associated publications
Sample collection dates are a range or a season (e.g., winter 2017-2018)	Sample collection date may not be identifiable to a specific year; unclear which samples were collected in which date range	Denote/report which year <u>each</u> sample was collected in (dates that also include month & day are ideal)

No metadata on BioSample relevance to genetic diversity of wild populations/species provided	Unclear which BioSamples (if any) were collected from wild populations versus e.g. brood stock, laboratory stock, domesticated species, artificial selection experiments, non-wild collections in seed banks	Provide metadata denoting which BioSamples were collected from wild populations, using Darwin Core term <i>establishmentMeans</i>
Metadata provided for some but not all BioSamples	Some BioSamples lack metadata, unclear why metadata are incomplete	Provide metadata for all BioSamples or list a specific reason for missing metadata (e.g. not collected, metadata lost, sample excluded from study due to misidentification) using Darwin Core <i>informationWithheld</i> .
Sampling location provided, but only at a coarse geographic scale i.e., state, province, or country name	Sample lacks spatial metadata at a resolution useful for future monitoring & macrogenetic research questions	Provide geospatial coordinates for sample collection locations. Specific place, state, & country names can be helpful additions to confirm the geospatial coordinates are correct (& to programmatically filter by broader geographic locations)

813  
814  
815  
816  
817

**Table 1:** Alphabetized list of required (in bold) and recommended metadata terms for individual organisms and/or derived tissues or DNA sequences included in the datathon. Square brackets in the definition column denote the metadata standard from which the definition comes. Terms with multiple definitions are in order of decreasing specificity. The importance column indicates which terms support the identification of Indigenous provenance and can therefore inform Access and Benefit-Sharing (ABS), and those that can support sample or Digital Sequence Information (DSI) re-use in conservation, according to the study approach definitions of Leigh et al. (2021). Class I studies generate new sequence data, requiring precise information regarding the spatiotemporal context of the collected sample, a unique materialSampleID, as well as the preservative the tissue is held in; Class II studies compile genetic diversity values from published studies, generally requiring less precise spatiotemporal information, but this needs to be associated with a publication (associatedReferences); Class III studies re-analyse digital sequence information, or derived genetic data, requiring precise spatiotemporal information, and a unique materialSampleID. Depending on the objective of re-use, habitat and environmental\_medium may also be important for sample/DSI re-use in conservation. Controlled vocabularies refer to standardized lists of acceptable entries, often defined by a standards organization.

Term	Definition	Importance	Controlled Vocabulary	Example
associatedReferences	[GEOME] <sup>1</sup> Any associated publications/references pertaining to this individual or its derivative tissues or sequences. The first place it was published is particularly relevant. DOIs in format: <a href="https://doi.org/10.1007/s10530-007-9196-8">https://doi.org/10.1007/s10530-007-9196-8</a> . Multiple DOIs separated by  . [Darwin Core] <sup>2</sup> A list (concatenated and separated) of identifiers (publication, bibliographic reference, global unique identifier, URI) of literature associated with the Occurrence.	Indigenous provenance, Class II	None	" <a href="https://doi.org/10.1111/j.1365">https://doi.org/10.1111/j.1365</a> <a href="https://doi.org/10.5343/bms.20">https://doi.org/10.5343/bms.20</a>
<b>coordinateUncertaintyInMeters</b>	[Darwin Core] The horizontal distance (in meters) from the given decimalLatitude and decimalLongitude describing the smallest circle containing the whole of the locality where the sample could possibly have come from. Value empty if the uncertainty is unknown, cannot be estimated, or is not applicable (because there are no coordinates). Zero is not a valid value for this term.	Class I, III	None	1 km = "1000"

