1

**Title:** Don't make genetic data disposable: Best practices for genetic and
genomic data archiving

**Authors:** Deborah M. Leigh[1]*, Amy Vandergast[2,] Margaret E. Hunter[3], Eric Crandall[4], W. Chris
Funk[5], Colin J. Garroway[6], Sean Hoban[7], Sara J. Oyler-McCance[8], Christian Rellstab[1], Gernot
Segelbacher[9], Chloe Schmidt[10], Ella Vázquez-Domínguez[11], Ivan Paz-Vinas[12,13]

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


**Deborah M. Leigh**
    1) Swiss Federal Research Institute WSL, Birmensdorf, Switzerland
       ORCID: 0000-0003-3902-2568
       deborah.leigh@wsl.ch


**Amy Vandergast**
    2) U.S. Geological Survey, Western Ecological Research Center, 4165 Spruance Road,
       Suite 200, San Diego CA, 92101, USA
       ORCID: 0000-0002-7835-6571
       avandergast@usgs.gov


**Margaret E. Hunter**
    3) U.S. Geological Survey, Wetland & Aquatic Research, Center, Sirenia Project 7920 NW
       71st Street, Gainesville, Florida 32653, USA
       ORCID: 0000-0002-4760-9302
       mhunter@usgs.gov

**Eric D. Crandall**
    4) Department of Biology, Pennsylvania State University, 208 Mueller Laboratory,
       University Park, PA 16802, USA
       ORCID: 0000-0001-8580-3651
       ecrandall@psu.edu

**W. Chris Funk**
    5) Department of Biology, Graduate Degree Program in Ecology, Colorado State
       University, 1878 Campus Delivery, Fort Collins, CO 80523-1878, USA
       ORCID: 0000-0002-6466-3618
       Chris.Funk@colostate.edu

**Colin J. Garroway**
    6) Department of Biological Sciences, University of Manitoba, Winnipeg, Manitoba, Canada
       ORCID: 0000-0002-0955-0688
       colin.garroway@umanitoba.ca

**Sean Hoban**

　　7) Center for Tree Science, The Morton Arboretum, Lisle, IL, 60532, USA *and*

　　Committee on Evolutionary Biology, University of Chicago, Chicago, IL, 60637, USA

　　ORCID: 0000-0002-0348-8449

　　shoban@mortonarb.org

**Sara J. Oyler-McCance**

　　8) U.S. Geological Survey, Fort Collins Science Center, 2150 Centre Avenue, Building

　　C, Fort Collins, CO, 80526, USA

　　ORCID: 0000-0003-1599-8769

　　soyler@usgs.gov

**Christian Rellstab**

　　1) Swiss Federal Research Institute WSL, Birmensdorf, Switzerland

　　ORCID: 0000-0002-0221-5975

　　christian.rellstab@wsl.ch

**Gernot Segelbacher**

　　9) Wildlife Ecology and Management, University Freiburg, Tennenbacher Str. 4, 79106

　　Freiburg, Germany

　　ORCID: 0000-0002-8024-7008

　　gernot.segelbacher@wildlife.uni-freiburg.de

**Chloe Schmidt**

　　10) German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig,

　　Puschstrasse 4, Leipzig, 04103, Germany

　　ORCID 0000-0003-2572-4200

　　chloe.schmidt@idiv.de

**Ella Vázquez-Domínguez**

　　11) Departamento de Ecología de la Biodiversidad, Instituto de Ecología, Universidad

　　Nacional Autónoma de México, Coyoacán, Ciudad de México, 04510, México;

　　ORCID: 0000-0001-6131-2014

　　evazquez@ecologia.unam.mx

**Ivan Paz-Vinas**

　　12) Department of Biology, Colorado State University, 1878 Campus Delivery, Fort

　　Collins, CO, 80523-1878, USA

　　ORCID: 0000-0002-0043-9289

　　13) Université Claude Bernard Lyon 1, CNRS, ENTPE, UMR5023 LEHNA, F-69622,

　　Villeurbanne, France

　　ivan.paz-vinas@univ-lyon1.fr

90 *Don't make genetic data disposable:*
91 *Best practices for genetic and genomic data archiving*
92

# Abstract

94     In ecology and evolution, genetic and genomic data are commonly collected for a
95 vast array of scientific and applied purposes. Despite mandates for public archiving, such
96 data are typically used only once by the data-generating authors. The repurposing of
97 genetic and genomic datasets remains uncommon because it is often difficult, if not
98 impossible, due to non-standard archiving practices and lack of contextual metadata. But
99 as the new research field of macrogenetics is demonstrating, if genetic data and their
100 metadata were more accessible, they could be reused for many additional purposes, far
101 beyond their initial intended impact. In this review, we outline the main challenges with
102 existing genetic and genomic data archives, factors underlying the challenges, and
103 current best practices for archiving genetic and genomic data. Recognising that this is a
104 longstanding issue due to an absence of formal data management training within the
105 research field of ecology and evolution, we highlight key steps that universities, funding
106 bodies, and scientific publishers could take to ensure timely change towards good data
107 archiving.
108

# Introduction

110     Synthesis of Open Data (publicly archived data free to reuse) is a powerful tool
111 that is increasingly being used to test pressing big-picture questions at large scales in
112 ecology and evolution. However, it still remains common for valuable datasets to be
113 forgotten and mislaid after a single use (Vines et al. 2014; Roche et al. 2015; Tedersoo
114 et al. 2021). This is a missed opportunity and hinders scientific progress. Producing and
115 collecting scientific data is often expensive and time-consuming. Furthermore, most data
116 have numerous potential applications beyond their original use (Piwowar et al. 2011).
117     Public archiving of genetic and genomic sequence data (hereafter 'genetic data')
118 became standard practice in the 1980s (Cochrane et al. 2012), but notably, public
119 archiving of associated metadata (metadata are data that describe other data, including
120 species name, sampling coordinates, sampling year, etc.), still remains discretionary.
121 Nevertheless, genetic data repositories were some of the earliest Open Data projects and
122 databases (e.g. Genbank; Strasser et al. 2011) and continue to arise to meet the
123 increasing needs of genomic data archiving (e.g. International Human Genome Mapping
124 Consortium 2001; BOLD, Ratnasingham and Hebert 2007).
125     Repurposed 'open' population genetic data has only just begun to accumulate but
126 has facilitated reconstruction of, for example, endangered species' demographic histories
127 (e.g. orangutans; Nater et al. 2015), and inference of global invasion pathways (e.g.
128 *Trachemys scripta elegans*; Espindola et al. 2022). Multi-species genotype data are also
129 being synthesized across large spatial and temporal ranges for macrogenetic studies

