

1 Data rescue: saving environmental data from extinction

2

3 Ellen K. Bledsoe^{1,2,†,*}, Joseph B. Burant^{1,3,4,†,*}, Gracielle T. Higinio^{1,5,†}, Dominique G. Roche^{1,6},

4 Sandra A. Binning^{1,4}, Kerri Finlay^{1,2}, Jason Pither^{1,7}, Laura S. Pollock^{1,3}, Jennifer M. Sunday^{1,3},

5 Diane S. Srivastava^{1,6,*}

6 Author affiliations

1 The Living Data Project, Canadian Institute of Ecology and Evolution, Vancouver, British Columbia, Canada

2 Department of Biology, University of Regina, Regina, Saskatchewan, Canada

3 Department of Biology, McGill University, Montreal, Quebec, Canada

4 Département de Sciences Biologiques, Université de Montréal, Montréal, Québec, Canada

5 Department of Zoology and Biodiversity Research Centre, University of British Columbia, Vancouver, British Columbia, Canada

6 Department of Biology and Institute for Environment & Interdisciplinary Science, Carleton University, Ottawa, Ontario, Canada

7 Department of Biology and Okanagan Institute for Biodiversity, Resilience, and Ecosystem Services, University of British Columbia, Kelowna, British Columbia, Canada

7

† These co-authors contributed equally to this work.

* Corresponding authors: ellen.bledsoe@weecology.org (EKB), joseph.burant@mcgill.ca (JBB), srivast@zoology.ubc.ca (DSS)

8 Running Headline

9 Data rescue: saving environmental data

10 Abstract

- 11 1. Historical and long-term environmental datasets are imperative to understanding how
12 natural systems respond to our changing world, setting baselines and establishing
13 trajectories of change. Although immensely valuable, these data are ultimately at risk of
14 being lost unless they are actively managed, curated, and eventually archived on data
15 repositories.
- 16 2. The practice of data rescue, which we define as identifying, preserving, and sharing
17 valuable data and associated metadata at risk of loss, is an important means of ensuring
18 the long-term viability and accessibility of such datasets. Improvements in policies and
19 best practices around data management will hopefully limit the future need for data
20 rescue; these changes, however, do not apply retroactively. While the concept of rescuing
21 data is not new, the term lacks a formal definition, is often conflated with other terms
22 (i.e., data reuse), and lacks general recommendations.
- 23 3. Here, we outline seven key guidelines for effective rescue of historically-collected and
24 unmanaged datasets. We discuss how to prioritize which datasets to rescue, form
25 effective data rescue teams, prepare the data and related metadata, and ultimately archive
26 and share the rescued data.
- 27 4. In an era of rapid environmental change, the best policy solutions will require evidence
28 from both contemporary and historical sources. It is, therefore, imperative that we
29 identify and preserve valuable, at-risk environmental data before they are lost to science.

30 Keywords

31 Data archiving, historic data, long-term ecological data, long-term studies, open data, open
32 science, reproducibility, transparency

33 Author contributions

34 EKB, JBB, DGR, and DSS proposed the initial idea for the manuscript; all authors contributed to
35 developing the methods of data rescue we describe and subsequent discussions about the paper.
36 EKB, JBB, and GTH wrote the first draft. DSS created the first draft of the figure. All authors
37 revised the manuscript for publication. The authors declare no competing interests.

38

39 Abstract word count = 219 words

40 Main text word count (excluding abstract, figure captions, boxes, and references) = 4,931 words

41 Number of figures: 2

42 Figure 1 = 80 words

43 Figure 2 = 114 words

44 Number of boxes: 3

45 Box 1 = 519 words (no images)

46 Box 2 = 492 words (including captions; 2 images)

47 Box 3 = 1,391 words (including captions; 4 images)

48 Number of references: 57

49 Why Rescue Data?

50 Data are among the primary units of research and scholarship. Not only are data used to
51 help answer important questions, but they can also be used to inform new lines of inquiry, new
52 testable hypotheses, and future data collection efforts. Observational and experimental data
53 derived from ecology, evolution, conservation and environmental sciences (hereafter
54 environmental data) are essential to establishing historical trajectories of ecosystems (i.e.,
55 baselines; McClenachan et al., 2012), understanding how species and communities respond to
56 environmental change (Gatti et al., 2015), and designing and evaluating the outcomes of
57 management efforts (Hawkin et al., 2013; Willis et al., 2007). Moreover, while data collection is
58 often targeted to a particular population, community, or location, the reuse (i.e., aggregation,
59 collation, and synthesis) of data from different systems and contexts is essential to establishing
60 broader ecological knowledge and informing conservation management. Yet, despite their high
61 value and central role in research, data are often misplaced, filed away, or otherwise rendered
62 unusable, often through poor data management practices (Vines et al., 2014). In their unusable
63 and “at-risk” state, these data represent an egregious waste of resources expended on their
64 collection (Buxton et al., 2021; Box 1). Languishing data, however, also offer an enormous
65 opportunity. **Data rescue**—defined here as the identification, preservation, and sharing of
66 valuable data and associated metadata at risk of loss—has the potential to realize huge benefits
67 for society, especially considering the crucial roles that baseline data play in informing
68 management and policy decisions. The ultimate goal of data rescue is to make previously
69 inaccessible or poorly preserved data available for (re)use, ideally through archiving them in a
70 permanent, publicly accessible, and reusable format.

71 In recent years, there has been a strong push from within the scientific and scholarly
72 communities for increased transparency and openness in the practice of science, including in
73 ecology and evolution (e.g., O’Dea et al., 2021). Calls for more transparency and accessibility in
74 science are not new (e.g., Eamon, 1985); the term “open science” itself was coined more than 20
75 years ago. However, the last decade has seen a surge in general awareness and promotion of
76 open science practices (e.g., open access publishing and open data, code, software, and peer-
77 review) and their benefits (Powers & Hampton, 2019). These initiatives have not been without
78 criticism, with many researchers unsure about sharing their data due to real or perceived
79 concerns about data misuse and loss of control (Roche et al., 2014; Mills et al., 2015; Smith &
80 Roberts, 2016; Stieglitz et al., 2020). Others have acknowledged important caveats to the general
81 appeal for openness (e.g., valid considerations about security, confidentiality, equity, and
82 Indigenous data sovereignty and governance; Borgman, 2018; Walter & Suina, 2018; Lennox et
83 al., 2020; Buck, 2021). Despite the legitimacy of (some of) these concerns, the benefits of data
84 sharing to individuals, the scientific community and the general public are apparent (Powers &
85 Hampton, 2019; Soeharjono & Roche, 2021). And yet, large amounts of data remain private,
86 unavailable for reuse by other scientists, and inaccessible to researchers and the public who
87 ultimately provided the funding and infrastructure for the data’s collection. For example, in a
88 sample of more than 4,000 ecology and evolution papers, only one in five papers (21.5%) had a
89 data availability statement or associated open data (Roche et al., 2021), and less than half of
90 archived datasets in ecology and evolution are reusable (Roche et al., 2015; Roche et al., 2021).

