

1 **Title:** Best practices for genetic and genomic data archiving

2

3 **Authors:** Deborah M. Leigh^{1*}, Amy G. Vandergast², Margaret E. Hunter³, Eric Crandall⁴, W.
4 Chris Funk⁵, Colin J. Garroway⁶, Sean Hoban⁷, Sara J. Oyler-McCance⁸, Christian Rellstab¹,
5 Gernot Segelbacher⁹, Chloe Schmidt¹⁰, Ella Vázquez-Domínguez¹¹, Ivan Paz-Vinas^{5,12}

6

7 ***corresponding author:** deborah.leigh@wsl.ch

8

9 **Deborah M. Leigh**

10 1) Swiss Federal Research Institute WSL, Birmensdorf, Switzerland

11 ORCID: 0000-0003-3902-2568

12 deborah.leigh@wsl.ch

13

14 **Amy G. Vandergast**

15 2) U.S. Geological Survey, Western Ecological Research Center, 4165 Spruance Road,
16 Suite 200, San Diego CA, 92101, USA

17 ORCID: 0000-0002-7835-6571

18 avandergast@usgs.gov

19

20 **Margaret E. Hunter**

21 3) U.S. Geological Survey, Wetland & Aquatic Research Center, 7920 NW 71st Street,
22 Gainesville, Florida 32653, USA

23 ORCID: 0000-0002-4760-9302

24 mhunter@usgs.gov

25

26 **Eric D. Crandall**

27 4) Department of Biology, Pennsylvania State University, 208 Mueller Laboratory,
28 University Park, PA 16802, USA

29 ORCID: 0000-0001-8580-3651

30 ecrandall@psu.edu

31

32 **W. Chris Funk**

33 5) Department of Biology, Graduate Degree Program in Ecology, Colorado State
34 University, 1878 Campus Delivery, Fort Collins, CO 80523-1878, USA

35 ORCID: 0000-0002-6466-3618

36 Chris.Funk@colostate.edu

37

38 **Colin J. Garroway**

39 6) Department of Biological Sciences, University of Manitoba, Winnipeg, Manitoba,
40 Canada

41 ORCID: 0000-0002-0955-0688

42 colin.garroway@umanitoba.ca

43

44 **Sean Hoban**
45 7) Center for Tree Science, The Morton Arboretum, Lisle, IL, 60532, USA *and*
46 Committee on Evolutionary Biology, University of Chicago, Chicago, IL, 60637, USA
47 ORCID: 0000-0002-0348-8449
48 shoban@mortonarb.org
49

50 **Sara J. Oyler-McCance**
51 8) U.S. Geological Survey, Fort Collins Science Center, 2150 Centre Avenue, Building
52 C, Fort Collins, CO, 80526, USA
53 ORCID: 0000-0003-1599-8769
54 soyler@usgs.gov
55

56 **Christian Rellstab**
57 1) Swiss Federal Research Institute WSL, Birmensdorf, Switzerland
58 ORCID: 0000-0002-0221-5975
59 christian.rellstab@wsl.ch
60

61 **Gernot Segelbacher**
62 9) Wildlife Ecology and Management, University Freiburg, Tennenbacher Str. 4, 79106
63 Freiburg, Germany
64 ORCID: 0000-0002-8024-7008
65 gernot.segelbacher@wildlife.uni-freiburg.de
66

67 **Chloé Schmidt**
68 10) German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig,
69 Puschstrasse 4, Leipzig, 04103, Germany
70 ORCID 0000-0003-2572-4200
71 chloe.schmidt@idiv.de
72

73 **Ella Vázquez-Domínguez**
74 11) Departamento de Ecología de la Biodiversidad, Instituto de Ecología, Universidad
75 Nacional Autónoma de México, Coyoacán, Ciudad de México, 04510, México;
76 ORCID: 0000-0001-6131-2014
77 evazquez@ecologia.unam.mx
78

79 **Ivan Paz-Vinas**
80 12) Universite Claude Bernard Lyon 1, LEHNA UMR 5023, CNRS, ENTPE, F-69622,
81 Villeurbanne, France
82 ORCID: 0000-0002-0043-9289
83 ivan.paz-vinas@univ-lyon1.fr
84
85
86

87

88

Best practices for genetic and genomic data archiving

89

90 **Abstract:** Genetic and genomic data are collected for a vast array of scientific and
91 applied purposes. Despite mandates for public archiving, data are typically used only by
92 the generating authors. The reuse of genetic and genomic datasets remains uncommon
93 because it is difficult, if not impossible, due to non-standard archiving practices and lack
94 of contextual metadata. But as the new field of macrogenetics is demonstrating, if genetic
95 data and their metadata were more accessible and FAIR compliant, they could be reused
96 for many additional purposes. We discuss the main challenges with existing genetic and
97 genomic data archives, and suggest best practices for archiving genetic and genomic
98 data. Recognising that this is a longstanding issue due to little formal data management
99 training within the fields of ecology and evolution, we highlight steps that research
100 institutions and publishers could take to improve data archiving.

101

102 **Main:**

103 *A brief overview on the history and value of genetic data in ecology and evolution*

104 Synthesis of Open Data (publicly archived data, free to reuse) is a powerful tool that is
105 increasingly being used to test pressing questions in ecology and evolution. However, it
106 remains common for valuable datasets to be forgotten after a single use¹⁻³. This is a
107 missed opportunity and hinders scientific progress. Producing scientific data is often
108 expensive and time-consuming. Furthermore, most data have numerous potential
109 applications beyond their original use⁴.

110 Public archiving of genetic and genomic sequence data (hereafter 'genetic data')
111 became standard practice in the 1980s⁵ but notably, archiving associated *metadata* (data
112 that describe the sampling event, sample, and other derived data), still remains
113 discretionary. Nevertheless, genetic data repositories were some of the earliest Open
114 Data projects (e.g. National Center for Biotechnology Information 'NCBI' GenBank⁶) and
115 continue to arise in response to the increasing needs and volume of genomic data
116 archiving (e.g.^{7,8}).

117 Open population genetic data have now accumulated to the point where data can
118 be synthesized across broad scales (e.g., in macrogenetics^{9,10}), rapidly advancing the
119 fields of ecology and evolution by enabling characterization of global biodiversity patterns
120 and genetic diversity trends^{10,11}. Sequences within NCBI are now frequently reused for
121 taxonomic assignments, facilitating species discovery and environmental DNA method
122 development¹². Accessible raw genomic read datasets have also become central to
123 bioinformatic teaching and analysis development (e.g.¹³). Yet the future reuse potential
124 of genetic data extends further, as an abundance of unattempted and unknown uses
125 remain. Vitally, data reuse is one way for countries to help prevent genetic diversity loss
126 through reporting the required genetic indicators of the Convention of Biological Diversity

127 (CBD) Kunming-Montreal Global Biodiversity Framework (e.g. headline indicator A.4^{14–}
128 ¹⁶).