130  (Blanchet et al. 2017; Leigh et al. 2021), rapidly advancing molecular ecology and
131  evolution by characterizing global biodiversity patterns, genetic diversity trends, and
132  informing biodiversity conservation (see Leigh et al. 2021 and references therein; Schmidt
133  et al. 2023). Sequences within the National Center for Biotechnology Information (NCBI)
134  are frequently reused as a biodiversity reference database, facilitating new species
135  discovery and the emergence of environmental DNA methods (Ruppert et al. 2019).
136  Accessible raw genomic read datasets have been important for teaching bioinformatics
137  and developing genomic analyses (e.g. Günther and Coop, 2013). Yet the future
138  repurposing potential of genetic data extends further, as an abundance of unattempted
139  and unknown uses remain. Vitally, repurposing of public genetic data is one way for
140  countries to report genetic indicators required by the Convention of Biological Diversity
141  (CBD) post-2020 Kunming-Montreal global biodiversity framework (e.g. headline indicator
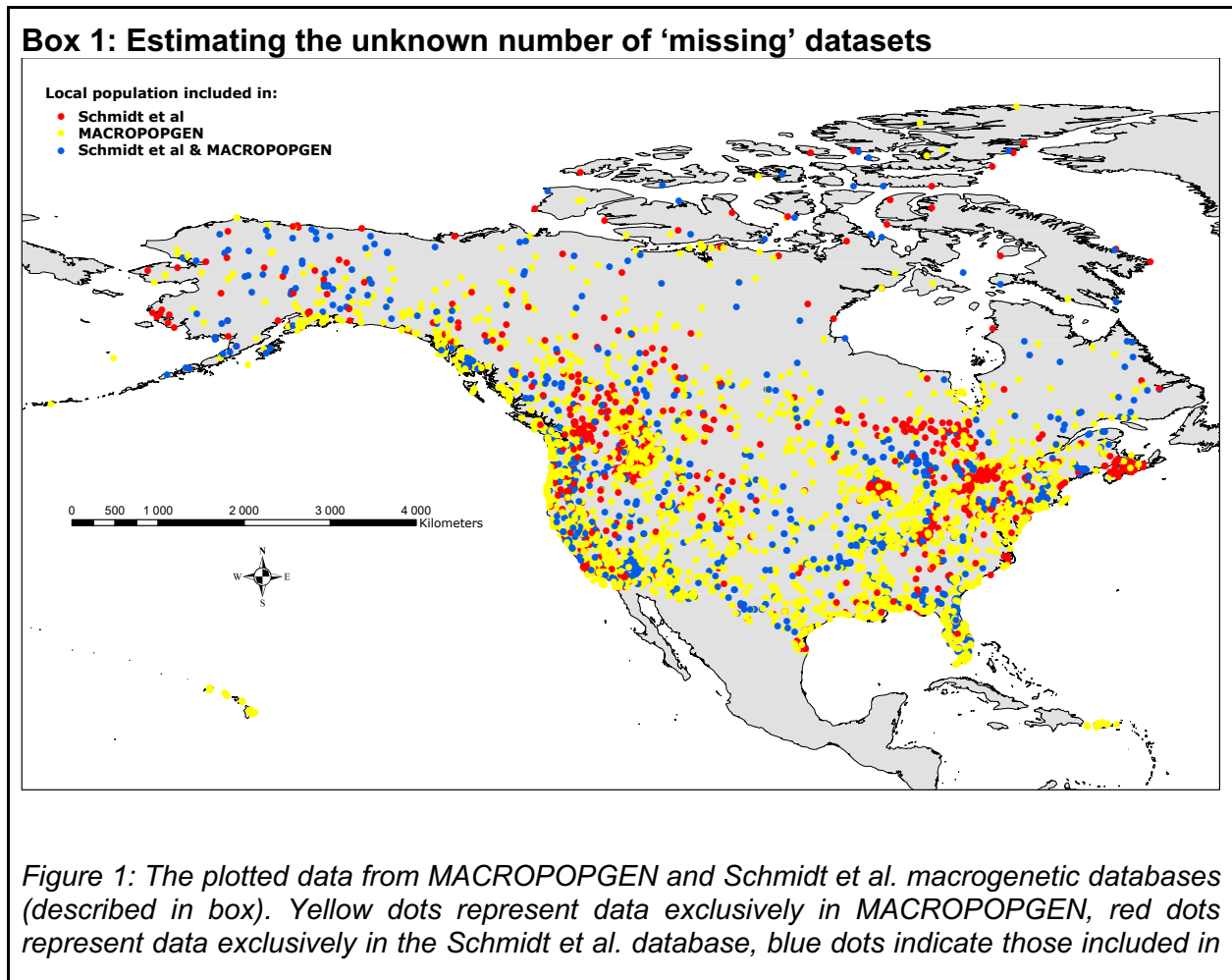142  A.4; CBD 2022; Hoban et al. 2023; Hoban et al. 2020).
143      Despite the abundance of genetic data in open repositories, long-standing
144  willingness of journals to mandate Open Data (e.g. JDAP Dryad 2011; Rieseberg et al.
145  2010; Moore et al. 2010; Whitlock 2011; Fairbairn 2011), and increasing popularity of the
146  FAIR principles (Findable, Accessible, Interoperable, and Reusable; Wilkinson et al.
147  2016), there are still numerous issues with genetic data archiving that inhibit
148  comprehensive repurposing. Many genetic datasets remain in private hands, often stored
149  on private storage devices, shared only on request (see Box 1). However, such devices
150  will quickly depreciate, leading to data loss. Furthermore, cross-disciplinary studies have
151  shown that many authors do not share data upon request, despite committing to do so in
152  data availability statements of their articles (Gabelica et al. 2022; Crandall et al. 2023).
153  Notably, ecologists (including molecular ecologists) were the most likely researchers to
154  ignore data request emails (Tedersoo et al. 2021). Data findability is also limited by the
155  repositories available, because genetic databases are built for nucleotide sequence data
156  and do not accept non-nucleotide data (e.g. processed Single Nucleotide Polymorphism
157  "SNP" genotypes and/or genotypic datasets based on microsatellites). These abundant
158  and valuable genotype data are subsequently archived, if at all, across an array of public
159  and private data repositories (e.g. Dryad, FigShare, ScienceBase, Zenodo, and personal
160  or institutional servers), making them hard to find (see Box 1).
161      Alongside issues of data accessibility and findability, genetic data that are publicly
162  archived are often in an unsuitable format. This is predominantly because they lack key
163  metadata. It is easy to underestimate the severity of poor metadata archiving: only ~6.5%
164  of published nucleotide data for land-living vertebrates are georeferenced in GenBank
165  (Gratton et al. 2017). Similarly, only 13% of biodiversity-relevant Sequence Read Archive
166  (SRA) BioProjects have spatiotemporal metadata (Toczydlowski et al. 2021), and less
167  than 33% of metagenomic data are archived with vital contextual environmental data
168  (Schriml et al. 2020). Beyond the absence of critical metadata, non-sequence genotypic
169  data are also archived in an array of formats often tailored for specific software packages,

170 which themselves change over time. Moreover, studies frequently differ in exactly what
171 they archive. This ranges from aligned sequences to raw data in genomic datasets, to
172 newly identified haplotypes only vs the entire set of sequences obtained in amplicon or
173 barcoding studies (Paz-Vinas et al. 2021).
174       Consequently, despite the long history of publicly archived genetic data in ecology
175 and evolution, most remain difficult to repurpose. Poor archiving represents a significant
176 waste of public funds and a loss of time, resources, and opportunities for scientists and
177 practitioners. This also represents an unnecessary ethical footprint because some genetic
178 studies require animal handling. It severely limits the development of promising Open
179 Data reliant research avenues and biodiversity monitoring. Collectively, poor archiving
180 lowers the impact of each dataset. There is a pressing need in ecology and evolution to
181 improve the archiving standards of genetic data and their metadata by establishing best
182 practices. In this review, we seek to address these challenges by offering guidance on
183 archiving different types of genetic data and their associated metadata. We also discuss
184 additional steps or infrastructure needed to improve the status quo. Ultimately, our goal
185 is to prevent data loss and facilitate data reuse.
186

---

**Box 1: Estimating the unknown number of 'missing' datasets**



**Local population included in:**
- Schmidt et al
- MACROPOPGEN
- Schmidt et al & MACROPOPGEN

*Figure 1: The plotted data from MACROPOPGEN and Schmidt et al. macrogenetic databases (described in box). Yellow dots represent data exclusively in MACROPOPGEN, red dots represent data exclusively in the Schmidt et al. database, blue dots indicate those included in*

*both databases.*

To estimate the number of genetic datasets that have not been publicly archived, we compared the datasets synthesized in two macrogenetic databases with similar taxonomic and spatial scopes that differed in the form of included data. The first macrogenetic database was MACROPOPGEN (Lawrence et al. 2019), which is a compilation of georeferenced vertebrate microsatellite-based genetic summary statistics and related metadata (e.g. taxonomic group) across the Americas for 897 species and 9,090 populations. This database was built by extracting data (e.g. genetic summary statistics) from published articles and reports, irrespective of whether raw genotypic datasets from which the compiled data were derived are publicly available. The second macrogenetic database (hereafter, Schmidt et al. database) was compiled by Schmidt and Garroway (2021; 2022) and Schmidt et al. (2020) and comprises repurposed microsatellite genotype datasets that were archived in open repositories (mostly those archived on Dryad) for terrestrial vertebrates across USA and Canada.

If all raw genetic datasets used in the studies identified in MACROPOPGEN were publicly available, we would expect the proportion of datasets overlapping with Schmidt et al. to be high (e.g. 80-90%), and if all data were clearly linked and/or archived with metadata we would expect the overlap to increase further (e.g. 90-100%).

We extracted datasets from birds, amphibians, mammals, and reptiles located in USA and Canada from both databases, and combined them to create a pooled database containing data for 5,395 populations from 412 species, with 68.48% and 31.52% of data originating from MACROPOPGEN and the Schmidt et al. database, respectively. Data were at the population level (Figure 1). We then crosschecked different metadata fields (e.g. species name, DOI identifier of the dataset and/or of the original article, author names) among data entries from each macrogenetic database to identify populations and datasets that were included in both macrogenetic databases.

Only 21.38% of the data entries in the combined database were found in both macrogenetic databases (Figure 1, blue dots), while 59.5% were included exclusively in MACROPOPGEN (Figure 1, yellow dots). While this does not comprehensively assess the number of publications missing from public data archives, this strongly indicates a large proportion of the studies behind this 59.5% are unlikely to have public genetic or genotype datasets, given the similarity in spatial and taxonomic scopes used by the authors of both databases.

187

188 *Why are we not archiving comprehensively?*

189       Several general issues drive variation in data archiving practices that ultimately
190 hinder data reuse (see Box 2; Tedersoo et al. 2021; Gomes et al. 2022; Huang et al. 2012;
191 Roche et al. 2014; Hostler et al. 2023). Most commonly, poor data archiving is driven by
192 researchers not having sufficient support or time to archive comprehensively (Hostler et
193 al. 2023). Furthermore, authors often may not realize their archives are incomplete and
194 challenging to reuse. A well-archived dataset should contain all of the relevant information
195 needed to reproduce or repurpose a study without the need to consult multiple sources.
196       Limited data archiving can sometimes be a requirement in ecology and evolution.
197 There are legal and ethical considerations in data archiving and reuse that can directly

198 limit what data can be made public (see Box 3). These limitations should be accepted
199 because it is essential that Open Data does not infringe on privacy, benefit sharing, or
200 species protection efforts (Frank et al. 2016). Unfortunately, however, limited data
201 archiving can sometimes be intentional to prevent repurposing; in contrast to sensitive
202 data intentionally poor archiving needs to be rectified as it inhibits FAIR compliance (see
203 Box 2 and 3).

204

205

**Box 2: Limiting competition through intentionally poor data archiving**

FAIR data archiving can fail due to a range of reasons (e.g. Tedersoo et al. 2021; Gomes et al. 2022; Huang et al. 2012; Roche et al. 2014; Hostler et al. 2023). These include a fear of scooping - when data are repurposed before the data-generating authors benefit through primary publication(s), collegial competition (Huang et al. 2012), or because researchers consider their study topic exclusive. Reluctance to archive data may also be related to a fear that reanalysis of data might reveal errors or lead to contradictory conclusions (Wicherts et al. 2011). Intentionally poor archiving can sadly inhibit reanalysis and synthesis that can lead to exciting new conclusions (e.g. Schumacher et al. 2022) and could harm biodiversity monitoring efforts.