91 Open science initiatives have developed rapidly, and the last few years have seen a rise in
92 the number of institutions, governments, funding agencies, and publishers who have
93 implemented policies that require the open, permanent, and accessible sharing of data (e.g.,

94 FAIR data principles (see *Data sharing* below; Wilkinson et al., 2016), the Ecological Society of
95 America's new [Open Research policy](#), and the European Commission's [OpenAIRE open access](#)
96 [and open data policy](#)). These requirements, and participation by scientists, will enhance our
97 ability to evaluate, reuse, and synthesize increasingly rich and complex ecological data.
98 However, open data policies are not retroactive and, therefore, do relatively little to address the
99 issue of access to and preservation of previously-collected data (Vines et al., 2014). Arguably,
100 data collected prior to the adoption of widespread sharing practices remain a public good, funded
101 by taxpayers and governments, so rescuing datasets to ensure their longevity and accessibility
102 should be seen as an ethical imperative.

103 Here, we present guidelines for implementing data rescue; although we focus on
104 environmental data, our guidelines are applicable more broadly. These guidelines are proposed
105 based on past and ongoing data rescue projects by the Living Data Project, an initiative of the
106 Canadian Institute of Ecology and Evolution (CIEE), which aims to identify and secure
107 vulnerable datasets and bring new life to them through collaborative analysis and synthesis. We
108 include examples using historical (Box 2) and recent data rescue efforts (Box 3). Our hope is that
109 these guidelines will (a) focus attention on the current threats to the usability and integrity of
110 previously-collected data, (b) stimulate broader consideration of the utility of previously-
111 collected datasets for current research efforts, (c) encourage people with access to or knowledge
112 of unarchived data to work towards their preservation, (d) provide a reference for those looking
113 to apply data rescue techniques in the context of their own work, and (e) help foster a strong
114 culture of data stewardship such that data rescue becomes unnecessary in the future.

115 Guidelines for data rescue

116 Imperiled data can be found nearly everywhere, such as non-profit organizations,
117 conservation councils, academic institutions, and government agencies (think: historical data
118 only available on paper records in basement filing cabinets, digitized data stored only on floppy
119 disks, etc.). Finding data to rescue is usually the easiest part of the process; how to implement a
120 successful data rescue mission, however, requires a more strategic approach (Fig. 1). Some of the
121 steps involved in data rescue are closely aligned with recommended practices in research data
122 management (see *Metadata, Data Compilation, Validation, Archiving* and *Sharing* sections
123 below). Several resources have already outlined “best” practices for data collection (Broman &
124 Woo, 2018), management (e.g., British Ecological Society’s “Data Management” Guide, 2018),
125 and archiving (Cook et al., 2001; Whitlock, 2011; White et al., 2013), yet these are written with
126 current or future data collection in mind and do not address historically-collected or unmanaged
127 data. Below, we outline seven key steps for data rescue, from identifying high-priority datasets to
128 archiving and sharing them for (re)use.

129 1. Data prioritization

130 Prioritizing data for rescue requires a consideration of both the scientific value of the data
131 and the potential risk that the data will be lost (Fig. 2). Data of high value and at high risk of
132 being lost should be given highest priority, while data which rank highly along just one of the
133 axes of value and risk should be considered moderate priorities. The concepts of value and risk
134 of loss are, of course, subjective, but there are some general factors to consider when
135 determining these characteristics of a dataset.

136 High-value environmental datasets have some common features. Scale is a key factor, as
137 datasets comprising long time series or covering a broad spatial extent are often important for
138 establishing temporal and spatial dynamics of change (e.g., population declines, range shifts,
139 etc.). The age of a dataset may be relevant, as older datasets can establish important baselines for
140 a species or system; the value of such datasets increases with time. The subject of the data is also
141 critical, as the societal value of the data may be higher when it involves species or ecosystems
142 with conservation, cultural, or economic value (e.g., datasets pertaining to species at risk have
143 higher conservation value). Additional considerations are the rarity of the data (e.g., data from an
144 under-sampled region or ecosystem), their uniqueness or irreplaceability (e.g., data from a
145 historical event, such as a natural disaster), and the potential costs of recollecting the data, if this
146 is possible (e.g., costs of re-running major experiments or extensive surveys). Finally, a key test
147 for the importance of a dataset is how it might be re-used in the future, with the most important
148 datasets having many immediate potential use scenarios. This is, perhaps, the most difficult (and
149 subjective) factor to assess.

150 The risks of data loss are similarly multifold. Data can be physically lost, and this risk is
151 highest for datasets for which there is only one copy (paper or digital). Data can also be
152 functionally lost when the datasets are unreadable because they are in older or defunct file
153 formats (e.g., Lotus 1-2-3) or in obsolete storage media (e.g., floppy disks). Data can also be
154 functionally lost when vital knowledge about collection or meaning of the data is lost (e.g.,
155 because the collector/creator of the data is deceased, retired, or otherwise unreachable).
156 Ultimately, successfully balancing the value of the data with the risk of its loss is essential for
157 effective prioritization of data rescue efforts.

158 2. Team creation

159 Data rescue takes a team, with different roles needed at different points in the rescue
160 process. We first consider those currently in possession of the data, which may include data
161 creators, data collectors, and data stewards: *data creators* are typically involved in generating the
162 ideas that lead to the data's collection and retain the intellectual property rights and
163 responsibilities for the data, even if not directly involved in collecting or managing the data
164 products; *data collectors* generate or collect the original data and, therefore, provide valuable
165 input for documenting the data (see *Metadata creation*); and *data stewards* are responsible for
166 managing and maintaining the data (i.e., organizing and keeping data safely archived, including
167 instances where researchers have been bequeathed data or organizations that act as custodians of
168 data collected by past employees). In ecology and evolution, these roles are often played by the
169 same person, though not always. For example, in a mentee-mentor relationship such as that
170 between a graduate student and supervisor, the student may play all three roles as data creator,
171 collector, and (temporary) steward, while the advisor may retain the data long-term as the
172 principal investigator, thereby acting as data creator and (long-term) steward. Having at least one
173 person who is a data creator, collector, or steward, if not more, as part of the data rescue team is
174 imperative for a successful data rescue mission.