129 Despite the long history and growing abundance of Open Genetic Data, journal
130 Open Data policies^{21–25}, and increasing awareness of the FAIR principles (Findable,
131 Accessible, Interoperable, and Reusable²⁶), there are still numerous issues that inhibit
132 comprehensive reuse. These range from issues general to ecology and evolution such
133 as inaccessible private data (comprehensively addressed elsewhere^{2,3,27–29}), to field-
134 specific issues we outline below. In this perspective, we suggest guidance on archiving
135 different types of genetic data and their associated metadata. We discuss additional steps
136 or infrastructure needed to improve the status quo. Ultimately, our goal is to prevent data
137 loss and facilitate data reuse.

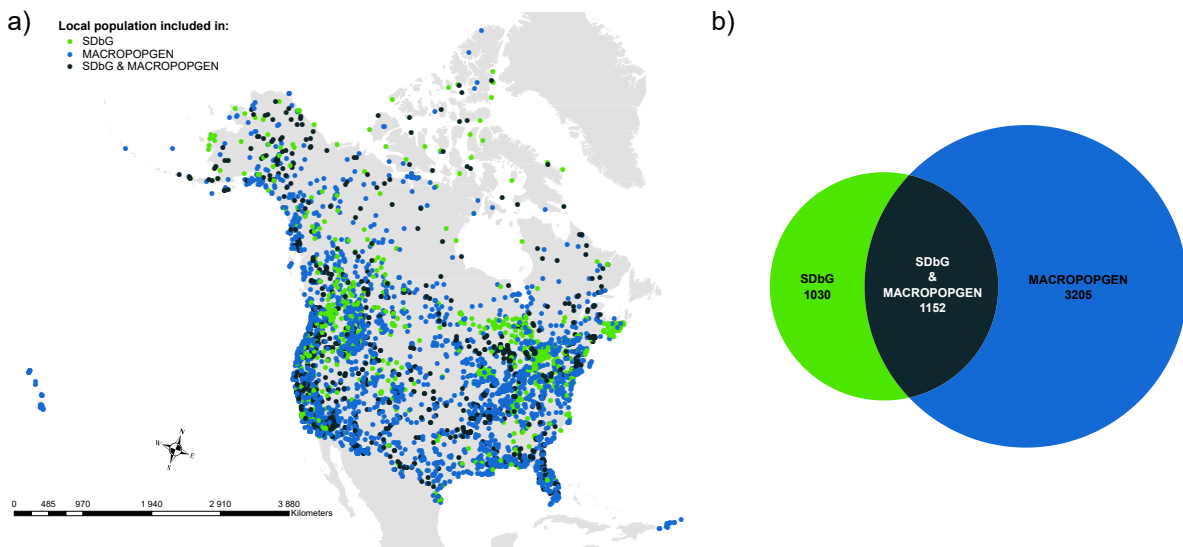


Figure 1: Estimating the unknown number of 'missing' datasets in open repositories. Spatial distribution (a) and proportion (b) of overlapping datasets available in two recently published macrogenetic databases for mammals, birds, reptiles, and amphibians from the USA and Canada: 1) MACROPOPGEN¹⁷, consists of georeferenced microsatellite-derived summary statistics extracted from published articles; 2) SDbG^{18–20} consists of raw microsatellite genotype datasets extracted directly from open repositories. After cross checking, only 21.38% of the data entries were found in both databases (black dots), while 59.5% were found exclusively in MACROPOPGEN (blue dots). Low overlap suggests a large proportion of genetic studies included in MACROPOPGEN did not have findable publicly archived data and/or sufficient metadata, and thus were not usable in the SDbG.

Recommendations for Genetic and Genomic Data Archiving

a) Gene sequences

- FASTA format
- use established databases (INSDC)
- archive all haplotypes (not only unique/new)
- minimum sample metadata (Box 2)

b) Microsatellite genotypes

- STRUCTURE single line format
- use any FAIR compliant database
- link data to publication
- use keywords on archive to enhance findability
- minimum sample metadata (Box 2)

c) Genomic read data

- demultiplexed read and least processed VCF files
- use established database (INSDC for read data) and FAIR compliant database (VCF)
- consistent sample names across all files
- minimum sample metadata (Box 2)
- comprehensive code archiving and enriched VCF header

d) Indigenous peoples' data

- collaboratively write data management plan(s)
- remove technology barriers when sharing data
- publically archive data carefully, respecting DSI and commercial value
- use an agreed-upon archiving database
- include essential sample metadata with clear reuse permissions/contacts in the archive (Box 2)

The need to improve genetic data archiving practices in ecology and evolution

The FAIR guiding principles are the foundation for good, transparent, and reproducible science. The clearest evidence for this is that Open Data have been used to identify scientific misconduct (e.g.³⁰). Datasets are often financed by taxpayers, making public releases of any data an ethical - often mandated³¹ - obligation to ensure the full value is obtained. Genomic data production has been particularly expensive, costing hundreds of millions of dollars that would have to be re-spent to regenerate data without archive enrichment³². Furthermore, due to the rapid pace of biodiversity loss³³, which includes allelic loss and population extirpation³⁴, recollecting data may be impossible, rendering existing data irreplaceable. Existing genetic data consequently represent an irreplaceable baseline against which to compare future measurements (i.e. for monitoring³⁵). Poor archiving represents a significant loss of time, resources, and opportunities (Figure 1). It is also an unnecessary ethical footprint when data re-generation requires animal handling. Further, increasing genetic data volumes and associated storage energy costs add urgency to the need to improve the archiving standards of genetic data and their metadata by establishing *best practices*.

Researchers can benefit in many ways from publicly archiving genetic data. Open datasets can enhance the scholarly recognition of individual research efforts, because data releases with DOIs and data papers can be cited (e.g. [MacroPopGen¹⁷](#)). As in many disciplines, data papers (e.g. [Darwin Tree of Life's Genome Notes](#)) are becoming increasingly popular. Synthesis can also test previously unanswerable big-picture questions in genetics, benefiting researchers through advancement of their field¹⁰.

Best practices for FAIR genetic data archiving

The most widely available genetic data types in molecular ecology and evolution are: a) barcoding/gene sequences (e.g. mitochondrial cytochrome oxidase, the major

Figure 2: Summary of Best Practice recommendations for each type of genetic data.

176

177

178

179

180 histocompatibility complex), b) microsatellite genotypes, and c) genomic read data (i.e.
181 raw high throughput sequences and Single Nucleotide Polymorphisms “SNPs”; Figure 2).
182 The latter two come in a constellation of software-specific formats^{36,37} and, due to lack of
183 standardization, repositories contain most of these formats. While format conversion tools
184 exist (PGDspider³⁶; Formatomatic³⁸; vcftools³⁹; plink⁴⁰; adegenet⁴¹), conversions are
185 time-consuming and often need customization. Mastering each file format also requires
186 specialist knowledge. Consequently, the lack of a standard archived format limits
187 interoperability and reusability. Due to fundamental differences in data types, file sizes
188 and formats used, a single genetic data file format is unrealistic. However, a single file
189 type for each data type is possible and would be a significant advancement.