We acknowledge that there are long standing valid concerns surrounding data archiving from long-term multi-grant studies (e.g. wild pedigreed populations) because reuse by external researchers at any point could have a disproportionately negative impact on data-generating scientists and the project (Mills et al. 2015; Whitlock et al. 2015). Similarly, data-repurposing from early career researchers (particularly matriculated students) can be harmful because they commonly have long delays prior to publication. In both such situations, to ensure FAIR compliance, data-generating researchers have an obligation to archive their data, but can use embargos to protect their planned analyses/publications (Mills et al. 2015; Whitlock et al. 2015). After embargos expire, ethical and sensitive data repurposing is vital in such situations to help maintain FAIR compliance; this includes discussion of data repurposing with generating authors and land owners/indigenous communities.

206
207

**Box 3: How to archive data from sensitive species**

Metadata from threatened or endangered species, as well as species that are commercially valuable or desirable, may need to be withheld or obscured to protect them (Frank et al. 2016). Withholding metadata facilitates species protection by mitigating the risk of poaching and/or habitat degradation caused by increased disturbance arising from species viewing or photography (Lindenmayer and Scheele 2017). Furthermore, for species on private land, this can protect the collaborations necessary for conservation (Lindenmayer and Scheele 2017). However, some argue that metadata must be published albeit with considerations (i.e. masking) or accessible upon request (Lowe et al. 2017).

The data of most concern for such species is location data – coordinates or specific habitat descriptions that would allow public access to these species. Best practices for generalizing sensitive species occurrence or geographic metadata have been developed (see Chapman 2020; Clarke 2016). For example, the Global Biodiversity Information Facility (GBIF) has mechanisms to incorporate location generalization and ways to document that information exists but is withheld for privacy (e.g. metadata field "informationWithheld"). Importantly the release of other data from sensitive species, such as genetic and genomic data, could facilitate conservation but their potential commercial value (e.g. for pharmaceutical or agricultural companies) should not be ignored. While we acknowledge some analyses cannot be performed without fairly accurate location data (e.g. genotype-environment association, macrogenetics), access can normally be arranged when needed. Release of metadata (e.g. number of individuals, age/reproductive status, sex, etc.), should be dependent on the potential risk to the species (i.e. providing age and size of a valuable tree, game species, or a medicinal herb may increase the likelihood harvest, but will be unimportant in other cases; Lowe et al. 2017). Chapman (2022) provides a decision tree to assist with such choices.

Examples of sensitive data archiving are the Greater Sage-grouse (*Centrocercus urophasianus*) and Gunnison Sage-grouse (*C. minimus*), both of which are species of significant conservation concern in North America. Males from both species gather in mating grounds (leks) to attract females and such places are often used for genetic sampling as well as observation by hobbyists. Yet human presence can disturb mating activities. Further, some leks occur on private land, requiring collaboration with landowners. While the genetic data from mating ground samples is publicly available, location information is either generalized (Zimmerman et al. 2019), or an averaged location for a group of mating grounds given (Row et al. 2018, Cross et al. 2018), or only on request (Oyler-McCance et al. 2022). This data masking step prevents disturbance increase and supports conservation.

In another example, the North American butternut (*Juglans cinerea L.*) is a species of conservation concern (IUCN Endangered), but is valuable for timber and traditional medicine (Pike et al. 2021). Sharing location data could lead to harmful timber harvest, thus population coordinates have sometimes been published with a random geographic offset (e.g. such as 10 km; Hoban et al. 2010). This simultaneously protects the species and allows for most genetic and geographic data repurposing, without the need for data access requests. Researchers also took care to remove location names (e.g. the name of a creek, landowner, or nearest town) to prevent location inference (Hoban et al. 2012).

208 *Why should we improve genetic data archiving practices in ecology and evolution?*
209        The FAIR guiding principles are the foundation for good, transparent, and
210 reproducible science. A straightforward demonstration of this is where open data have
211 been used to identify scientific misconduct, some of which impeded evolutionary
212 understanding (e.g. Kozlov 2022). Datasets are often ultimately financed by taxpayers,
213 making public releases an ethical - often even a legal - obligation to ensure the full value
214 of data is obtained. Collectively, data cost many hundreds of millions of dollars to produce,
215 which without archive enrichment, will have to be unnecessarily re-spent to generate the
216 data anew (Crandall et al. 2023). Furthermore, due to the rapid pace of biodiversity loss
217 (e.g. Ceballos et al. 2015), which include genetic diversity decline (Leigh et al. 2019),
218 local extinctions may make regeneration of data impossible, rendering the data
219 irreplaceable and priceless. Existing genetic data also represent an invaluable baseline
220 against which to compare future measurements (i.e. for monitoring of genetic diversity,
221 Jensen and Leigh 2023).
222        Further arguments for data archiving involve benefit sharing and the rights of local
223 communities and local scientists to access data generated from specimens within their
224 country or region. Thus, the CARE principles (Collective benefits, Authority to control,
225 Responsibility, and Ethics; Carroll et al. 2020) could be considered in data generation,
226 archiving, and repurposing of any genetic and genomic data. Emerging benefit-sharing
227 requirements, such as those put forth in the Nagoya Protocol and being developed by the
228 CBD (e.g. Digital Sequence Information or "DSI"; Scholz et al. 2022), are becoming a
229 legal requirement (Marden et al. 2021). This is particularly pertinent to ecology and
230 evolution where researchers often work internationally (e.g. Bhaumik 2023; Miller et al.
231 2023).
232        Researchers themselves can benefit professionally from publicly archiving data.
233 Open datasets can enhance the scholarly recognition of individual research efforts,
234 because data releases with DOI identifiers and data papers can be cited (e.g.
235 MacroPopGen, Lawrence et al. 2019). The increasing popularity of data papers, journals
236 publishing data releases (e.g. Chavan and Penev 2011), and meta-data papers (Raciti et
237 al. 2018) is an early sign that accurate data archiving can benefit individual scientists and
238 the community. Researchers could also benefit from the advancement of their field;
239 synthesis is a powerful tool that has successfully tested pressing big-picture questions in
240 ecology and evolution (Halpern et al. 2020).
241
242 *Best practices for FAIR genetic and genomic data archiving*
243        The most widely used and available genetic data types in molecular ecology and
244 evolution are: 1) barcoding/gene sequences (e.g. mitochondrial cytochrome oxidase, the
245 major histocompatibility complex), 2) microsatellite genotypes, and 3) genomic read data
246 (i.e. unaligned high throughput sequences and SNPs). These come in a constellation of
247 software-specific formats (Lischer and Excoffier 2011; Adamack & Gruber 2014) and, due
248 to lack of standardization, open genetic data repositories contain most of these formats.

249     While there are several tools to convert between file formats (GUI-tools: PGDspider,
250 Lischer and Excoffier 2011; Formatomatic, Manoukis 2007; command line: vcftools, plink,
251 R packages: 'adegenet', Jombart et al. 2008), conversions are time-consuming and often
252 need to be customized for each dataset. Understanding and working with each file format
253 also requires specialist knowledge. Consequently, the lack of a standard archived format
254 limits FAIR data reusability. Due to fundamental differences in data types, file sizes and
255 formats used, a single genetic data file format is unrealistic. However, a single file type
256 for each data type is possible and would be a significant advancement.

257     Unlike other genetic data types, gene sequences are somewhat standardised on
258 archives as FASTA files, and we recommend maintaining this approach. However, many
259 gene sequences lack essential metadata to allow their reuse. It is important that authors
260 archiving gene sequences include the minimum metadata needed to interpret their data
261 (Box 4) otherwise, archives are impossible to reuse (e.g. non-georeferenced sequences
262 in GenBank; Gratton et al. 2017).

263     For microsatellite data, based on its persistent popularity and flexibility, we
264 recommend that it is archived in a "STRUCTURE" input file format (Pritchard et al. 2000;
265 Lischer and Excoffier 2011). STRUCTURE input files also have the advantages that they
266 can handle both haploid and diploid genotype data and have an intuitive and simple
267 format that is conducive to editing in R (R Core Team, 2023) or common spreadsheet
268 software without generating formatting errors. Files can be saved as a comma- (.csv) or
269 tab-delimited text file (.txt) with missing alleles clearly coded as NA or "-9". This file format
270 can house (minimal) metadata (geographical coordinates, populations, sample name,
271 phenotypes though it is essential these match with those used in published papers), as
272 well as marker information (i.e. presence of recessive alleles, inter-marker distances,
273 phase information). We note that there are also variations within the STRUCTURE line
274 format, notably the 1 vs 2 lines *per* individual format; as both are accepted by major
275 conversion tools like PGDSpider (Lischer and Excoffier 2011), either is suitable for
276 archiving. However, we recommend use of the single line format to maximize similarity
277 with VCF files.