175 A *data management expert* is another key role in the data rescue process. Usually, a data
176 manager is the one that plans the data lifecycle, but in a data rescue project, this role is mainly
177 focused on organizing and documenting the digitized datasets. This person will have the skills to
178 connect different datasets, clean and manage data, and compile previously unwritten information
179 in detailed metadata files. Additionally, if there are any data that have not been entered into a
180 digital format, a *data entry technician* will be an integral part of the team, ensuring that all

181 necessary data have been digitized in the appropriate format and validated against the original
182 records.

183 3. Metadata creation

184 *Metadata* are information about the data, typically contained in a file separate from the
185 dataset (Michener et al., 1997). The metadata generally describe the data collection process
186 (including the types of data collected, methodology, and contributors, among other information),
187 a description of all the variables in the dataset (e.g., column headings for tabular data),
188 abbreviations, units of measurement, and other relevant information necessary to understanding
189 how the data were generated and how to (re)use them (e.g., why some measurements are lacking;
190 British Ecological Society, 2018). We recommend early creation of the metadata, as this will
191 often inform the rest of the data rescue process and, ultimately, the structure of the compiled
192 data.

193 For datasets with more than one associated file, the metadata should also include a
194 description of the database structure, which data are contained in each file, and how files or
195 tables relate to each other. For datasets which include ongoing data collection, detailed metadata
196 files are important to ensure that subsequent data added to the database conform to the
197 appropriate standards and match the existing structure (Yenni et al., 2019). The metadata will
198 likely need to be revised after *Data compilation* (Step 5) and before *Data archiving* (Step 6) to
199 incorporate details about the data rescue process (e.g., data manipulation, validation, or changes
200 to the structure of the dataset or database; see below; Fig. 1).

201 The metadata file format varies (often dependent on the type of data or chosen
202 repository). Metadata are often found in a “README” style text file. Another useful format is a
203 text file written in Extensible Markup Language (XML; some examples and basics of XML can

204 be found at <https://www.xmlfiles.com/xml/>). Tools like XML have been developed for the
205 express purpose of writing and storing metadata and other information in a format that is both
206 human *and* machine readable, which not only ensures that prospective end users understand the
207 data structure and how it was created but also facilitates use by other software/programming
208 tools (e.g., search engines) that may rely on metadata being available in a standardized form.
209 Each variable is stored as “tags,” and its description is stored between tags. In ecology, there is a
210 set of suggested tags that should be used in such files, forming a variation of XML called
211 Ecological Metadata Language (EML; Fegraus et al., 2005; Jones et al., 2019; see
212 <https://eml.ecoinformatics.org/>).

213 4. Data transfer and compilation

214 For the data rescue team to work most effectively, all team members should have access
215 to the data and metadata files. However, this might only be possible if all files are already in a
216 digital format; if there are physical copies, they should either be photographed or scanned first or
217 entrusted to the team member responsible for data entry and validation. From there, discussion
218 about how the data should be compiled most effectively can ensue. While the details of data
219 compilation will need to be tailored to each dataset, the workflow should be as reproducible as
220 possible. At a minimum, all major decisions should be documented in the metadata. For
221 example, any edits made to the data should be done in a file separate from the original; a digital
222 file with the untouched original data should always remain.

223 In structuring the data, we generally suggest following Wickham’s (2014) “tidy data”
224 principles, which consist of 3 main concepts: (1) each variable has its own column, (2) each
225 observation has its own row, and (3) each type of observational unit is in its own data table, (e.g.,
226 individual-level measurements from a population, such as mass, in one table and population-

227 level metrics, such as abundance, in another). If there are multiple data tables, they should be
228 connected to each other by one or more variables that uniquely identify individual observations
229 (i.e., a primary key in a relational database design; Codd, 1990). While we advocate for tidy data
230 principles, as they are most likely to generate a data structure that will be useful in subsequent
231 analyses, sometimes other formats may be more relevant or efficient (e.g., a species by site
232 matrix).

233 5. Data cleaning and validation

234 Following data entry and compilation, data cleaning can be one of the most time-
235 intensive steps of the data management process. Data cleaning refers specifically to the process
236 of identifying and fixing issues in the dataset, such as data entry errors or incomplete records.
237 The importance of thorough and accurate data cleaning should not be overlooked, since, as the
238 adage “*garbage in, garbage out*” suggests, the inference drawn from an analysis is only as strong
239 as the inputs. In addition to common steps like correcting typographical or data entry errors, data
240 cleaning commonly includes checking for data completeness (i.e., that the data from all records
241 are fully and correctly transcribed) and uniformity (i.e., that variables are recorded in a consistent
242 way for all records, ensuring common measurement units, etc.), and otherwise ensuring the data
243 conform to expected standards. For environmental data, other common data cleaning steps
244 include checking for common date formats (e.g., the International Organization for
245 Standardization (ISO) 8601 standard recommends date-time objects be recorded as YYYY-MM-
246 DD hh:mm:ss + UTC offset), ensuring geographic coordinates are complete and standardized
247 (e.g., ISO 6709 applies to the representation of spatial information), and correcting misspellings
248 or synonyms in taxonomic information. Many tools have been developed to help with specific

249 aspects of data cleaning (e.g., the *taxize* package in R can be used to check and correct
250 taxonomies; Chamberlin & Szocs, 2013).