190 Unlike other genetic data types, gene sequences are somewhat standardized on
191 archives as FASTA files, and we recommend maintaining this approach (Figure 2a).
192 However, many gene sequences lack essential metadata to allow their reuse. It is
193 important that authors include the minimum metadata needed to interpret their archived
194 data (briefly summarised in Box 1), otherwise archives are challenging to reuse (e.g. non-
195 georeferenced sequences in GenBank⁴²).

196 For microsatellite data (Figure 2b), we suggest archiving in the popular and flexible
197 STRUCTURE format⁴³. STRUCTURE input files can handle genotype data of varying
198 ploidy and have a simple format that is conducive to editing in R, spreadsheet software
199 or on the command line, without generating formatting errors. This file format can house
200 metadata (Box 1), as well as marker information (i.e. presence of recessive alleles, inter-
201 marker distances, phase information). We note that there are also variations within the
202 STRUCTURE line format, notably the 1 vs 2 lines *per* individual. Either is suitable for
203 archiving as both are accepted by conversion tools like PGDSpider³⁶. However, we
204 recommend use of the single line format to maximize similarity with VCF (“Variant Call
205 Format”) files.

206 Genomic data are often mandated to be publicly archived as raw read data on
207 INSDC servers (“INSDC” International Nucleotide Sequence Database Collaboration)⁴⁴,
208 or as aligned BAM files for model organisms⁴⁵. Raw read data can be highly variable,
209 ranging from completely unprocessed files containing several individuals, demultiplexed
210 read files, cleaned files (i.e. with low-quality reads or individuals removed), to error-
211 corrected files (e.g. in ancient DNA)⁴⁶. In contrast to microsatellite data, the variable
212 archiving of genomic data means basic error removal, sample delimitation, and genotype
213 calls are not expected to be present in archived data. We recommend sequencing read
214 data are archived as demultiplexed read files to ensure that separation from key barcode
215 metadata will not render the read data unusable (Figure 2c). This also facilitates archiving
216 of individual sample metadata which has higher reuse potential than study-level
217 metadata. A bioinformatic pipeline can also be challenging to reproduce because there
218 are chronic issues surrounding open code archiving that make it hard to identify what
219 parameters were applied, tool versions used, or even to have access to custom scripts

220 (further detailed in^{27,47}). Even if a pipeline is accessible, version changes of reference
221 genomes or software programs quickly make reproducing it impossible. Thus, we
222 recommend archiving genotype files (detailed below) in addition to demultiplexed
223 sequencing reads to improve the Open Data compliance and reusability of genomic data.
224 Processed VCF files containing genotype calls (or genotype
225 likelihoods/probabilities) are standard for genomic analyses and we recommend archiving
226 them in parallel with raw read files (notably, this is not possible on INSDC). Although such
227 processed files are not currently widely archived, the practice is becoming more common.
228 Standardization of exactly which variant file is archived also needs consideration.
229 Maximum reusability would be achieved if the archived file represents the least processed
230 genotype (i.e. unfiltered and pruned only for basic errors like technical faults or
231 contamination, with file headers retained to indicate the bioinformatic steps applied and
232 versions used). Notably, archiving VCF files could allow reuse by researchers, managers,
233 or benefit-holders without High-Performance Computing capabilities (e.g. for
234 conservation). Furthermore, VCF archiving would reduce the non-negligible energy,
235 storage and ultimately emissions costs associated with reanalysing raw genomic data⁴⁸.

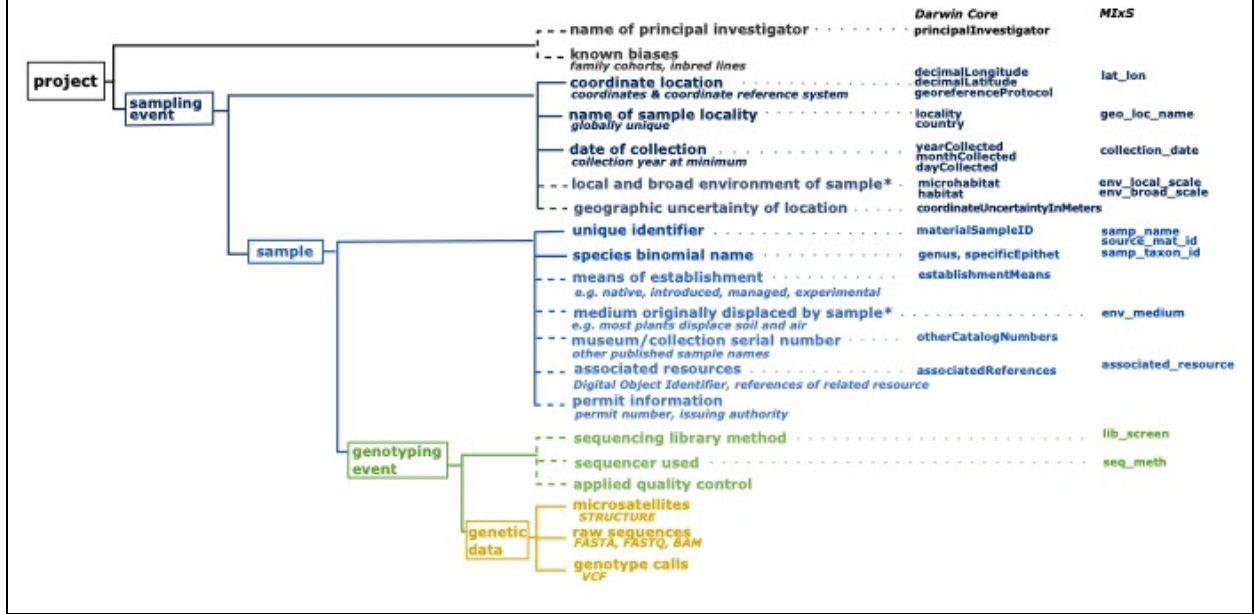
Box 1: Archiving metadata for genetic studies

In genetic studies metadata describing where, when, how and by whom genotype or sequence data were created are invaluable for making data FAIR. There are currently two genomic metadata standards: the Darwin Core standard for biodiversity data⁴⁹ and the Minimum Information about any(x) Sequence (MIxS) standard⁵⁰. Both have cross-mapped terms that overlap⁵¹. The box Figure summarizes term overlap.