278     Genomic data is often mandated to be publicly archived as raw read data on
279 INSDC servers ("INSDC" International Nucleotide Sequence Database Collaboration)
280 (Cochrane et al. 2016), or as aligned BAM files for model organisms (Li et al. 2009). What
281 constitutes "raw" read data can be highly variable, ranging from completely unprocessed
282 files containing several individuals, demultiplexed read files, cleaned files (i.e. with low-
283 quality reads or individuals removed), to error-corrected files (e.g. in ancient DNA)
284 (Mallick et al. 2023). In contrast to microsatellite data, the variable archiving of genomic
285 data means basic error removal, sample delimitation, and genotype calls are not expected
286 to be present in archived data. Ideally, sequencing read data should be archived as
287 demultiplexed read files. A bioinformatic pipeline can also be challenging to reproduce
288 because there are chronic issues surrounding open code archiving that make it hard to

289 know exactly what parameters were applied, tool versions used, or even to have access
290 to custom scripts (further detailed in: Gomes et al. 2022; Jenkins et al. 2023). Even if a
291 pipeline is accessible, version changes of reference genomes or software programs
292 quickly make reproducing a pipeline impossible. Thus, archiving FAIR genomic genotype
293 files in addition to demultiplexed sequencing reads would greatly improve the Open Data
294 compliance and reusability of genomic data.
295        Processed VCF files containing genotype calls (or genotype likelihoods) are
296 standard for genomic analyses and could be archived in parallel with raw read files
297 (though notably this is not possible on INSDC). Though such processed files are not
298 currently widely archived, the practice is becoming more common. Standardization of
299 exactly which variant file is archived also needs consideration. Maximum reusability would
300 be achieved if the archived file represents the least processed SNP or genotype likelihood
301 calls. Specifically, unfiltered genotypes pruned only for basic errors (e.g. technical faults,
302 known contaminated samples), with headers retained to allow for easy assessment of the
303 bioinformatic steps applied. Notably, archiving genomic genotype files could allow non-
304 bioinformatic wildlife managers to repurpose genomic data for analyses and enable
305 researchers without High-Performance Computing access to work with genomic-derived
306 data. Furthermore, this would limit the non-negligible energy, storage and ultimately
307 emissions costs associated with reanalysing genomic data (Grealey et al. 2022).
308

**Box 4: Archiving metadata**

Metadata describing where, when, how and by whom genotype or sequence data was created are invaluable for making genetic data FAIR. There are currently two genomic metadata standards: the Darwin Core standard for biodiversity data (Wieczorek et al. 2012) and the Minimum Information about any(x) Sequence (MIxS) standard (Field et al. 2008). Both standards have cross-mapped terms that overlap (summarized below). What metadata to archive will vary by sample type, project goals, and what researchers deem important (Figure 2). At a minimum, we suggest that authors provide the required (solid lines) and recommended (dotted lines) categories represented in Figure 2 to ensure valuable context. To report metadata not covered here, we also recommend using Darwin Core or MixS standards terms to guarantee FAIRness. Note: Darwin Core contains terms that can handle geologic context of special samples, such as ancient DNA, where metadata related to sampling events generally does not reflect the conditions of the sampled individual before death. As discussed in Box 3, sensitive data should be withheld to ensure it is protected. This can be denoted with the terms "informationWithheld", "dataGeneralizations" or "coordinateUncertaintyInMeters".

The 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. They can thus be used to quickly join the two data types. 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 on

file repositories unless it is archived with the genotype data directly. 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, as described above. To aid automated retrieval, authors should avoid using symbols, special characters, and/or colour-based cell codes.
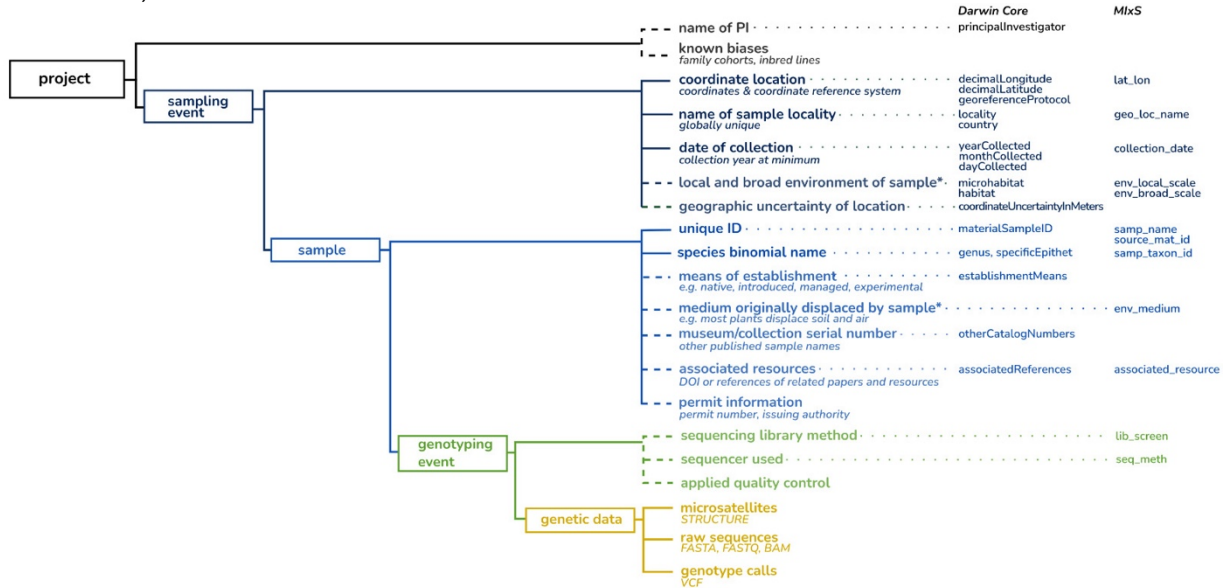


*Figure 2: Metadata can be viewed in a hierarchical manner based on how they were created. The required (solid line) and recommended (dashed line) metadata terms that would improve publicly archived genetic and genomic data reuse potential. Terms denoted with * should use controlled vocabulary from the Environment ontology ("ENVO", Buttigieg et al. 2013). Note: these 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.*

309

310   *Context is key: missing metadata renders most archived data useless*

311       Metadata are a crucial aspect of ensuring genetic data adhere to the FAIR
312   principles (Wilkinson et al. 2016) because the context they provide vastly increases their
313   potential reuses. Genetic metadata record the material and processes that were used for
314   the creation of the genetic data, and can be viewed in a hierarchical manner based on
315   how they were created: 1) sampling events, which include temporal, spatial, and
316   methodological metadata (e.g. year, coordinates). Each sampling event can give rise to
317   many 2) biological samples, each of which have taxonomic, biological, and
318   methodological metadata (e.g. genus, environmental medium, sample preservative).
319   Samples may have many 3) tissues, which might have different biological attributes (e.g.
320   different expression of genes) and may be 4) subject to several different genotyping
321   protocols, which have methodological metadata (e.g. library protocol, Field et al. 2008,
322   Deck et al. 2017, Hassenrück et al. 2021, Crandall et al. 2023).

323    It is standard to include taxonomic metadata (species and genus) in archives, but
324  this is often not sufficient for reuse. The minimum required and recommended metadata
325  are shown in Box 4, without which archived data are often functionally useless and could
326  lead to incorrect inferences. Should key metadata be unavailable to authors we suggest
327  they provide as much information as possible to increase the chance that data will be
328  found and re-used profitably. Currently, publicly accessible metadata are often housed in
329  non-standardized file formats, archived with non-standard terms or present only in
330  published manuscripts and supplementary files. These can take a significant amount of
331  time to access, reformat, or convert for reuse (Crandall et al. 2023). As a result, great
332  efforts have been made to retrospectively georeference existing genetic data to improve
333  their reusability (e.g. Miraldo et al. 2016; Crandall et al. 2023), but this often relies on
334  inference (e.g. inferring coordinates from place names) leaving significant room for error
335  or lost resolution. Thus, we would encourage authors to enrich the public metadata of
336  their data archives and ensure that the metadata included in publications is also present
337  in the data archive.
338
339  *Special considerations when working with important species or indigenous communities*
340    CARE principles (described above, Carroll et al. 2020) are important
341  considerations for data archiving when data are from a culturally important species or
342  indigenous community territory. What steps researchers need to follow will be situation-
343  specific and should be developed in conjunction with interested parties. To ensure these
344  requirements are upheld, data-generating authors should include specific benefit-sharing
345  statements in the publications themselves and in the data archives. This should contain
346  contextual metadata within the statement, for instance provenance information,
347  community names, and also clearly outline community-granted permissions for reuse and
348  circulation. Links to biocultural notices created by researchers and endorsement labels
349  issued by indigenous peoples should be stored within each sample's metadata. When
350  reusing such data, researchers should also follow the ethical repurposing guidelines
351  above and discuss planned analyses with interested parties. Attribution and citation of the
352  original datasets in resulting manuscripts, and dissemination of results to the communities
353  involved could help ensure that cultural authority and sovereignty over such data are
354  recognized (e.g. McCartney et al. 2023), and that data are not reused inappropriately.
355
356  *Which data repository should researchers use?*
357    Centralized infrastructure already exists for genetic sequence data storage
358  (INSDC) that makes finding and accessing data straightforward. Such databases are now
359  impossible to replace and should continue to be used. However, these databases are
360  designed to store only sequence data (e.g. raw reads, gene sequences, whole genomes
361  or transcriptomes) and their metadata (e.g. BioSamples). Genotype data are not stored

362 in sequence databases and there is limited established guidance or storage conventions
363 for them.