country	[Darwin Core] The name of the country or major administrative unit or exclusive economic zone (for marine samples) in which the locality occurs.	Indigenous provenance, Class II	ISO 3166-1	"Indonesia"
decimalLatitude	[Darwin Core] The geographic latitude (in decimal degrees, using the spatial reference system given in geodeticDatum) of the geographic center of a Location. Positive values are north of the Equator, negative values are south of it. Legal values lie between -90 and 90, inclusive.	Indigenous provenance, Class I, II, III	None	"-6.147183"
decimalLongitude	[Darwin Core] The geographic longitude (in decimal degrees, using the spatial reference system given in geodeticDatum) of the geographic center of a Location. Positive values are east of the Greenwich Meridian, negative values are west of it. Legal values lie between -180 and 180, inclusive.	Indigenous provenance, Class I, II, III	None	"105.46326"
derivedDataFilename	[GEOME] A list (concatenated and separated with  ) of the file names for datasets that include data derived from this tissue that are accessible via the 'derivedDataURI'. Could be a compressed archive.	Class III	None	"SDM_snps.tar.gz"
derivedDataFormat	[GEOME] A list (concatenated and separated with  ) of the dataset formats relating to the 'derivedDataType' that include data derived from this tissue.	Class III	{microsatellites, Sequence alignment, SNPs,OTUs,ASVs, Other}	"SNPs"
derivedDataType	[GEOME] A list (concatenated and separated with  ) of the dataset types that include data derived from this tissue.	Class III	{genepop,FASTA,VCF, nexus,PHYLIP, structure,Other}	"VCF"
derivedDataURI	[GEOME] A URI (preferably a DOI in this format: <a href="https://doi.org/10.1007/s10530-007-9196-8">https://doi.org/10.1007/s10530-007-9196-8</a> ) for any datasets that include data derived from this tissue. Multiple URIs/DOIs can be separated by  .	Class III	None	<a href="https://doi.org/10.5061/dryad.">"https://doi.org/10.5061/dryad."</a>
environmental_medium	[MIXS] <sup>3</sup> Terms that identify the material displaced by the entity at time of sampling. Recommend subclasses of environmental material [ENVO:00010483]. Multiple terms can be separated by pipes e.g.: a duck might displace fresh water air.		ENVO <sup>4</sup> Environmental N ENVO_00010483	"sea water"

<b>habitat</b>	[MIXS: Broad-scale environmental context] In this field, report which major environmental system your sample or specimen came from. The systems identified should have a coarse spatial grain, to provide the general environmental context of where sampling was done. [Darwin Core] A category or description of the habitat in which the Event occurred.		ENVO Biomes: ENVO_00000428	"marine benthic biome"
<b>locality</b>	[Darwin Core] The specific name or description of the site or place where the sample was taken as given by the original researchers. This would be the place name that appears in a table next to the coordinates, or the labels for sampling sites on a map. Less specific geographic information can be provided in other geographic terms (continentOcean, country, stateProvince, island). This term may contain information modified from the original to correct perceived errors or standardize the description.	Indigenous provenance, Class I, II, III	None	"Rakata"
<b>materialSampleID</b>	[GEOME] The collector's specimen number. This number must be unique among the IDs within the sheet. [Darwin Core]. An identifier for the MaterialSample (as opposed to a particular digital record of the material sample). In the absence of a persistent global unique identifier, construct one from a combination of identifiers in the record that will most closely make the materialSampleID globally unique.	Class I, III	None	"Rakata_1190.01"
<b>permitInformation</b>	[GEOME] Information regarding the permit acquired to collect this sample. At least the permit number and issuing authority. Multiple values separated by	Indigenous provenance, Class I	None	"Indonesian Institute of Sciences Permit #'s 1187/ SU/KS/2006 and 04239/SU.3/KS/2006"
<b>preservative</b>	[GEOME] Preservative used on the specimen.	Class I	GEOME List of preserva	"95% Ethanol"
<b>yearCollected</b>	[Darwin Core] The year the collecting event took place	Class I, II, III	None	March 24, 2006 = "2006"

1. Deck et al. (2017); Riginos et al. (2020)

2. Wiczorek et al. (2012)
3. Yilmaz et al. (2011)
4. Buttigieg et al. (2013, 2016)