Metadata can be viewed in a hierarchical manner based on how they were created starting from the sampling events moving down to the genotyping event. While metadata will vary by sample type and project goals, at a minimum, we suggest that authors provide the required (solid lines) and recommended (dotted lines) categories in the Figure below to improve publicly archived genetic data reuse potential. Terms denoted with * could use controlled vocabulary from the Environment ontology (“ENVO”⁵²). To report metadata not covered here, we also recommend using Darwin Core or MixS standards terms to guarantee FAIRness. Sensitive data should be withheld to ensure it is protected. This can be denoted with the terms “informationWithheld”, “dataGeneralizations” or “coordinateUncertaintyInMeters”. Note, metadata fields might not be adapted for ancient DNA, for which metadata related to sampling events generally does not reflect the age and the environmental conditions of the sampled individual before death. Geological context names may be needed.

An often-overlooked key to FAIR metadata lies in the sample identifier (materialSampleID or samp_name in Darwin Core and MIxS respectively). These identifiers should be unique within the project, and identical between the genetic data and the metadata. To protect genetic data being separated from metadata and help spot errors during complex uploads to databases, we recommend introducing metadata-enriched unique sample names enriched with core metadata like species name, coordinates and/or sampling year (i.e. Capra.ibex_46.97.8.25 or

Capra.ibex.pilatus.2014). Samples that need to be linked across files or studies must be named consistently. We also discourage archiving metadata in separate file repositories from genotype or sequencing data. If unavoidable, we recommend that metadata are stored in a simple table (CSV or text format) with clearly labelled columns (e.g. using MIxS or Darwin Core terms), and consistent sample identifiers. To aid automated retrieval, authors can ensure metadata are machine actionable by avoiding the use of symbols, special characters, and/or colour-based cell codes.



236

237 *Missing metadata renders most archived data useless*

238 Metadata are a crucial aspect of ensuring genetic data adhere to the FAIR
 239 principles²⁶ because data context vastly increases name potential reuses. It was historically
 240 standard to only include taxonomic metadata (species and genus) in genetic data
 241 archives. Recent mandates have expanded to include the country of collection and
 242 collection dates^{53,54}. However, neither is sufficient for comprehensive reuse. The
 243 minimum required and recommended metadata are shown in Box 1, without which
 244 archived data are often functionally useless and could foster incorrect inferences.

245 We note that metadata from at-risk species, and species that are commercially
 246 valuable or desirable, may need to be withheld or obscured to protect them⁵⁵. The most
 247 concerning data for such species is location data – coordinates or specific habitat
 248 descriptions that would allow public access to these specimens or locations. This is
 249 particularly pertinent in genetic studies as individual-level high accuracy coordinates are
 250 often collected. Recommended best practices for generalizing sensitive species
 251 occurrence or geographic metadata involve masking, controlled access, or not reporting
 252 such metadata^{56,57}. These limitations should be accepted because it is essential that
 253 Open Data do not infringe on privacy, benefit sharing, or species protection efforts⁵⁵.

254

255 *Special considerations when working internationally, with sensitive species, or*
256 *Indigenous communities*

257 Genetic data and metadata archiving are also key to benefit sharing, including the
258 rights of local communities and local scientists to access data generated from specimens
259 within their country or region. Emerging benefit-sharing requirements, such as those put
260 forth in the Nagoya Protocol and being developed by the CBD (e.g. Digital Sequence
261 Information ‘DSI’⁵⁸), are becoming a legal requirement⁵⁹. This is particularly pertinent to
262 ecology and evolution where researchers often work transnationally^(e.g.60,61) and steps are
263 needed to overcome parachute science⁶². Importantly, while the release of genetic data
264 from sensitive or commercially relevant species could facilitate conservation or
265 evolutionary understanding, protecting their potential commercial value (e.g.
266 pharmaceutical or agricultural) should be carefully considered during archiving.

267 For work involving Indigenous communities, the CARE principles (Collective
268 benefits, Authority to control, Responsibility, and Ethics^{63,64}) could be considered in
269 archiving and reuse of genetic data (Figure 2d). What steps researchers should follow will
270 be situation-specific and developed in conjunction with the benefit holders⁶⁵. Data-
271 generating authors can include specific benefit-sharing statements in publications and in
272 data archives to ensure data sovereignty is upheld. The statement should contain
273 contextual metadata, for instance provenance information, community names, and also
274 clearly outline community-granted permissions for reuse and circulation. Links to
275 biocultural notices created by researchers and endorsement labels issued by Indigenous
276 peoples can also be stored as sample metadata. Authors should also be aware that
277 respecting data sovereignty can influence the data repository used (e.g. Aotearoa
278 Genomic Data Repository, which allows for access only once permission is granted⁶⁶),
279 data storage location⁶⁷, and archiving formats (e.g. archiving VCF files is key to limiting
280 technology barriers which otherwise may inhibit reuse by Indigenous researchers or
281 communities). When reusing Indigenous owned genetic data, researchers should also
282 discuss or co-design planned reanalyses with Indigenous communities. Attribution and
283 citation of the original datasets in resulting manuscripts and dissemination of results to
284 the communities involved could further help ensure that cultural authority and sovereignty
285 over reused data remain recognized (e.g.⁶⁸), and that data are not reused inappropriately.
286 Overall, best practices involve improving respect and compliance with the rights
287 Indigenous peoples have for agency over their data.

288

289 *Key features of data repositories for FAIR data*

290 Currently, genotype data are often stored in generalist Open Data repositories
291 (DRYAD, Zenodo, and increasingly FigShare). However, genetic data can quickly get lost
292 among many other data types archived, where researchers can find everything from non-
293 peer-reviewed ecological survey data (e.g.⁶⁹) to violent crime statistics (e.g.⁷⁰). Local rules
294 and repository fees make it impossible to advocate for a single database for all genotype

295 data. While there are interoperable search platforms that facilitate simultaneous cross-
296 repository search (e.g. [DataONE](#)), their functionality is not guaranteed and database
297 linking has failed in the past (Chloé Schmidt *pers. comm.*). Thus, there is a need for a
298 free⁷¹ inter-government supported public database specifically for archiving genotype
299 data (e.g. microsatellites calls, SNPs, etc). Note that for data involving Indigenous
300 communities, particularly that with restricted reuse, special repositories may be needed
301 to prevent automated retrieval and improper reuse⁶⁶.

302 In lieu of a dedicated repository, researchers can take a few key steps to ensure
303 genotypic data findability. At a minimum, the repositories used must clearly link data to
304 publications and provide citable DOIs. Key metadata fields (Box 1) should be included in
305 the database description to aid findability. Marker type (e.g. “microsatellite” or “SNP”) and
306 key geographical descriptors (e.g. “Kruger National Park”) can also be used as keywords
307 to aid search functions. Researchers could also link genotypic data to “metadatabases”
308 that track samples through metadata and can facilitate upload to the SRA (INSDC
309 BioProjects and BioSamples⁷²; Genomic Observatories MetaDatabase ‘GEOME’⁷³;
310 Collaborative OPen Omics ‘COPO’⁷⁴).