364      Currently, genotype data are often stored in multi-purpose Open Data repositories
365 (DRYAD, Zenodo, and increasingly FigShare). However, genetic data can quickly get lost
366 among many of the other data types archived in multi-purpose repositories, where
367 researchers can find everything from non-peer-reviewed ecological survey data (e.g.
368 Shaikh 2014) to violent crime statistics (e.g. Gonzales 2010). Local rules and repository
369 cost barriers (i.e. archiving fees) make it currently impossible to advocate for a single
370 existing database for all genotype data. We note that there are cross-database
371 interoperable search platforms that enable users to search multiple data repositories at
372 once (for example DataONE). However, this functionality is not guaranteed and database
373 linking has failed in the past (Chloé Schmidt *pers. comm.*). There is a need for a free inter-
374 government supported public database specifically for archiving genotype data
375 (microsatellites and SNPs).

376      In lieu of a dedicated repository, researchers can take a few key steps to ensure
377 the findability of genotype data. At a minimum, researchers should ensure the database
378 links their data to their publication. The archiving researcher should include the key
379 metadata fields in Box 4 in the database description and/or the title to aid findability. We
380 also encourage including marker type as a keyword (e.g. "microsatellite" or "SNP") and
381 key geographical descriptors (e.g. "Kruger National Park") to make searching for data
382 more straightforward. Researchers could also consider linking genotypic data to
383 "metadatabases" that keep sample-level metadata in structured, searchable format,
384 enabling users to track samples from the point of collection. These tools can also facilitate
385 upload to the SRA, thereby making the data much more FAIR through structured queries
386 of the metadatabases or INSDC (INSDC BioProjects and BioSamples, Barrett et al. 2021;
387 Genomic Observatories MetaDatabase (GEOME), Deck et al. 2017; Collaborative Open
388 Plant Omics (COPO), Shaw et al. 2020).

389      There are also important database features that researchers should seek out when
390 archiving their genotype data. Researchers should look for a free (or affordable) FAIR
391 compliant Trustworthy Digital Repository (Wilkinson et al. 2016) because they capture
392 key accessibility criteria by definition (compliant repositories are listed on the Registry of
393 Research Data Repositories; Pampel et al. 2013). Institution-specific databases (i.e.
394 university or research institute level) are less desirable because they rarely produce DOIs
395 for data citation, are not easily accessible (e.g. require a password), and might suddenly
396 become depreciated or unsupported.

397      The researcher community could also request new features within existing
398 databases that facilitate genotype data accessibility. The Web of Science's "associated
399 data" link is a notable advance (Web of Science, 2018), as is the soon mandatory
400 metadata (sample location, collection date) for BioSample packages (DDBJ 2023). A
401 desirable additional feature, which would benefit multiple scientific disciplines, is an

402  automatic identifier for retracted data and/or data associated with retracted articles. As of
403  writing, datasets found to be fraudulent from retracted papers remain on servers with no
404  clear notification that the publications was retracted (e.g. Dryad, Costa-Pereira and Pruitt
405  2019). Similarly data found to be erroneous remains on sequence databases (e.g.
406  Genbank, compiled by van den Burg and Vieites 2022) posing a huge challenge to
407  researchers that automate data collection for repurposing. Researchers could benefit
408  from an easy and anonymous way to notify data editors or database curators if they
409  encounter incomplete non-FAIR compliant archives, who should then be responsible for
410  formally rectifying in a harmonious manner.
411
412  *The role of funding bodies and universities in increasing data archiving*
413      University libraries, funding bodies, scientific journals, and data repositories could
414  also take on a greater responsibility to ensure FAIR data archiving. Funding bodies can
415  facilitate data archiving by continuing to mandate Open Data (e.g. National Institute of
416  Health, 2020; European Commission 2016), which have undoubtedly driven an increase
417  in accessibility. However, funding bodies need to support researchers by reviewing or
418  assisting in data management plans (e.g. UKRI 2013), reviewing archived data
419  accessibility and integrity, paying data archiving fees, and offering data archiving
420  educational resources or training. We would specifically encourage funding bodies to
421  ensure future projects budget time for data archiving in their project plan and reward
422  researchers with an established history of Open Data in any field through positively
423  valuing data products or dataset citations.
424      Universities and other science organizations could play a greater role in Open Data
425  through hiring data "stewards" or "librarians" familiar with ecological and evolutionary
426  genetic data. The tasks of data stewards include supporting researchers writing data
427  management plans, identifying suitable databases for archiving, and ensuring dataset
428  longevity through file format conversion (Peng et al. 2016). Notably, data stewards may
429  not be able to archive data directly due to lack of resources and the specialist knowledge
430  required.
431      Science organizations and funding bodies can further foster Open Data by offering
432  data management education (e.g. short courses and training) for both students and
433  career scientists of all disciplines (e.g. Toelch and Ostwald 2018). Few ecology, evolution,
434  or life sciences researchers have received any formal introduction into the importance of
435  Open Data nor in correct data archive practices. Scientific departments could also reward
436  researchers who archive their data or whose data have been reused. Datasets and their
437  reuse (number of views/downloads) can be credited as scientific products on a
438  researcher's *Curriculum Vitae or Research Record* and used during hiring, promotion and
439  tenure decisions.
440
441

*Scientific journal and reviewer roles to ensure Open Data compliance*

Scientific journals can facilitate Open Data by ensuring data are archived on a FAIR compliant Trustworthy Digital Repository before final acceptance of an article (Jenkins et al. 2023, Wilkinson et al. 2016). Journals could also check that data links are activated and that authors have not added reuse clauses or unjustifiable embargos that impede the repurposing of Open Data (see Thrall et al. 2023).

Data are often made *accessible upon publication* with links activated when papers are *in press*. However, this makes it impossible for journals to assess data presence and support archiving. A shift to data accessible *upon submission* is needed, particularly at the resubmission stage when papers are close to acceptance (Thrall et al. 2023). Alternatively, journals could make the final acceptance dependent on data accessibility and FAIR data compliance. Scientists concerned with data being accessed prior to publication should note that several databases offer non-public shareable links that can prevent reuse before publication acceptance.

Journals also have a role to play in improving essential metadata accessibility, which can be easily implemented by having a table of standardized terms that authors must fill out and/or ensure sample names are meaningful (see Box 4). While sample information can be included in supplementary material, versioning issues may arise if metadata are in multiple places. Thus, data editors for journals could ensure all data derived from the same sample are linked (i.e. same name) and key differences (e.g. resequencing with a new technology) highlighted. Importantly, as stated above, we recommend that journals also inform data repositories if papers have been retracted so that the dataset can be demarcated as such (though not removed).

While peer reviewers should not be tasked with ensuring data archiving, they are in a key position to help advance Open Data through a small number of tasks. Initially with novel genetic datasets, authors often need to check the integrity of the data (e.g. checking for contamination) and reviewers could consider asking for evidence of this (e.g. mapping statistics or quality, van den Burg and Vieites 2022). Reviewers could also check archived data files to ensure they are not corrupt, contain the correct number of markers or loci, and contain basic metadata. For repurposed data papers, reviewers can ensure that datasets are cited correctly (see Cousjin et al. 2018). Reviewers could also examine author contribution statements and report to the editor cases where the data-generating authors have not received equal accreditation (Nature 2022).

*Rectifying past mistakes - enriching archived data*

An important step many of us can take to advance Open Data, is to improve metadata archives or archiving inaccessible genetic data. For example, GEOME has successfully run remote datathons to enrich genetic metadata archives (Crandall et al. 2023). We encourage authors to similarly enrich metadata in old data archives, to archive inaccessible genetic datasets, and/or expand on what was archived (e.g. archive all

482  mtDNA haplotypes rather than only unique haplotypes). Although old genetic marker
483  types may be regarded as being of low value to some authors, when combined with other
484  datasets (as in macrogenetics, Leigh et al., 2021), they can be highly informative and can
485  even provide baselines for important biodiversity protection assessments (e.g. Figuerola-
486  Ferrando et al. 2023; Schmidt et al. 2022).
487        Data enrichment initiatives could be run at the Department (similar to MoveBank,
488  Max Planck Institute of Animal Behavior 2023), University library, or country level (e.g.
489  GenDiB and CIEE Living Data), with support from students or technicians to upload data.
490  Such retroactive data archives could even be collaboratively published as a "resource"
491  paper (similar to those in Box 1). These datasets could then support mandated CBD
492  reporting (Hoban et al. 2020), inform local conservation (e.g. Beninde et al. 2022), and
493  identify interesting scientific opportunities (e.g. resampling populations after extreme
494  events, Jensen & Leigh 2022).
495
496  **Perspectives**
497

---

**Box 5: Five take-home messages to improve genetic data archives**

1) Archiving genetic and genomic data in standardized file formats will facilitate reuse (i.e. microsatellites in STRUCTURE; sequences or barcodes in FASTA; SNPs in VCF; Genotype likelihoods in VCF; raw genomic data as demultiplexed FASTQ files).
2) 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*). Include additional contextual metadata when needed to interpret data correctly.
3) Carefully archive data from sensitive species and those affected by the CARE principles to ensure that metadata do not endanger the species, their habitats, or landowner relationships.
4) There is no centralized database for genotype data but this data has great value. Use keywords on FAIR compliant databases (e.g. Dryad) to improve data accessibility.
5) To help more colleagues follow the FAIR principles, request both formalized data management support and a higher value of open data from research institutes and journals.