251 Related to data cleaning, data validation involves the comparison of the dataset against a
252 set of assertions determined *a priori* (e.g., dry body mass of an organism should be less than its
253 wet mass) or *post hoc* (e.g., the ratio of dry to wet mass should be similar among replicates).
254 Data validation is important for ensuring data quality and integrity by evaluating the data against
255 a set of expectations to confirm the structure and content of the data are appropriate. In the case
256 of data rescue projects, unlike most recently or currently collected data, data validation may
257 come with the extra challenge that the original data creator or collector may be unreachable or
258 deceased. As such, having as many original members of the data team (Fig. 1, Step 2; see *Team*
259 *creation*) is particularly beneficial for effective data validation. Common data validation
260 techniques include plotting the data in various ways to assist with identifying incorrect or
261 improbable values, checking that the contents (e.g., number of unique values in a column) or
262 dimensions of the data are in line with expectations following data manipulation, cross-checking
263 data from different columns or tables for mutual compatibility (i.e., to ensure that combinations
264 of data are within the realm of possibility), and evaluating summary statistics or other outputs
265 that characterize the data. In addition, many tools exist to help with the data validation process,
266 including open-source, “point-and-click” software (e.g., OpenRefine) as well as a number of
267 programming tools (e.g., the *assertr* and *validate* packages in R; Fischetti, 2020; van der Loo &
268 de Jonge, 2021).

269 Although the exact implementation of data cleaning and validation steps will vary
270 depending on the nature of the dataset, many of the same general principles described in the
271 *Data transfer and compilation* section are also relevant here. Validation should be conducted in

272 as reproducible a way as possible (e.g., in a script file that can be run on the original or cleaned
273 data files), and any errors identified during the validation process should be corrected without
274 manipulating the original (raw) data files. Importantly, any changes made based on these checks
275 should be well documented (e.g., as comments in the script or as notes in the metadata), as
276 should the rationale behind the corrections. More generally, the metadata are a critical source of
277 information for understanding the provenance of the data—that is, documentation of where the
278 data came from, how they were collected or generated, and the steps taken to clean and compile
279 the final dataset. Hence, thorough documentation of data validation steps is a key component of
280 *Metadata creation* and open and reproducible data sharing.

281 6. Data archiving

282 Archiving data in non-proprietary formats is imperative for longevity and future
283 accessibility. Non-proprietary software or file formats are those which do not have a copyright or
284 trademark and are, therefore, part of the public domain. Using non-proprietary formats ensures
285 that anyone can access the data without needing a specific (and often expensive) software
286 program or in the event that the program becomes defunct. For example, tabular data should be
287 stored in comma-separated values (.csv) format or text files (.txt) rather than in proprietary
288 formats such as Microsoft Excel files (.xls or .xlsx).

289 There is a strong and growing movement to archive data on public (and open) data
290 repositories rather than, or in addition to, private or institutional systems (e.g., a lab hard drive).
291 Indeed, many governments and funding agencies have recently implemented new data
292 management protocols that either encourage or mandate the archiving—though not necessarily
293 sharing—of all data generated using their resources (see below; e.g., Canada’s [Tri-agency](#)
294 [Research Data Management Policy](#)). The benefits of public archiving are clear. With each year

295 that passes after a publication, data that have not been publicly archived are 17% less likely to be
296 recoverable (Vines et al., 2014; see also Tedersoo et al., 2021). As a result, we also consider
297 public archiving to be an essential part of data rescue, since private archiving does not mitigate
298 the possibility that the data will need to be “re-rescued” in the future. Once the data and metadata
299 are compiled and validated, they should be placed in a data repository to maintain the data in a
300 secure and retrievable format for the future. Importantly, the push for public archiving does not
301 contradict the need for privacy or sensitivity associated with some datasets; it is possible to
302 publicly archive data while maintaining restrictions on when and how the data are accessed (see
303 below). In general however, we suggest that most environmental data should be openly
304 accessible upon archiving; exceptions include, for example, data pertaining to threatened species
305 and considerations of Indigenous data sovereignty.

306 There are now a number of excellent data repositories from which to choose, with some
307 being very generalized (e.g., Dryad, Dataverse, Figshare, Zenodo) while others cater to specific
308 types of data (e.g., DataONE for environmental data, GenBank for genetic sequences, Global
309 Biodiversity Information Facility (GBIF) for biodiversity data). Data repositories tend to use a
310 distributed (i.e., decentralized) approach to storing the data and have contingency plans in place
311 to ensure the longevity of the archived datasets (see r3data.org for a comprehensive list). Which
312 repository to choose will also be influenced by whether the data will remain private or be made
313 openly—and publicly—available upon upload or sometime in the near future (i.e., following an
314 embargo period; Roche et al., 2014). Some repositories allow for the long-term preservation of
315 datasets regardless of whether they are made openly available (e.g., Dataverse); others require
316 that the data be open access if they are to be hosted by the repository (e.g., Dryad). Many
317 archives also offer an option to place an embargo, or delay, on the publication of data. Most data

318 repositories will establish a Digital Object Identifier (DOI), a unique identifier which will remain
319 constant for the lifetime of the object, even if the object or metadata change. If the data will be
320 openly available, we recommend explicitly stating the terms of use for said data, such as noting
321 that authors should be contacted if the data are to be included in a publication or adding a
322 copyright statement, such as those from Creative Commons (e.g., CC0, CC-BY, etc.).

323 7. Data sharing

324 A final step to the data rescue workflow is to ensure that the data meet the open science
325 standards and that their use can be tracked. Open science principles and values entail
326 transparency, participation, and accessibility (Bartling & Friesike, 2014). These values can be
327 addressed in different ways and, because of that, ensuring a dataset meets these standards can be
328 overwhelming for researchers who are not trained in data management. These values can be met
329 with a combination of actions, some of which are summarized in the FAIR and CARE principles;
330 the first focuses on how data can be made useful and the second on how we can promote justice
331 through responsibly sharing open data.

332 The **FAIR** principles aim to improve **F**indability, **A**ccessibility, **I**nteroperability and
333 **R**eusability of datasets (Wilkinson et al., 2016). Providing human- and machine-readable
334 metadata improves both the findability and accessibility of a dataset. Combined with proper
335 archiving and identification, strong metadata may also help with the automatic discoverability of
336 datasets. As mentioned in the *Data archiving* section, tagging a dataset with a DOI makes it
337 trackable and citable, which improves the reproducibility of analyses. Many online data
338 repositories provide DOIs for datasets, and they are crucial to connect the actual dataset to its
339 metadata (which will also be registered under the same DOI). A comprehensive metadata file
340 also allows interoperability, or the ability of the data to be combined with other datasets in

341 different ways and in different systems. Additionally, accessibility and reusability can be
342 achieved through licenses, which explicitly describe the usage and attribution rights of the data.