311 Researchers can also request new features within existing databases to facilitate
312 data accessibility. The Web of Science’s “associated data” link is a notable advance⁷⁵, as
313 is the increased mandatory metadata (sample location, collection date) for BioSample
314 packages and European Nucleotide Archive archives^{54,76}. An additional feature, which
315 would benefit multiple disciplines, is the implementation of an automatic identifier for data
316 associated with retracted articles. While datasets from retracted papers should remain
317 online as important records of technical errors, or even fraud, as of writing fraudulent data
318 remain on data repositories with no notice of retraction (e.g.⁷⁷). Similarly, data found to
319 be erroneous remains on sequence databases (e.g. GenBank⁷⁸). Collectively, this poses
320 a huge challenge for studies based on automated data reuse. Archives could flag data
321 with “concerns raised”, “under evaluation”, “technical errors present”, or “retracted” for
322 clarity. Researchers could also benefit from an easy and anonymous way to notify
323 database curators if they encounter incomplete non-FAIR compliant archives to improve
324 database integrity.

325

326 *The role of scientific institutions and journals in improving data archiving*

327 Funding bodies could take on a greater responsibility to ensure cross-discipline
328 FAIR data archiving (Figure 3). While mandating Open Data has undoubtedly increased
329 data accessibility (e.g.^{79,80}), funding bodies could also support researchers with data
330 management plans and confirm implementation (e.g.⁸¹), check data accessibility, pay
331 archiving fees, and offer general data archiving educational resources or training. For
332 genetic projects, funding bodies can ensure sufficient time is budgeted for archiving
333 because it can take several days.

334 University libraries or research organizations could support Open Data by hiring
 335 data “stewards” or “librarians” familiar with ecological and evolutionary genetic data. Data
 336 stewards can help write data management plans, identify suitable databases for genotype
 337 file archiving, and ensure dataset longevity through best practice compliance⁸². Few
 338 ecology and evolution researchers receive formal training in Open Data or data archiving.
 339 Thus they could also offer data management education (e.g. short courses and training)
 340 for both students and career scientists (e.g.⁸³).

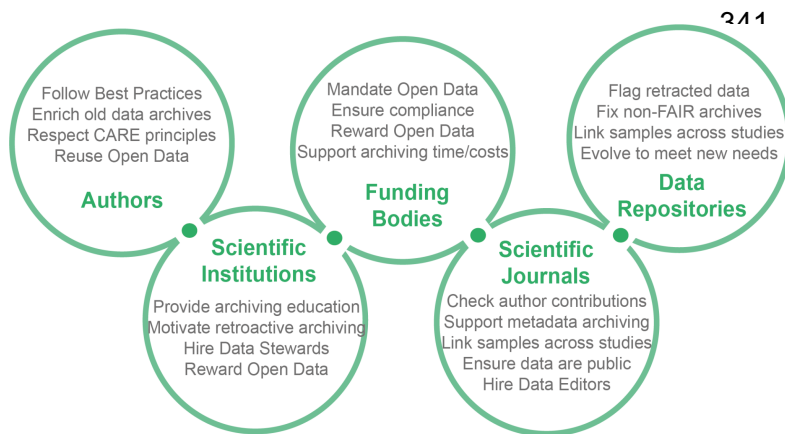


Figure 3: A brief summary of the roles that would improve Open Genetic Data

355 activated when papers are *in press*. However, this makes it impossible for journals to
 356 assess data presence and support archiving. A shift to making data accessible *upon*
 357 *submission* to journals is thus needed, particularly at the resubmission stage when papers
 358 are close to acceptance⁸⁴. Journal data editors could also check data archives to ensure
 359 files are not corrupt, contain the reported number of markers or loci, and contain basic
 360 metadata. Scientists concerned with data being accessed prior to publication should note
 361 that several databases offer non-public shareable links. Alternatively, journals could make
 362 the final acceptance dependent on evidence of FAIR data compliance⁸⁴.

363 Journals could also improve genetic data reuse potential by establishing a
 364 mandatory table of standardized metadata terms (see Box 1). However, we note that
 365 versioning issues may arise if metadata are in multiple places (e.g. supplementary
 366 materials, INSDC, Dryad). Journal data editors could help prevent this by ensuring data
 367 from the same sample are linked (i.e. same name) and key differences highlighted (e.g.
 368 resequencing with a new technology). As noted above, journals (or monitoring efforts like
 369 [Retraction Watch](#)) can also inform data repositories if papers have been retracted to
 370 advance data reuse. For papers reusing data, journal editors can ensure that datasets
 371 are cited correctly (see⁸⁵) and that generating authors receive equal accreditation for their
 372 work⁸⁶. We note that peer reviewers should not be tasked with these jobs, because this
 373 may be out of their realm of expertise and would increase their already high burden.

211 All scientific institutions could reward researchers with an established history of Open Data products or dataset citations on *Curriculum Vitae* and in grant proposals.

Scientific journals can also help improve genetic data archiving by ensuring genotype and read data accessibility on FAIR Trustworthy Digital Repositories before final article acceptance^{26,47}. Data are often made accessible *upon publication* with links

374

375 *Rectifying past mistakes by enriching archived data*

376 An important step many can take to advance Open Data is to improve metadata
377 for existing archives (e.g. GEOME datathons to enrich genetic metadata archives³²).
378 Publicly accessible metadata are often housed in non-standardized file formats, archived
379 with non-standard terms, or present only in published manuscripts and supplementary
380 files, consequently these can take a significant amount of time to convert for reuse³².
381 Retrospective georeferencing is often needed in such situations (e.g.^{32, 53,87}), but this often
382 relies on inference (e.g. coordinates derived from place names) leaving significant room
383 for error or lost resolution. We therefore encourage authors to enrich metadata in their
384 old data archives. We would also encourage public archiving of currently inaccessible
385 genetic datasets, and expansion of what was archived (e.g. archive all mtDNA haplotypes
386 rather than only unique haplotypes). Although older datasets may be regarded as being
387 of low value to some authors, when combined with other datasets they can be highly
388 informative and can even provide baselines for important biodiversity protection
389 assessments (e.g.^{11,88}).

390 Data enrichment initiatives could be run at the Department (similar to MoveBank⁸⁹),
391 University library, or country level (e.g. [GenDiB](#) and [CIEE Living Data](#)). Such retroactive
392 data archives could even be collaboratively published as a “resource” paper (similar to
393 those in Figure 1). These datasets could then support mandated CBD reporting on
394 genetic indicators¹⁶, inform local conservation (e.g.⁹⁰), and identify interesting scientific
395 opportunities (e.g. resampling populations after extreme events³⁵).