---

498
499        We close on the note that genetic diversity is the most fundamental component of
500  biodiversity (Hoban et al., 2023). Despite underlying all levels of biodiversity, the
501  biogeographic patterns in intra-specific genetic diversity are largely understudied and
502  poorly protected (Leigh et al. 2019; Figuerola-Ferrando et al. 2023). Consequently,
503  perhaps the most exciting potential of improved archiving is that we can reach research
504  scales beyond what any single research group could achieve. With data spanning such
505  vast spatial and taxonomic scales, open genetic data is pivotal to whole new areas of
506  research and conservation that would have previously been unimaginable due to logistic,
507  cost, or expertise issues. Similar to data collected as part of long-term ecological

508 monitoring programs, publicly archived genetic data is likely to only become more
509 valuable and versatile as it accumulates. The potential of public genetic data is pertinent
510 and timely due to the recently signed United Nations Kunming-Montreal Global
511 Biodiversity Framework which includes commitments by 192 countries to conserve and
512 restore genetic diversity within and among species' populations, and to monitor and report
513 on progress towards that commitment within the decade (Hoban et al. 2023b). Better
514 archiving practices are likely to be central to meet these targets. Although new archiving
515 infrastructure would undoubtedly enhance our ability to do this research, we feel the steps
516 we propose (see Box 5) are achievable with the currently available resources and in the
517 rapid timescale needed.
518
519
520
521

535

536 **Data accessibility statement:**
537
538 The data underpinning Box 1 is available for reviewers and will be accessible upon publication.
539

540 **Conflicts of interest:**
541
542 The authors declare no conflicts of interest.
543

544 **References**
545   Adamack A.T., Gruber B. (2014) Pop Gen Report: Simplifying Basic Population Genetic
546   Analyses in R. Edited by Stephane Dray. *Methods in Ecology and Evolution* 5(4): 384–87.
547   https://doi.org/10.1111/2041-210X.12158.
548

549    Baack S. (2015) Datafication and Empowerment: How the Open Data Movement Re-
550    Articulates Notions of Democracy, Participation, and Journalism. *Big Data & Society* 2(2):
551    205395171559463. https://doi.org/10.1177/2053951715594634.
552

553    Barrett T., Clark K., Gevorgyan R., Gorelenkov V., E. Gribov, Karsch-Mizrachi I., Kimelman M.,
554    et al. (2012) BioProject and BioSample Databases at NCBI: Facilitating Capture and
555    Organization of Metadata. *Nucleic Acids Research* 40(1):57–63.
556    https://doi.org/10.1093/nar/gkr1163.
557

558    Beninde J., Toffelmier E.M., Andreas A., Nishioka C., Slay M., Soto A., Bueno J.P., et al.
559    (2022) CaliPopGen: A Genetic and Life History Database for the Fauna and Flora of
560    California. *Scientific Data* 9(1):380. https://doi.org/10.1038/s41597-022-01479-z.
561

562    Blanchet S., Prunier J.G., De Kort H. (2017) Time to Go Bigger: Emerging Patterns in
563    Macrogenetics. *Trends in Genetics* 33(9)579–80. https://doi.org/10.1016/j.tig.2017.06.007.
564

565    CBD (2022) Decision Adopted By The Conference Of The Parties To The Convention On
566    Biological Diversity. https://www.cbd.int/doc/decisions/cop-15/cop-15-dec-05-en.pdf
567

568    Chapman AD (2020) Current Best Practices for Generalizing Sensitive Species Occurrence
569    Data. Copenhagen: GBIF Secretariat. https://doi.org/10.15468/doc-5jp4-5g10.
570

571    Chavan V., Penev L. (2011) The Data Paper: A Mechanism to Incentivize Data Publishing in
572    Biodiversity Science. *BMC Bioinformatics* 12(S15): S2. https://doi.org/10.1186/1471-2105-12-
573    S15-S2.
574

575    Clarke K. C. A. (2016) Multiscale Masking Method for Point Geographic Data. *International*
576    *Journal of Geographical Information Science* 30(2): 300–315.
577    https://doi.org/10.1080/13658816.2015.1085540.
578

579    Cochrane G., Cook C.E., Birney E. (2012) The Future of DNA Sequence Archiving.
580    *GigaScience* 1(1): 2047–217X–1–2. https://doi.org/10.1186/2047-217X-1-2.
581

582    Cochrane G., Karsch-Mizrachi I., Takagi T., International Nucleotide
583    Sequence Database Collaboration. (2016) The International Nucleotide Sequence Database
584    Collaboration, Nucleic Acids Research, 44(1) 48–D50, https://doi.org/10.1093/nar/gkv1323
585

586    Costa-Pereira R., Pruitt J. (2020) Data from: Behavior, morphology, and microhabitat use:
587    what drives individual niche variation? Dryad, Dataset, https://doi.org/10.5061/dryad.bd26mq0
588

589    Crandall E.D., Toczydlowski R.H., Liggins L., Holmes A.E., Ghoojaei M., Gaither M.R, Wham
590    B.E., et al. (2023) Importance of Timely Metadata Curation to the Global Surveillance of
591    Genetic Diversity. *Conservation Biology* 37(4): e14061. https://doi.org/10.1111/cobi.14061.
592

593    Cross T.B., Schwartz M.K., Naugle D.E., Fedy B.C., Row J.R., Oyler-McCance S.J. (2018)
594    The Genetic Network of Greater Sage-grouse: Range-wide Identification of Keystone Hubs of
595    Connectivity. *Ecology and Evolution* 8(11): 5394–5412. https://doi.org/10.1002/ece3.4056.
596

597    DDBJ (2023) Including Sample Location and Collection Date and Time for BioSample
598    submissions Including Sample Location. https://www.ddbj.nig.ac.jp/news/en/2023-05-02-e.html
599

600    Deck J., Gaither M.R., Ewing R., Bird C.E., Davies N., Meyer C., Riginos C., Toonen R.J.,
601    Crandall E.D. (2017) The Genomic Observatories Metadatabase (GeOMe): A New Repository
602    for Field and Sampling Event Metadata Associated with Genetic Samples'. *PLOS Biology*
603    15(8): e2002925. https://doi.org/10.1371/journal.pbio.2002925.
604

605    European Commission (2016) Guidelines on FAIR data management in Horizon 2020
606    https://ec.europa.eu/research/participants/data/ref/h2020/grants_manual/hi/oa_pilot/h2020-hi-
607    oa-data-mgt_en.pdf
608

609    Espindola S., Vázquez-Domínguez E., Nakamura M., Osorio-Olvera L., Martínez-Meyer E.,
610    Myers E.A., Overcast I., Reid B.N., Burbrink F.T. (2022) Complex Genetic Patterns and
611    Distribution Limits Mediated by Native Congeners of the Worldwide Invasive Red-eared Slider
612    Turtle. *Molecular Ecology* 31(6): 1766–82. https://doi.org/10.1111/mec.16356.
613

614    Fairbairn, D.J. (2011) The Advent Of Mandatory Data Archiving. *Evolution* 65(1): 1–2.
615    https://doi.org/10.1111/j.1558-5646.2010.01182.x.
616

617    Field D., Amaral-Zettler L., Cochrane G., Cole J.R., Dawyndt P., Garrity G.M., Gilbert J., et al.
618    (2011) The Genomic Standards Consortium. *PLoS Biology* 9(6): e1001088.
619    https://doi.org/10.1371/journal.pbio.1001088.
620

621    Frank R. D., Kriesberg A., Yakel E., Faniel I. M. (2016) Looting Hoards of Gold and Poaching
622    Spotted Owls: Data Confidentiality among Archaeologists & Zoologists. *Proceedings of the*

623     *Association for Information Science and Technology* 52(1): 1–10.
624     https://doi.org/10.1002/pra2.2015.145052010037.
625
626     Gabelica M., Bojčić R., Puljak L. (2022) Many Researchers Were Not Compliant with Their
627     Published Data Sharing Statement: A Mixed-Methods Study. *Journal of Clinical Epidemiology*
628     150: 33–41. https://doi.org/10.1016/j.jclinepi.2022.05.019.
629
630     Gomes D.G.E., Pottier P., Crystal-Ornelas R., Hudgins E.J., Foroughirad V., Sánchez-Reyes
631     L.L., Turba R., et al. (2022) Why don't we share data and code? Perceived barriers and
632     benefits to public archiving practices. *Proceedings of the Royal Society B*. 289(1987): 2022111
633     https://doi.org/10.1098/rspb.2022.1113
634
635     Gonzalez L. (2010). Sexual crime in Colombia 2010-2022. figshare. Dataset.
636     https://doi.org/10.6084/m9.figshare.21937154.v1
637
638     Gratton P., Marta S., Bocksberger G., Winter  M., Trucchi E., Kühl  H. (2017) A World of
639     Sequences: Can We Use Georeferenced Nucleotide Databases for a Robust Automated
640     Phylogeography? *Journal of Biogeography* 44(2): 475–86. https://doi.org/10.1111/jbi.12786.
641
642     Grealey J., Lannelongue L., Saw W-Y, Marten J., Méric G., Ruiz-Carmona S., Inouye M.
643     (2022) The Carbon Footprint of Bioinformatics. *Molecular Biology and Evolution* 39(3):
644     msac034. https://doi.org/10.1093/molbev/msac034.
645
646     Günther T., Coop G. (2013) Robust Identification of Local Adaptation from Allele Frequencies.
647     *Genetics* 195(1): 205–20. https://doi.org/10.1534/genetics.113.152462.
648
649     Hassenrück C., Poprick T., Helfer V., Molari M., Meyer R., Kostadinov I. (2021) FAIR Enough?
650     A Perspective on the Status of Nucleotide Sequence Data and Metadata on Public Archives,
651     bioRxiv 2021.09.23.461561; doi: https://doi.org/10.1101/2021.09.23.461561
652
653     Hoban S.M., Borkowski D.S., Bros S.L., McCleary T.S., Thompson L.M., McLachlan J.S.,
654     Pereira M.A., Schlarbaum  S.E., Romero-Severson  J. (2010) Range-Wide Distribution of
655     Genetic Diversity in the North American Tree *Juglans cinerea*: A Product of Range Shifts, Not
656     Ecological Marginality or Recent Population Decline. *Molecular Ecology* 19(22):4876–91.
657     https://doi.org/10.1111/j.1365-294X.2010.04834.x.
658
659     Hoban S., McCleary T.S., Schlarbaum S.E., Anagnostakis  S.L, Romero-Severson J. (2012).
660     Human-Impacted Landscapes Facilitate Hybridization between a Native and an Introduced

661    Tree. *Evolutionary Applications* 5(7): 720–31. https://doi.org/10.1111/j.1752-
662    4571.2012.00250.x.

663

664    Hoban S., Bruford M., D'Urban Jackson J., Lopes-Fernandes  M., Heuertz  M., Hohenlohe  P.
665    A, Paz-Vinas  I., et al. (2020) Genetic Diversity Targets and Indicators in the CBD Post-2020
666    Global Biodiversity Framework Must Be Improved. *Biological Conservation* 248: 108654.
667    https://doi.org/10.1016/j.biocon.2020.108654.