343 The **CARE** principles focus on datasets that used traditional knowledge or benefited
344 somehow from Indigenous lands, promoting transparency and participation of open data (Carroll
345 et al., 2020). They aim to address and encourage consideration of the **Collective** benefit for
346 **Indigenous Peoples**, **Authority** to control (recognizing Indigenous data sovereignty),
347 **Responsibility** to be respectful with Indigenous Peoples involved in the dataset collection, and
348 **Ethics** (by assuring participation of Indigenous Peoples in the assessment of benefits, harms and
349 usability of the data; Carroll et al. 2020). These principles are meant to begin addressing the
350 larger, complicated history of colonialism in ecology, evolution, and related disciplines. While
351 these guidelines were written with current and future data collection in mind, they are equally
352 applicable to and important for previously collected data, and we recommend that all researchers
353 who are rescuing datasets take these principles into consideration.

354 Ongoing data rescue initiatives

355 Data rescue is not a new concept (e.g., Hawkins et al., 2013; Specht et al., 2018), and a
356 number of examples have been noted in both the scientific and grey literature (e.g., Box 2;
357 Norton et al., 2000; Hawkins et al., 2013; Kelly et al., 2016; Specht et al., 2018; Knockaert et
358 al., 2019). That said, the approach lacks a formal definition and can be conflated with other
359 terms, such as data reuse. It also lacks general guidelines and best practices, of which we have
360 offered a brief overview in this paper.

361 Some data rescue efforts have embraced community science, using crowdsourcing
362 platforms such as Zooniverse to facilitate data rescue (e.g., [Unearthing Michigan Ecological](#)
363 [Data](#)). Additionally, there are currently a few organizations focused on preserving ecological

364 data at risk. For example, the Canadian Institute for Ecology and Evolution (CIEE/ICEE) started
365 the [Living Data Project \(LDP\)](#) in 2018 with a mission to rescue and breathe new life into
366 languishing ecological, evolutionary, and environmental datasets (Box 3). Other organizations
367 practicing organized data rescue include (but are not limited to) the [Atmospheric Circulation](#)
368 [Reconstructions over the Earth \(ACRE\)](#) and the [International Environmental Data Rescue](#)
369 [Organization \(IEDRO\)](#).

370 Conclusion

371 Ultimately, we hope to reach a point where data rescue is no longer needed. This requires
372 researchers, funding agencies, and publishers to align their views around ethical and professional
373 obligations to archive data and make them publicly accessible where appropriate. It also requires
374 a culture change that sees best practices in data managing, archiving, and publicly sharing data
375 become the default in publicly funded research. While there has been movement in this direction,
376 we are still far from the ideal. To achieve this goal, data sharing and accessibility need to be
377 prioritized as a critical component of the scientific enterprise. We believe that the solution to
378 shifting the culture around data sharing is two-fold. First, there must be continued, long-term
379 investment in data management (Mons, 2020; Ritchie, 2021). Such investment includes not only
380 infrastructure but also training and support for students and personnel (Soeharjono & Roche,
381 2021). Additionally, publishers, employers and funding agencies must require some level of
382 accountability from researchers to preserve data in accessible, non-proprietary formats and, if
383 appropriate, make those data openly available to anyone interested (Mons, 2020). Until these
384 large, institutional-level paradigm shifts occur, however, smaller-scale and innovative data
385 rescue is an integral part of environmental data curation.

386 Currently, training in data management and shifting regulations regarding data
387 availability have, rightfully, focused on present and future data and data practices. With such a
388 strong eye to the future, however, much of the data of the past is being left behind. Data rescue
389 presents an opportunity to mitigate this loss of past data while also providing additional, less
390 tangible benefits. In the LDP, our mission of breathing life into languishing data is concomitant
391 with training the next generations of scientists in data management best practices and forging
392 connections amongst researchers across a wide variety of career stages and trajectories, thus
393 ensuring the longevity of scientific knowledge and preparing students for a data-rich future.

394 Acknowledgements

395 The Living Data Project (LDP) is a collaborative initiative by researchers at multiple
396 institutions across Canada: University of British Columbia-Vancouver, University of British
397 Columbia-Okanagan, University of Regina, McGill University, and Université de Montréal. The
398 authors recognize that we live and work on the traditional, ancestral, treaty, and unceded
399 territories of many Indigenous peoples, including the Coast Salish Peoples, xwməθkwəy̓əm
400 (Musqueam), Syilx (Okanagan), nêhiyawak (Cree), anihšīnāpēk (Saulteaux), Dakota, Lakota,
401 Nakoda, Attawanderon, Mississaugas, kanien'kehà:ka (Mohawk), and Haudenosaunee, and the
402 homeland of the Métis/Michif Nation.