396

Box 2: Five take-home messages to improve genetic data archives

- 1) Archive genetic data in standardized file formats to facilitate reuse (i.e. sequences or barcodes in FASTA; microsatellites in STRUCTURE; SNPs or genomic genotypes in VCF; Genotype likelihoods in VCF; raw genomic data as demultiplexed FASTQ files).
- 2) Although no centralized database for genotype data exists, these data have great value and should be retroactively archived on FAIR compliant databases to facilitate data rescue.
- 3) Publicly archive key metadata with the genetic or genomic data, and use enriched sample names (including *a study identifier, species name, coordinates, and sampling year*).
- 4) To help more colleagues follow the FAIR principles, request both formalized data management support and a higher value of Open Data from research institutions, journals, and funding bodies.
- 5) For work involving Indigenous communities, carefully archive data affected by the CARE principles so data sovereignty is maintained.

397

398 We close on the note that genetic diversity is the most fundamental component of
399 biodiversity¹⁶. Despite underlying all levels of biodiversity, the biogeographic patterns in

400 intra-specific genetic diversity are largely understudied and poorly protected^{34,88}.
401 Improved archiving (summarised in Box 2) would expand genetic research scales far
402 beyond what any single study or research group could achieve due to logistic, cost, or
403 expertise issues. With data spanning such vast spatial and taxonomic scales, open
404 genetic data will be pivotal to new previously unimaginable areas of research and
405 conservation. Similar to data collected as part of long-term ecological monitoring
406 programs, publicly archived genetic data are likely to only become more valuable and
407 versatile when used in aggregation. This potential is pertinent and timely, due to the
408 recently signed CBD Framework which includes commitments by 192 countries to
409 conserve and restore genetic diversity within and among species' populations, and to
410 monitor and report on progress towards that commitment within the next decades⁹¹.
411 Better archiving practices will be central to meeting these targets.

412

413

414 **Corresponding Author**

415

416 `deborah.leigh@wsl.ch`

417

418

419 **Acknowledgements**

420

421 D.M.L. was funded by the BiodivERsA project "ACORN" granted by the Swiss National Science
422 Foundation (SNSF Project 31BD30_193900). I.P-V. was supported by the U.S. Geological Survey
423 John Wesley Powell Center for Analysis and Synthesis. Thanks to Torsten Günther and Bastiaan
424 Star for their comments on ancient DNA archiving practices and considerations. Thanks also to
425 Jennifer Gibson for her helpful discussions about FAIR databases. Thanks to Felix Gugerli and
426 Corine Buser-Schoebel for their helpful feedback on the manuscript. This work was conducted as
427 a part of the Standardizing, Aggregating, Analyzing and Disseminating Global Wildlife Genetic
428 and Genomic Data for Improved Management and Advancement of Community Best Practices
429 Working Group supported by the John Wesley Powell Center for Analysis and Synthesis, funded
430 by the U.S. Geological Survey. Any use of trade, firm, or product names is for descriptive purposes
431 only and does not imply endorsement by the U.S. Government.

432

433 **Data accessibility statement:**

434

435 The data underpinning Figure 1 are available for reviewers at this private link
436 [https://docs.google.com/spreadsheets/d/14RyPOCwmZL8LdqH4t-
437 xPlwy2u9QDoWiK/edit?usp=sharing&oid=110189835630079798646&rtpof=true&sd=true](https://docs.google.com/spreadsheets/d/14RyPOCwmZL8LdqH4t-xPlwy2u9QDoWiK/edit?usp=sharing&oid=110189835630079798646&rtpof=true&sd=true) and
438 will be uploaded to a Dryad when a journal link is supplied.

439

440 **Conflicts of interest:**

441

442 The authors declare no conflicts of interest.

443

444 **Author contribution statement**

445

446 All authors contributed to the inception and writing of this work. DML supervised this work and
447 conducted the editing, with support from IPV, AGV, and MEH. IPV conducted the database
448 included analysis.

449

450 **References**

- 451 1. Vines, T. H. *et al.* The Availability of Research Data Declines Rapidly with Article Age. *Curr.*
452 *Biol.* **24**, 94–97 (2014).
- 453 2. Roche, D. G., Kruuk, L. E. B., Lanfear, R. & Binning, S. A. Public Data Archiving in Ecology
454 and Evolution: How Well Are We Doing? *PLOS Biol.* **13**, e1002295 (2015).
- 455 3. Tedersoo, L. *et al.* Data sharing practices and data availability upon request differ across
456 scientific disciplines. *Sci. Data* **8**, 192 (2021).
- 457 4. Piwowar, H. A., Vision, T. & Whitlock, M. C. Data archiving is a good investment. *Nature* **473**,
458 (2011).
- 459 5. Cochrane, G., Cook, C. E. & Birney, E. The future of DNA sequence archiving. *GigaScience*
460 **1**, 2 (2012).
- 461 6. Strasser, B. J. The Experimenter’s Museum: GenBank, Natural History, and the Moral
462 Economies of Biomedicine. *Isis* **102**, 60–96 (2011).
- 463 7. International Human Genome Mapping Consortium *et al.* A physical map of the human
464 genome. *Nature* **409**, 934–941 (2001).
- 465 8. Ratnasingham, S., Hebert P,D. bold: The Barcode of Life Data System
466 (<http://www.barcodinglife.org>). *Mol Ecol Notes.* **7**, 355-364 (2007)
- 467 9. Blanchet, S., Prunier, J. G. & De Kort, H. Time to Go Bigger: Emerging Patterns in
468 Macrogenetics: Trends in Genetics. *Trends Genet.* **33**, 579–580 (2017).
- 469 10. Leigh, D. M. *et al.* Opportunities and challenges of macrogenetic studies. *Nat. Rev. Genet.*
470 **22**, 791–807 (2021).

- 471 11.Schmidt, C., Hoban, S. & Jetz, W. Conservation macrogenetics: harnessing genetic data to
472 meet conservation commitments. *Trends Genet.* **39**, 816–829 (2023).
- 473 12.Ruppert, K. M., Kline, R. J. & Rahman, M. S. Past, present, and future perspectives of
474 environmental DNA (eDNA) metabarcoding: A systematic review in methods, monitoring, and
475 applications of global eDNA. *Glob. Ecol. Conserv.* **17**, e00547 (2019).
- 476 13.Günther, T. & Coop, G. Robust Identification of Local Adaptation from Allele Frequencies.
477 *Genetics* **195**, 205–220 (2013).
- 478 14.CBD. Decision Adopted By The Conference Of The Parties To The Convention On Biological
479 Diversity. (2022) doi:<https://www.cbd.int/doc/decisions/cop-15/cop-15-dec-05-en.pdf>.
- 480 15.Hoban, S. *et al.* Genetic diversity targets and indicators in the CBD post-2020 Global
481 Biodiversity Framework must be improved. *Biol. Conserv.* **248**, 108654 (2020).
- 482 16.Hoban, S. *et al.* Monitoring status and trends in genetic diversity for the Convention on
483 Biological Diversity: An ongoing assessment of genetic indicators in nine countries. *Conserv.*
484 *Lett.* **16**, e12953 (2023).
- 485 17.Lawrence, E. R. *et al.* Geo-referenced population-specific microsatellite data across
486 American continents, the MacroPopGen Database. *Sci. Data* **6**, 14 (2019).
- 487 18.Schmidt, C., Domaratzki, M., Kinnunen, R. P., Bowman, J. & Garroway, C. J. Continent-wide
488 effects of urbanization on bird and mammal genetic diversity. *Proc R Soc B* **287**, (2020).
- 489 19.Schmidt, C. & Garroway, C. J. Systemic racism alters wildlife genetic diversity. *Proc. Natl.*
490 *Acad. Sci.* **119**, e2102860119 (2022).
- 491 20.Schmidt, C. & Garroway, C. J. The population genetics of urban and rural amphibians in
492 North America. *Mol. Ecol.* **30**, 3918–3929 (2021).
- 493 21.Rieseberg, L., Vines, T. & Kane, N. Editorial and retrospective 2010. *Mol. Ecol.* **19**, 1–22
494 (2010).
- 495 22.Moore, A. J., Mcpeek, M. A., Rausher, M. D., Rieseberg, L. & Whitlock, M. C. The need for
496 archiving data in evolutionary biology. *J. Evol. Biol.* **23**, 659–660 (2010).