668

669    Hoban S., da Silva, J.M., Mastretta-Yanes A., Grueber C.E., Heuertz, M., Hunter M.E.,
670    Mergeay, J., et al. (2023). Monitoring status and trends in genetic diversity for the Convention
671    on Biological Diversity: an ongoing assessment of genetic indicators in nine countries.
672    *Conservation Letters*, 16, e12953. https://doi.org/10.1111/conl.12953

673

674    Hostler T.J. (2023) The Invisible Workload of Open Research. *Journal of Trial and Error*
675    https://doi.org/10.36850/mr5.

676

677    Huang X., Hawkins B.A., Lei F., Miller G. L, Favret C., Zhang R., Qiao G. (2012) Willing or
678    Unwilling to Share Primary Biodiversity Data: Results and Implications of an International
679    Survey. *Conservation Letters* 5: 399–406. https://doi.org/10.1111/j.1755-263X.2012.00259.x.

680

681    Jenkins G.B., Beckerman A.P., Bellard C., Benítez-López A., Ellison A.M., Foote C.G., Hufton
682    A.L., et al. (2023) Reproducibility in Ecology and Evolution: Minimum Standards for Data and
683    Code. *Ecology and Evolution* 13(5): e9961. https://doi.org/10.1002/ece3.9961.

684

685    Jensen E. L., Leigh D. M. (2022) Using Temporal Genomics to Understand Contemporary
686    Climate Change Responses in Wildlife. *Ecology and Evolution* 12:e9340
687    https://doi.org/10.1002/ece3.9340.

688

689    Jombart T. (2008) *Adegenet*: A R Package for the Multivariate Analysis of Genetic Markers.
690    *Bioinformatics* 24(11): 1403–5. https://doi.org/10.1093/bioinformatics/btn129.

691

692    Kozlov M. (2022) How A Spider-Biology Scandal Upended Researchers' Lives. *Nature* 608:
693    658-659  doi: https://doi.org/10.1038/d41586-022-02156-2

694

695    Lawrence E. R., Benavente J. N., Matte J.-M., Marin K., Wells Z.R.R., Bernos T. A., Krasteva
696    N., et al. (2019) Geo-Referenced Population-Specific Microsatellite Data across American
697    Continents, the MacroPopGen Database. *Scientific Data* 6(1):14.
698    https://doi.org/10.1038/s41597-019-0024-7.

699

Leigh D.M., Hendry A.P., Vázquez-Domínguez E., Friesen V.L. (2019) Estimated Six per Cent Loss of Genetic Variation in Wild Populations since the Industrial Revolution. *Evolutionary Applications* 12(8):1505–12. https://doi.org/10.1111/eva.12810.

Leigh D.M., van Rees C.B., Millette K.L., Breed M.F., Schmidt C., Bertola L.D., Hand B.K., et al. (2021) Opportunities and Challenges of Macrogenetic Studies. *Nature Reviews Genetics* 22(12): 791–807. https://doi.org/10.1038/s41576-021-00394-0.

Li H., Handsaker B., Wysoker A., Fennell T., Ruan J., Homer N., Marth G., Abecasis G., Durbin R., and 1000 Genome Project Data Processing Subgroup. (2009) The Sequence Alignment/Map Format and SAMtools. *Bioinformatics* 25(16): 2078–79. https://doi.org/10.1093/bioinformatics/btp352.

Lindenmayer D., Scheele B. (2017) Do Not Publish. *Science* 356(6340): 800–801. https://doi.org/10.1126/science.aan1362.

Lischer H.E.L., Excoffier L., (2012) PGDSpider: An Automated Data Conversion Tool for Connecting Population Genetics and Genomics Programs. *Bioinformatics* 28(2): 298–99. https://doi.org/10.1093/bioinformatics/btr642.

Lowe A.J., Smyth A.K., Atkins K., Avery R., Belbin L., Brown N., Budden A.E., et al. (2017) Publish Openly but Responsibly. *Science* 357(6347): 141–141. https://doi.org/10.1126/science.aao0054.

Mallick S., Micco A., Mah M., Ringbauer H, Lazaridis I., Olalde I., Patterson N, Reich D. (2023) The Allen Ancient DNA Resource (AADR): A Curated Compendium of Ancient Human Genomes'. biorxiv. https://doi.org/10.1101/2023.04.06.535797.

Marden E., Abbott R.J., Austerlitz F., Ortiz-Barrientos D., Baucom R.S., Bongaerts P., Bonin A., et al. (2021) Sharing and Reporting Benefits from Biodiversity Research. *Molecular Ecology* 30(5): 1103–7. https://doi.org/10.1111/mec.15702.

Max Planck Institute of Animal Behavior. 2023. Movebank Attribute Dictionary. British Oceanographic Data Centre, Natural Environment Research Council Vocabulary Server. http://vocab.nerc.ac.uk/collection/MVB/current

736  Mills J.A., Teplitsky C., Arroyo B., Charmantier A., Becker P. H., Birkhead T.R., Bize P., et al.
737  (2015) Archiving Primary Data: Solutions for Long-Term Studies. *Trends in Ecology &*
738  *Evolution* 30(10): 581–89. https://doi.org/10.1016/j.tree.2015.07.006.
739
740  Miraldo A., Li S., Borregaard M. K., Flórez-Rodríguez A., Gopalakrishnan S., Rizvanovic M.,
741  Wang Z. et al. (2016) An Anthropocene Map of Genetic Diversity. *Science* 353(6307): 1532–
742  35. https://doi.org/10.1126/science.aaf4381.
743
744  Moore A.J., Mcpeek M.A., Rausher M.D., Rieseberg L., Whitlock M.C. (2010) The Need for
745  Archiving Data in Evolutionary Biology. *Journal of Evolutionary Biology* 23(4): 659–60.
746  https://doi.org/10.1111/j.1420-9101.2010.01937.x.
747
748  Nater A., Greminger M.P., Arora N., van Schaik C.P., Goossens B., Singleton I., Verschoor
749  E.J., et al. (2015) Reconstructing the demographic history of orang-utans using Approximate
750  Bayesian Computation. *Molecular Ecology* 24(2): 310-327. https://doi.org/10.1111/mec.13027
751
752  National Institute for Health (2020) Final NIH Policy for Data Management and Sharing: NOT-
753  OD-21-013 https://grants.nih.gov/grants/guide/notice-files/NOT-OD-21-013.html
754
755  Nature (2022) Time to recognize authorship of open data. *Nature* 604:8 doi:
756  https://doi.org/10.1038/d41586-022-00921-x
757
758  Oyler-McCance S.J., Cross T.B., Row J.R., Schwartz M.K., Naugle D.E., Fike J.A., Winiarski
759  K., Fedy B.C. (2022) New Strategies for Characterizing Genetic Structure in Wide-Ranging,
760  Continuously Distributed Species: A Greater Sage-Grouse Case Study. *PLoS ONE* 17(9):
761  e0274189. https://doi.org/10.1371/journal.pone.0274189.
762
763  Pampel H., Vierkant P., Scholze F., Bertelmann R., Kindling M., Klump, J., Goebelbecker H.-
764  J., et al., (2013) Making Research Data Repositories Visible: The Re3data.Org Registry *PLoS*
765  *ONE* 8,(11): e78080. https://doi.org/10.1371/journal.pone.0078080.
766
767  Paz-Vinas I., Jensen E.L., Bertola L.D., Breed M.F., Hand B.K., Hunter M.E., Kershaw F., et
768  al. (2021) Macrogenetic Studies Must Not Ignore Limitations of Genetic Markers and Scale.
769  *Ecology Letters* 24(6): 1282–84. https://doi.org/10.1111/ele.13732.
770
771  Peng G., Ritchey N.A., Casey K.S., Kearns E.J., Prevette J.L., Saunders D., Jones P.,
772  Maycock T., Ansari S. (2016) Scientific Stewardship in the Open Data and Big Data Era  Roles

773     and Responsibilities of Stewards and Other Major Product Stakeholders. *D-Lib Magazine* 22,

774     no. 5/6. https://doi.org/10.1045/may2016-peng.