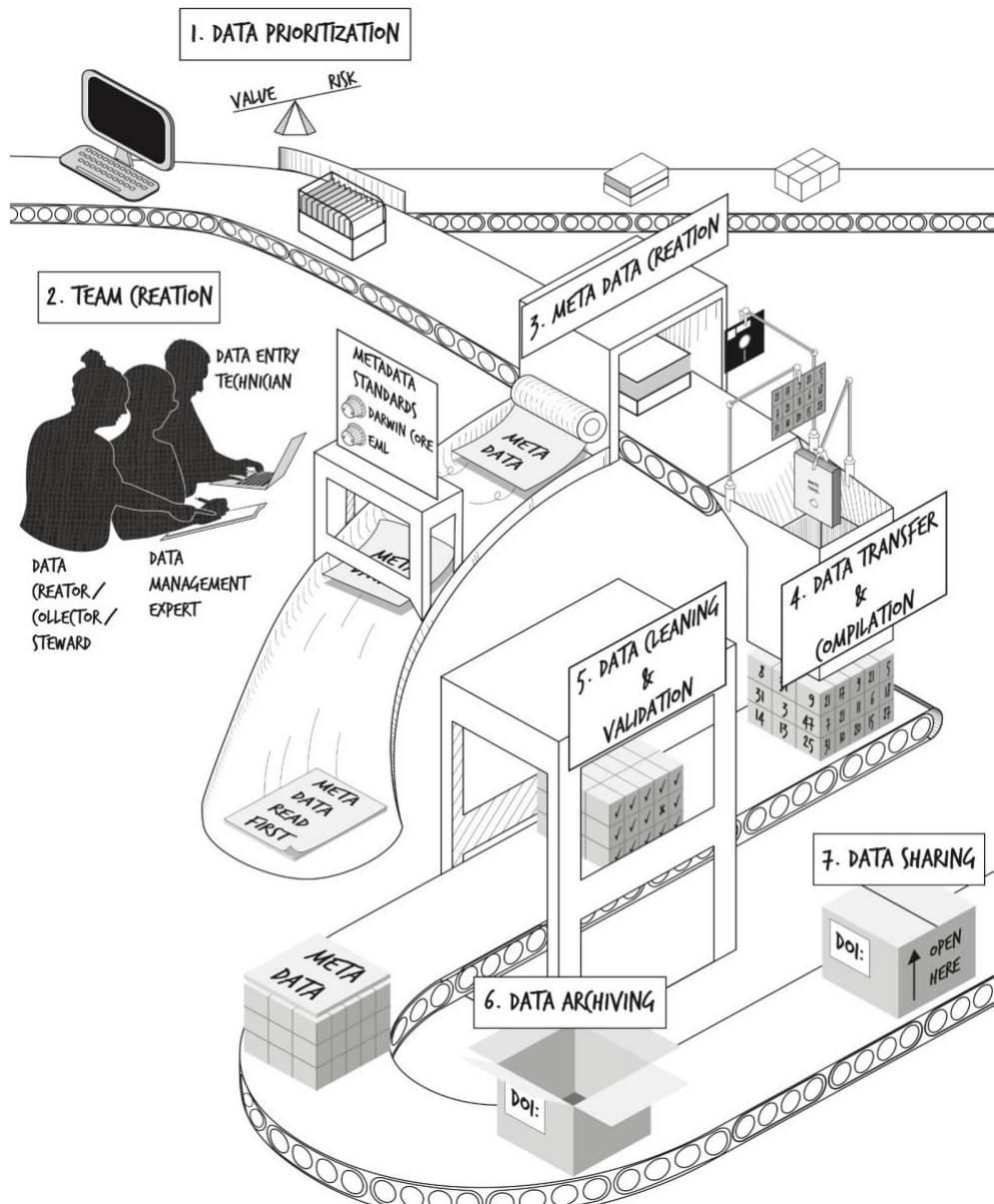
403 We would like to thank LDP members and partner organizations for their thoughtful
404 discussions on data rescue, best practices in data management, and fostering a more open and
405 transparent science. We are grateful to all of our project partners and the data rescue interns, and
406 in particular to those whose work we reference in this manuscript: Jenna Loesberg and Amelia
407 Hesketh, who worked in partnership with Drs. Ellen Macdonald and Justine Karst at the

408 University of Alberta, and Andrea Brown, who worked with Dr. Harold Eyster at the University
409 of British Columbia.

410 Funding for the LDP is provided through a Collaborative Research and Training
411 Experience (CREATE) grant to the Canadian Institute of Ecology and Evolution from the
412 Natural Sciences and Engineering Research Council of Canada. EKB, JBB, and GTH were
413 supported by the CREATE grant; EKB was also funded in part by the University of Regina.
414 DGR was supported by the European Union's Horizon 2020 research and innovation program,
415 under the Marie Skłodowska-Curie grant agreement (No. 838237-OPTIMISE).

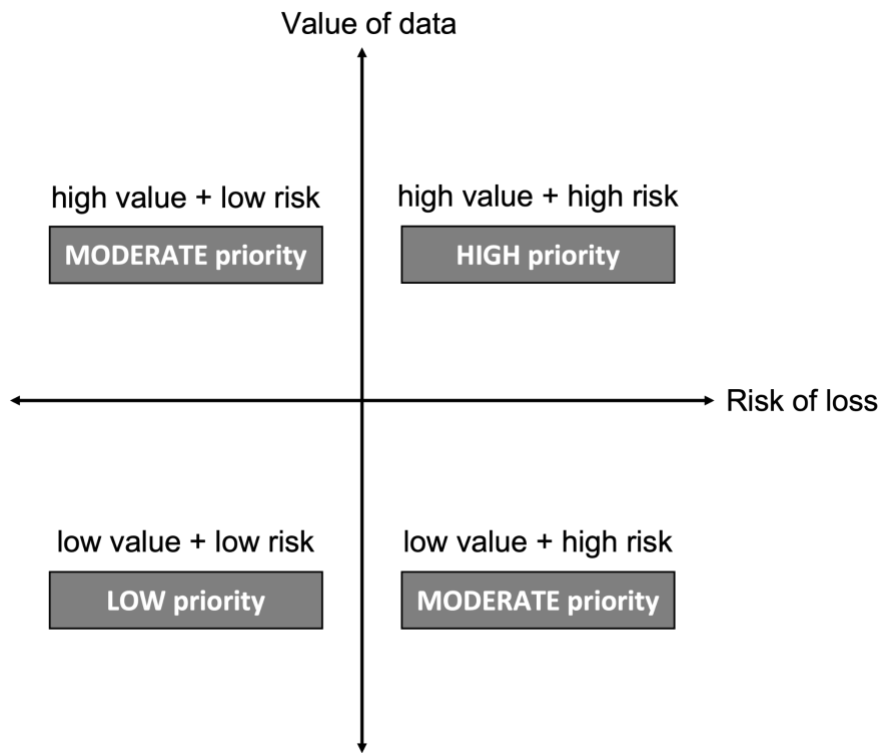
416

417 **Figure 1.** Steps in the data rescue process.



418
419 Figure 1. Steps in the data rescue process. First, data must be prioritized for rescue (Step 1).
420 After team creation (Step 2) and metadata creation (Step 3), the data must be transferred and
421 compiled into a consistent and effective format (Step 4). After data cleaning and validation (Step
422 5) is complete, the finalized data and metadata should be archived on a long-term data repository
423 (Step 6). The ultimate goal is to have the rescued data openly available for reuse (Step 7).

424 **Figure 2.** Prioritizing data for rescue: balancing the value of the
425 data and its risk of loss.



426
427 Figure 2. Prioritizing data for rescue: balancing the value of the data and its risk of loss. With
428 many datasets in need of preservation and limited resources, the first step in the data rescue
429 process requires developing a list of priorities for consideration and identifying relevant datasets
430 (Fig. 1). We consider data prioritization to be a balance between the assessed value of a dataset
431 in question and the potential risk of its loss in the absence of intervention (see *Data prioritization*
432 under *Guidelines*). While the process of prioritization is inherently subjective, we suggest that
433 considerations of value and risk can provide a useful heuristic for practitioners looking to best
434 target their time and effort.