- 497 23. Whitlock, M. C. Data archiving in ecology and evolution: best practices. *Trends Ecol. Evol.*
498 **26**, 61–65 (2011).
- 499 24. Fairbairn, D. J. The advent of mandatory data archiving. *Evolution* **65**, 1–2 (2011).
- 500 25. Berberi, I. & Roche, D. G. No evidence that mandatory open data policies increase error
501 correction. *Nat. Ecol. Evol.* **6**, 1630–1633 (2022).
- 502 26. Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and
503 stewardship. *Sci. Data* **3**, 160018 (2016).
- 504 27. Gomes, D. G. E. *et al.* Why don't we share data and code? Perceived barriers and benefits
505 to public archiving practices. *Proc R Soc B* **289**, 2022111 (2022).
- 506 28. Huang, X. *et al.* Willing or unwilling to share primary biodiversity data: results and
507 implications of an international survey: Biodiversity data sharing and archiving. *Conserv. Lett.*
508 **5**, 399–406 (2012).
- 509 29. Hostler, T. J. The Invisible Workload of Open Research. *J. Trial Error* (2023)
510 doi:10.36850/mr5.
- 511 30. Kozlov, M. How A Spider-Biology Scandal Upended Researchers' Lives. *Nature* **608**, (2022).
- 512 31. European Commission. H2020 Programme: AGA – Annotated Model Grant Agreement.
513 (2019).
- 514 32. Crandall, E. D. *et al.* Importance of timely metadata curation to the global surveillance of
515 genetic diversity. *Conserv. Biol.* **37**, e14061 (2023).
- 516 33. Ceballos, G. *et al.* Accelerated modern human-induced species losses: Entering the sixth
517 mass extinction. *Sci. Adv.* **1**, e1400253 (2015).
- 518 34. Leigh, D. M., Hendry, A. P., Vázquez-Domínguez, E. & Friesen, V. L. Estimated six per cent
519 loss of genetic variation in wild populations since the industrial revolution. *Evol. Appl.* **12**,
520 1505–1512 (2019).
- 521 35. Jensen, E. L. & Leigh, D. M. Using temporal genomics to understand contemporary climate
522 change responses in wildlife. *Ecol. Evol.* **12**, (2022).

- 523 36.Lischer, H. E. L. & Excoffier, L. PGDSpider: an automated data conversion tool for
524 connecting population genetics and genomics programs. *Bioinformatics* **28**, 298–299 (2012).
- 525 37.Adamack, A. T. & Gruber, B. P OP G EN R EPORT : simplifying basic population genetic
526 analyses in R. *Methods Ecol. Evol.* **5**, 384–387 (2014).
- 527 38.Manoukis, N. C. formatomatic: a program for converting diploid allelic data between common
528 formats for population genetic analysis. *Mol. Ecol. Notes* **7**, 592–593 (2007).
- 529 39.Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158
530 (2011).
- 531 40.Purcell, S. *et al.* PLINK: A Tool Set for Whole-Genome Association and Population-Based
532 Linkage Analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 533 41.Jombart, T. *adegenet* : a R package for the multivariate analysis of genetic markers.
534 *Bioinformatics* **24**, 1403–1405 (2008).
- 535 42.Gratton, P. *et al.* A world of sequences: can we use georeferenced nucleotide databases for
536 a robust automated phylogeography? *J. Biogeogr.* **44**, 475–486 (2017).
- 537 43.Pritchard, J. K., Stephens, M. & Donnelly, P. Inference of Population Structure Using
538 Multilocus Genotype Data. *Genetics* **155**, 945–959 (2000).
- 539 44.Cochrane, G., Cook, C. E. & Birney, E. The future of DNA sequence archiving. *GigaScience*
540 **1**, 2 (2012).
- 541 45.Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–
542 2079 (2009).
- 543 46.Mallick, S., *et al.* The Allen Ancient DNA Resource (AADR) a curated compendium of ancient
544 human genomes. *Sci Data* **11**, 182 (2024).
- 545 47.Jenkins, G. B. *et al.* Reproducibility in ecology and evolution: Minimum standards for data
546 and code. *Ecol. Evol.* **13**, e9961 (2023).
- 547 48.Grealey, J. *et al.* The Carbon Footprint of Bioinformatics. *Mol. Biol. Evol.* **39**, msac034
548 (2022).

549 49. Wieczorek, J. *et al.* Darwin Core: An Evolving Community-Developed Biodiversity Data
550 Standard. *PLOS ONE* **7**, e29715 (2012).

551 50. Field, D. *et al.* The Genomic Standards Consortium. *PLoS Biol.* **9**, e1001088 (2011).

552 51. Meyer, R. *et al.* Aligning Standards Communities for Omics Biodiversity Data: Sustainable
553 Darwin Core-MIxS Interoperability. *Biodivers. Data J.* **11**, e112420 (2023).

554 52. Buttigieg, P. *et al.* The environment ontology: contextualising biological and biomedical
555 entities. *J. Biomed. Semant.* **4**, 43 (2013).

556 53. Böhne, A. *et al.* Contextualising samples: Supporting reference genomes of European
557 biodiversity through sample and associated metadata collection. *bioRxiv*.
558 <https://doi.org/10.1101/2023.06.28.546652> (2024).

559 54. Stroe, O. ENA to introduce mandatory spatiotemporal annotations.
560 [https://www.ebi.ac.uk/about/news/updates-from-data-resources/ena-spatiotemporal-](https://www.ebi.ac.uk/about/news/updates-from-data-resources/ena-spatiotemporal-metadata/)
561 [metadata/](https://www.ebi.ac.uk/about/news/updates-from-data-resources/ena-spatiotemporal-metadata/) (2023).

562 55. Frank, R. D., Kriesberg, A., Yakel, E. & Faniel, I. M. Looting hoards of gold and poaching
563 spotted owls: Data confidentiality among archaeologists & zoologists. *Proc. Assoc. Inf. Sci.*
564 *Technol.* **52**, 1–10 (2015).