775

776     Pike V.L., Cornwallis C.K., Griffin A.S. (2021) Why Don't All Animals Avoid Inbreeding?

777     *Proceedings of the Royal Society B: Biological Sciences* 288(1956): 20211045.

778     https://doi.org/10.1098/rspb.2021.1045.

779

780     Piwowar H.A., Vision T.J., Whitlock M.C. (2011) Data archiving is a good investment. *Nature.*

781     473(7347): 285. doi: 10.1038/473285a..

782

783     Pritchard J.K., Stephens M., Donnelly P. (2000) Inference of Population Structure Using

784     Multilocus Genotype Data. *Genetics* 155(2): 945–59.

785     https://doi.org/10.1093/genetics/155.2.945.

786

787     R Core Team (2023). R: A language and environment for statistical computing. R Foundation

788     for Statistical Computing, Vienna, Austria. https://www.R-project.org/.

789

790     Raciti D., Yook K., Harris T.W., Schedl T., Sternberg P.W. (2018) *Micropublication* :

791     Incentivizing Community Curation and Placing Unpublished Data into the Public Domain.

792     *Database* 2018: bay013 https://doi.org/10.1093/database/bay013.

793

794     Rieseberg L., Vines T., Kane N. (2010) Editorial and Retrospective 2010. *Molecular Ecology*

795     19(1): 1–22. https://doi.org/10.1111/j.1365-294X.2009.04450.x.

796

797     Roche D.G., Kruuk L.E.B., Lanfear  R., Binning S.A. (2015) Public Data Archiving in Ecology

798     and Evolution: How Well Are We Doing? *PLOS Biology* 13(11): e1002295.

799     https://doi.org/10.1371/journal.pbio.1002295.

800

801     Roche D.G., Lanfear R., Binning S.A., Haff T.M., Schwanz L.E., Cain K.E., Kokko H., Jennions

802     M.D., Kruuk L.E.B. (2014) Troubleshooting Public Data Archiving: Suggestions to Increase

803     Participation. *PLoS Biology* 12(1): e1001779. https://doi.org/10.1371/journal.pbio.1001779.

804

805     Row J.R., Doherty K.E., Cross T.B., Schwartz M.K., Oyler-McCance S.J., Naugle D.E., Knick

806     S.T., Fedy B.C. (2018) Quantifying Functional Connectivity: The Role of Breeding Habitat,

807     Abundance, and Landscape Features on Range-wide Gene Flow in Sage-grouse.

808     *Evolutionary Applications* 11(8): 1305–21. https://doi.org/10.1111/eva.12627.

809

810     Ruppert K.M., Kline R.J., Rahman S. (2019) Past, Present, and Future Perspectives of
811     Environmental DNA (EDNA) Metabarcoding: A Systematic Review in Methods, Monitoring,
812     and Applications of Global EDNA. *Global Ecology and Conservation* 17: e00547.
813     https://doi.org/10.1016/j.gecco.2019.e00547.
814

815     Scholz A.H., Freitag J., Lyal C.H.C., Sara R., Cepeda M.L., Cancio I., Sett S., et al. (2022)
816     Multilateral Benefit-Sharing from Digital Sequence Information Will Support Both Science and
817     Biodiversity Conservation. *Nature Communications* 13(1): 1086.
818     https://doi.org/10.1038/s41467-022-28594-0.
819

820     Schmidt C., Domaratzki M., Kinnunen R.P., Bowman J., Garroway C.J. (2022) Continent-Wide
821     Effects of Urbanization on Bird and Mammal Genetic *Diversity. Proceedings of the Royal*
822     *Society B* 287(1920): 20192497. http://dx.doi.org/10.1098/rspb.2019.2497
823

824     Schmidt C., Garroway C.J. (2022) Systemic Racism Alters Wildlife Genetic Diversity.
825     *Proceedings of the National Academy of Sciences* 119(43): e2102860119.
826     https://doi.org/10.1073/pnas.2102860119.
827

828     Schmidt C., Garroway C. J. (2021) The Population Genetics of Urban and Rural Amphibians in
829     North America. *Molecular Ecology* 30(16): 3918–29. https://doi.org/10.1111/mec.16005.
830

831     Schmidt, C; Hoban, S; Jetz, W. (2023) Conservation macrogenetics: harnessing genetic data
832     to meet conservation commitments. Trends in Genetics. In press.
833     https://doi.org/10.1016/j.tig.2023.08.002
834

835     Schumacher, E., Brown, A., Williams, M., Romero-Severson, J., Beardmore, T. & Hoban, S.
836     (2022). Range shifts in butternut, a rare, endangered tree, in response to past climate and
837     modern conditions. *Journal of Biogeography*, 49(5): 866–878. https://doi.org/10.1111/jbi.14350
838

839     Schriml L.M., Chuvochina M., Davies N., Eloe-Fadrosh E.A., Finn R.D., Hugenholtz  P., Hunter
840     C. I., et al. (2020) COVID-19 Pandemic Reveals the Peril of Ignoring Metadata Standards.
841     *Scientific Data* 7(1): 188. https://doi.org/10.1038/s41597-020-0524-5.
842

843     Shaikh A. (2014). Ecology week 4: field sample with animals. figshare. Dataset.
844     https://doi.org/10.6084/m9.figshare.1194651.v1
845

846    Shaw, F., Etuk A., Minotto A., Gonzalez-Beltran A., Johnson D., Rocca-Serra P.,  Laporte  M.-
847    A., et al. (2020) COPO: A Metadata Platform for Brokering FAIR Data in the Life Sciences.
848    *F1000Research* 9: 495. https://doi.org/10.12688/f1000research.23889.1.
849
850    Tedersoo, L., Küngas R., Oras E., Köster  K., Eenmaa H., Leijen Ä., Pedaste M., et al. (2021)
851    Data Sharing Practices and Data Availability upon Request Differ across Scientific Disciplines.
852    *Scientific Data* 8(1): 192. https://doi.org/10.1038/s41597-021-00981-0.
853
854    Thrall P.H., Chase J., Drake J., Espuno N., Hello S., Ezenwa V., Han B., Mori A., Muller-
855    Landau H. (2023) From raw data to publication: Introducing data editing at Ecology Letters.
856    *Ecology Letters* 26: 829-830. https://doi.org/10.1111/ele.14210
857
858    Toczydlowski R. H., Liggins L., Gaither M. R., Anderson T. J., Barton R. L., Berg J. T., Beskid
859    S. G., et al. (2021) Poor Data Stewardship Will Hinder Global Genetic Diversity Surveillance.
860    *Proceedings of the National Academy of Sciences* 118(34): e2107934118.
861    https://doi.org/10.1073/pnas.2107934118.
862
863    Toelch U., Ostwald D. (2018) Digital open science—Teaching digital tools for reproducible and
864    transparent research. *PLoS Biology* 16(7): e2006022. https://doi.org/10.1371/journal.
865
866    UKRI (2013) Data Management Plan: Guidance for Peer Reviewers https://www.ukri.org/wp-
867    content/uploads/2021/07/ESRC-200721-DataManagementPlan-
868    GuidanceforPeerReviewers.pdf
869
870    van den Burg M. P., Vieites D. R. (2023) Bird Genetic Databases Need Improved Curation and
871    Error Reporting to NCBI. *Ibis* 165(2): 472–81. https://doi.org/10.1111/ibi.13143.
872
873    Vines, T.H., Albert A.Y.K., Andrew R.L., Débarre F., Bock D.G., Franklin M.T., Gilbert K.J.,
874    Moore J.-S., Renaut S., Rennison D.J. (2014) The Availability of Research Data Declines
875    Rapidly with Article Age. *Current Biology* 24(1): 94–97.
876    https://doi.org/10.1016/j.cub.2013.11.014.
877
878    Web of Science (2018) Associated Data.
879    https://images.webofknowledge.com/images/help/WOK/hp_associated_data.html
880
881    Whitlock M.C. Data Archiving in Ecology and Evolution: Best Practices. *Trends in Ecology &*
882    *Evolution* 26(2): 61–65. https://doi.org/10.1016/j.tree.2010.11.006.
883

884    Whitlock M.C., Bronstein J.L., Bruna E.M., Ellison A.M., Fox C.W., McPeek M.A., Moore A.J.,
885    et al. (2015) A Balanced Data Archiving Policy for Long-Term Studies. *Trends in Ecology &*
886    *Evolution* 31(2): 84–85. https://doi.org/10.1016/j.tree.2015.12.001.
887
888    Wicherts J.M., Bakker M., Molenaar D. (2011) Willingness to Share Research Data Is Related
889    to the Strength of the Evidence and the Quality of Reporting of Statistical Results. *PLoS ONE*
890    6(11): e26828. https://doi.org/10.1371/journal.pone.0026828.
891
892    Wilkinson M.D., Dumontier M.I, Aalbersberg I.J., Appleton G., Axton M., Baak A., Blomberg N.,
893    et al. (2016) The FAIR Guiding Principles for Scientific Data Management and Stewardship.
894    *Scientific Data* 3(1): 160018. https://doi.org/10.1038/sdata.2016.18.
895
896    Zimmerman S.J., Aldridge C.L., Oh K.P., Cornman R.S., Oyler-McCance S.J. (2019)
897    Signatures of Adaptive Divergence among Populations of an Avian Species of Conservation
898    Concern. *Evolutionary Applications* 12(8): 1661–77. https://doi.org/10.1111/eva.12825.
899
900
901                                   ****** **END**\*****
902
903
904
905
906
907