435 **Box 1.** Spilt oil, spent money, and lost data: Exxon-Valdez oil spill as a
436 case study on the costs of data loss.

437 In 1989, the oil tanker *Exxon Valdez* struck the Bligh Reef in Prince William Sound, less
438 than 2.5 km from the Alaskan shore. As a result, approximately 37,000 tonnes of crude oil
439 spilled into the sound, leading to catastrophic short- and long-term ecological consequences. The
440 *Exxon Valdez* Oil Spill Trustee Council (EVOSTC) was established in 1991 to oversee the
441 spending of funds from a civil settlement in 1991 between *Exxon*, the United States federal
442 government and the state government of Alaska. A large portion of the funds were directed
443 towards determining and monitoring the impacts of the oil spill on oceanographic,
444 environmental, and ecological conditions. Prior to 2003, there was no requirement for data
445 preservation or availability; afterwards, all projects were awarded under explicit conditions from
446 EVOSTC that data be preserved and made publicly available (Jones et al., 2018). In their annual
447 report from 2010, the EVOSTC notes that the amount of funds spent on “Research, Monitoring,
448 and General Restoration” during 1992-2010 fiscal years was \$151.2 million USD (EVOSTC,
449 2012). The majority of funding went to state and federal agencies, though a few projects were
450 awarded to universities, professional societies, consultants, and other private entities (EVOSTC,
451 2018).

452 From 2012-2014, a group of researchers from the National Center for Ecological
453 Analysis & Synthesis (NCEAS) worked to recover the historical datasets funded by EVOSTC,
454 focusing specifically on data collected between 1989-2010 (Jones et al., 2018). Of the 419
455 projects determined to have been funded by EVOSTC during this time, only 27% of the datasets
456 were able to be recovered; after a total of 5 years hunting down datasets, this grew to 30% (Jones

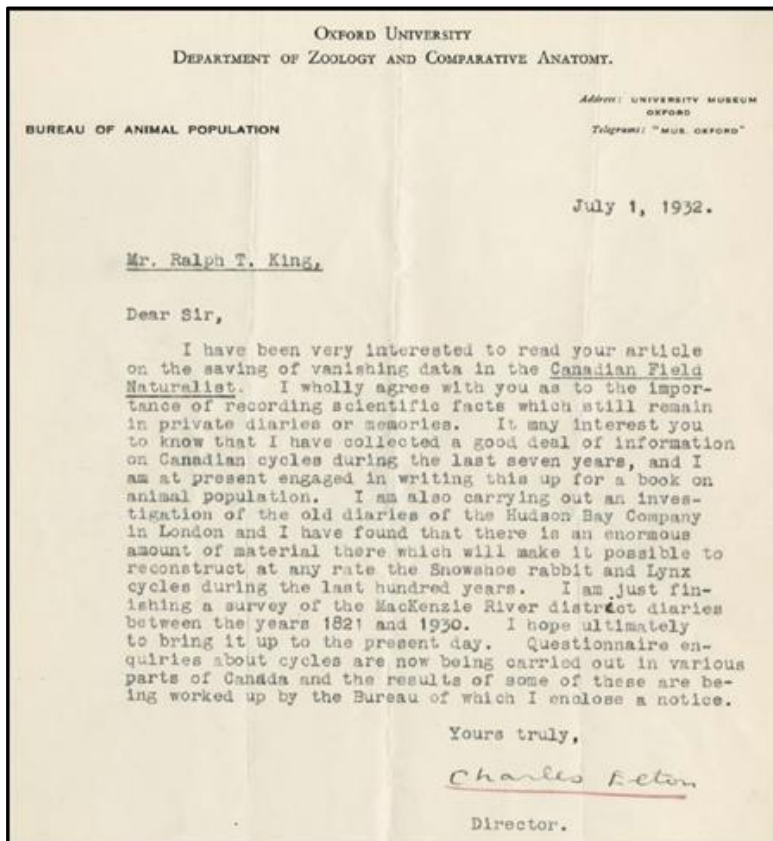
457 et al., 2018). Using these numbers, we can calculate a rough estimate of money spent on research
458 for which the data are not recoverable (70% of datasets): approximately **\$105 million USD was**
459 **spent collecting data which are no longer recoverable and, therefore, effectively lost to**
460 **science.** While we do not know the distribution of years from which data were recovered or how
461 money was divided by year, \$105 million USD is likely a conservative estimate, given that the
462 original cost does not include the first 3 years following the spill, when extensive ecological
463 assessments would have been completed. Similarly, this valuation does not include any of the
464 nearly \$50 million USD spent towards “Scientific Management, Public Information &
465 Administration” (EVOSTC, 2012).

466 The group tasked with recovering these historic datasets also noted the reasons for their
467 inability to recover the data. Instances in which data collectors specifically stated that the data
468 were lost or unrecoverable were rare (Jones et al., 2018). Instead, over 80% of datasets which
469 were unrecovered were lost due to a lack or failure of communication (~50% categorized as
470 “communication lost”); the authors of the final report, however, interpret much of this lack of
471 communication as an unwillingness or inability by the data owners to share data (Jones et al.,
472 2018), highlighting the importance of proper documentation and putting datasets in publicly
473 available data repositories for longevity.