565 56. Chapman, A. D. Current Best Practices for Generalizing Sensitive Species Occurrence Data.

566 57. Clarke, K. C. A multiscale masking method for point geographic data. *Int. J. Geogr. Inf. Sci.*
567 **30**, 300–315 (2016).

568 58. Scholz, A. H. *et al.* Multilateral benefit-sharing from digital sequence information will support
569 both science and biodiversity conservation. *Nat. Commun.* **13**, 1086 (2022).

570 59. Marden, E. *et al.* Sharing and reporting benefits from biodiversity research. *Mol. Ecol.* **30**,
571 1103–1107 (2021).

572 60. Bhaumik, V. Global inequities in local science. *Nat. Ecol. Evol.* **7**, 793–793 (2023).

573 61. Miller, J., White, T. B. & Christie, A. P. Parachute conservation: Investigating trends in
574 international research. *Conserv. Lett.* **16**, e12947 (2023).

- 575 62.de Vos, A. & Schwartz, M. W. Confronting parachute science in conservation. *Conserv. Sci.*
576 *Pract.* **4**, e12681 (2022).
- 577 63.Carroll, S. R. The CARE Principles for Indigenous Data Governance - Data Science Journal.
578 *Data Sci. J.* **19**, 1–12.
- 579 64.Carroll, S. R., Herczog, E., Hudson, M., Russell, K. & Stall, S. Operationalizing the CARE
580 and FAIR Principles for Indigenous data futures. *Sci. Data* **8**, 108 (2021).
- 581 65.Kukutai, T. Indigenous data sovereignty—A new take on an old theme. *Science* **382**,
582 eadl4664 (2023).
- 583 66.Te Aika, B. *et al.* Aotearoa genomic data repository: An āhuru mōwai for taonga species
584 sequencing data. *Mol. Ecol. Resour.* 1-14 (2023)
- 585 67.Hudson, M. *et al.* Indigenous Peoples' Rights in Data: a contribution toward Indigenous
586 Research Sovereignty. *Front. Res. Metr. Anal.* **8**, (2023).
- 587 68.Mc Cartney, A. M. *et al.* Indigenous peoples and local communities as partners in the
588 sequencing of global eukaryotic biodiversity. *Npj Biodivers.* **2**, 1–12 (2023).
- 589 69. Shaikh A. (2014). Ecology week 4: field sample with animals. figshare. Dataset.
590 <https://doi.org/10.6084/m9.figshare.1194651.v1>
- 591 70. Gonzalez L. (2010). Sexual crime in Colombia 2010-2022. figshare. Dataset.
592 <https://doi.org/10.6084/m9.figshare.21937154.v1>
- 593 71.Roche, D. G., Jennions, M. D. & Binning, S. A. Fees could damage public data archives.
594 *Nature* **502**, 171–171 (2013).
- 595 72.Barrett, T. *et al.* BioProject and BioSample databases at NCBI: facilitating capture and
596 organization of metadata. *Nucleic Acids Res.* **40**, D57–D63 (2012).
- 597 73.Deck, J. *et al.* The Genomic Observatories Metadatabase (GeOMe): A new repository for
598 field and sampling event metadata associated with genetic samples. *PLOS Biol.* **15**,
599 e2002925 (2017).

600 74. Shaw, F. *et al.* COPO: a metadata platform for brokering FAIR data in the life sciences.
601 *F1000Research* **9**, 495 (2020).

602 75. Web of Science. Associated Data. (2018).
603 https://images.webofknowledge.com/images/help/WOK/hp_associated_data.html

604 76. DDBJ Including Sample Location and Collection Date and Time for BioSample submissions
605 Including Sample Location. (2023) <https://www.ddbj.nig.ac.jp/news/en/2023-05-02-e.html>

606 77. Costa-Pereira, R. & Pruitt, J. Retraction: Behaviour, morphology and microhabitat use: what
607 drives individual niche variation? *Biol. Lett.* **16**, 20200588 (2020).

608 78. van den Burg, M. P. & Vieites, D. R. Bird genetic databases need improved curation and
609 error reporting to NCBI. *Ibis* **165**, 472–481 (2023).

610 79. National Institute for Health. Final NIH Policy for Data Management and Sharing: NOT-OD-
611 21-013 <https://grants.nih.gov/grants/guide/notice-files/NOT-OD-21-013.htm> (2020)

612 80. European Commission. Guidelines on FAIR data management in Horizon 2020. (2016).

613 81. UKRI. Data Management Plan: Guidance for Peer Reviewers. (2013)
614 [https://www.ukri.org/wp-content/uploads/2021/07/ESRC-200721-DataManagementPlan-](https://www.ukri.org/wp-content/uploads/2021/07/ESRC-200721-DataManagementPlan-GuidanceforPeerReviewers.pdf)
615 [GuidanceforPeerReviewers.pdf](https://www.ukri.org/wp-content/uploads/2021/07/ESRC-200721-DataManagementPlan-GuidanceforPeerReviewers.pdf) .

616 82. Peng, G. *et al.* Scientific Stewardship in the Open Data and Big Data Era Roles and
617 Responsibilities of Stewards and Other Major Product Stakeholders. *D-Lib Mag.* **22**, (2016).

618 83. Toelch, U. & Ostwald, D. Digital open science—Teaching digital tools for reproducible and
619 transparent research. *PLOS Biol.* **16**, e2006022 (2018).

620 84. Thrall, P. H. *et al.* From raw data to publication: Introducing data editing at Ecology Letters.
621 *Ecol. Lett.* **26**, 829–830 (2023).

622 85. Cousijn, H. *et al.* A data citation roadmap for scientific publishers. *Sci. Data* **5**, 180259
623 (2018).

624 86. Nature. Time to recognize authorship of open data. *Nature* **604**, 8–8 (2022).

625 87. Miraldo, A. *et al.* An Anthropocene map of genetic diversity. *Science* **353**, 1532–1535 (2016).

- 626 88.Figuerola-Ferrando, L. *et al.* Global patterns and drivers of genetic diversity among marine
627 habitat-forming species. *Glob. Ecol. Biogeogr.* **32**, 1218–1229 (2023).
- 628 89.Kays, R. *et al.* The Movebank system for studying global animal movement and
629 demography. *Methods Ecol. Evol.* **13**, 419–431 (2022).
- 630 90.Beninde, J. *et al.* CaliPopGen: A genetic and life history database for the fauna and flora of
631 California. *Sci. Data* **9**, 380 (2022).
- 632 91.Hoban, S. *et al.* Genetic diversity goals and targets have improved, but remain insufficient for
633 clear implementation of the post-2020 global biodiversity framework. *Conserv. Genet.* **24**,
634 181–191 (2023).

635
636
637
638
639
640
641
642
643
644
645

***** END*****