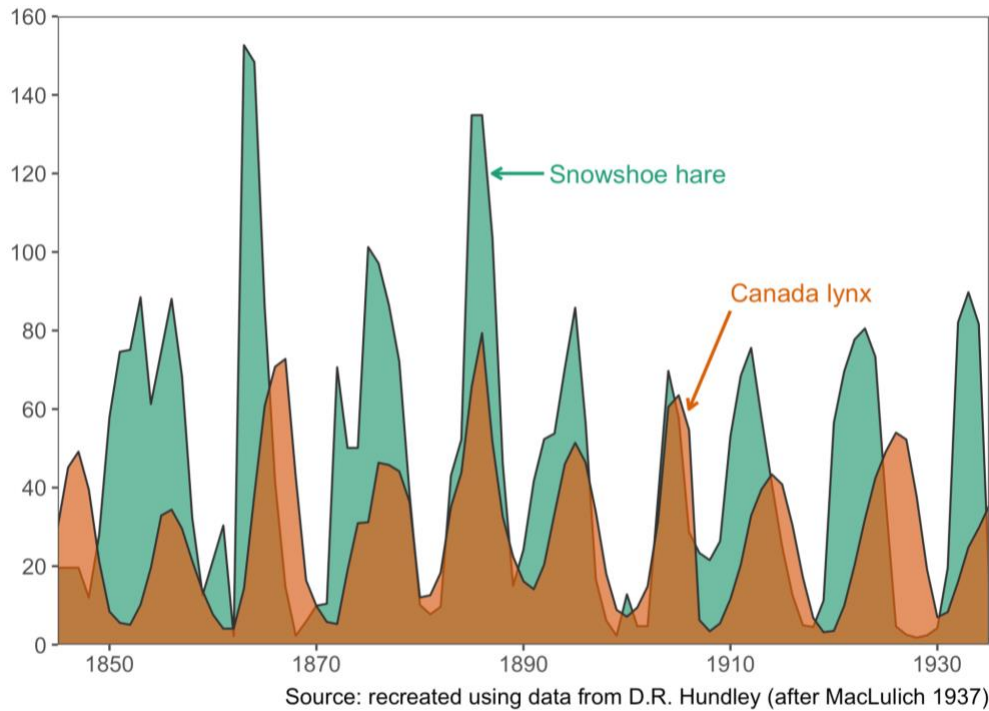
474 **Box 2.** From fur trappers to fundamental ecological theory: how a data
475 rescue effort shaped our understanding of population cycling.

476 Charles Sutherland Elton (29 May 1900 – 01 May 1991) was a British ecologist, whose major
477 contributions included work on population cycling and community dynamics. In 1932, Elton
478 established and became the first director of the Bureau of Animal Population (BAP) at Oxford
479 University and the inaugural editor of the *Journal of Animal Ecology* later the same year. As part
480 of their work on population cycling, Elton and his colleague Mary Nicholson endeavoured to
481 recover historical records on the number of Canada lynx (*Lynx canadensis*) furs collected by
482 Hudson’s Bay Company (HBC) trappers in Canada (Box 2.S1; Elton & Nicholson, 1942). In an
483 effort that spanned more than 15 years, Elton and Nicholson used these and other records to
484 collate information on trapping activities across the whole of Canada from 1886 to 1940 and as
485 far back as 1821 for areas in the Mackenzie River District, Northwest Territories. Much of this
486 work was akin to data rescue, including correspondence between Elton and the original data
487 owners (see Elton’s description of the process in a letter to Ralph King in Box. 2.1; Elton &
488 Nicholson, 1942), collation of data from different sources, as well as data cleaning and
489 validation. This work was not only central to compiling one of the longest time series of animal
490 populations and revealing the now classical example of ~12-year population cycles in snowshoe
491 hare and Canada lynx abundances (Box 2.2) but has spurred an entire field of ecology
492 (population/community cycling) and many decades of ecological research in the Canadian
493 Arctic. This is just one, elegant example of the immense value of historical data—even those
494 from unconventional places (like the ledgers of a colonial fur-trading company)—and the
495 importance of working to identify and preserve them.

496 **Box 2.1.** Letter from Charles S. Elton (Bureau of Animal Population, Oxford University) to
497 Ralph T. King (SUNY College of Environmental Science and Forestry), dated 01 July 1932.
498 In this letter, Elton expresses his interest in King's recent article on "saving vanishing data"
499 (King, 1932), which regarded many aspects of what we are calling "data rescue" and was itself
500 based on a paper of the same name written some three decades earlier (Haddon, 1903). Elton
501 goes on to describe his efforts to reconstruct time series of hares and lynx from HBC records.
502 This letter is not only an important historical artifact, but also highlights a "tradition" of data
503 rescue that dates to the formalization of ecology as a discipline. This letter was provided to us
504 courtesy of Dr. Adam T. Ford and is available through the Elton Archive at Oxford (Elton,
505 1932). A transcription of the letter's text is available in the Supporting Information (Box 2.S2).
506



507 **Box 2.2.** Time series of the numbers (in thousands) of Canada lynx and snowshoe hare pelts
508 provided to the Hudson's Bay Company.



509

510

511 **Box 3.** Recent data rescue examples from the Living Data Project.

512 As part of its core mission to contribute to and preserve ecological knowledge, the Living
513 Data Project (LDP) aims to rescue valuable ecological and environmental data at risk of being
514 lost. To achieve this objective, the LDP provides training opportunities for graduate students at
515 Canadian universities, including courses on topics and skills related to data rescue (data
516 management, reproducibility, and collaboration), and opportunities to put these skills into
517 practice through paid, short-term internships. The LDP partners with a variety of external
518 organizations, including government agencies, universities, and non-profits. These partners
519 propose potential data rescue projects, which are prioritized by a selection committee and
520 matched to graduate student interns with the relevant skills specific to each project (e.g., with
521 considerations for coding, database design, geospatial software, and language skills). Interns
522 work as part of a team comprised of representatives from the partner organization as well as
523 postdoctoral and faculty mentors from the LDP. Below we describe two recent data rescue
524 projects completed by LDP interns.

525 *Seeing the Forest Data for the Trees*

526 As researchers retire, they often think about the legacies they leave behind. Frequently,
527 however, curating the data they have collected in order to cement their legacies is not at the
528 forefront of their minds. Upon the retirement or death of a professor, students or colleagues often
529 must take the reins and piece together documents and data from decades-old research projects to
530 ensure the data's own legacy.

531 Dr. George H. La Roi was a professor of forest ecology at the University of Alberta for
532 35 years. In 2016, he composed an email to colleagues asking for help archiving his extensive

533 long-term survey data from the boreal forests in Alberta. Before this could be accomplished,
534 however, Dr. La Roi passed away in 2018. Upon his passing, La Roi's children bequeathed much
535 of his legacy of highly valuable data to his former colleague, Dr. Ellen Macdonald, who had
536 earlier taken over sampling some of his long-term plots. With no living data creator and with
537 much of the data in unorganized boxes containing unsorted datasheets, various documents, CD-
538 ROMs, and picture slides, the data was at high risk of being lost. Macdonald knew she would not
539 be able to tackle the boxes of materials on her own and joined forces with another University of
540 Alberta colleague, Dr. Justine Karst, who had also come into possession of some of La Roi's
541 boxes of data by way of University of Alberta's Botanic Garden. Together, they wrote an
542 application for an LDP data rescue internship. With the data being highly valuable long-term
543 data and also at a high risk of loss, this dataset was deemed to be of high priority for rescue.

544 Over the course of two data rescue internships, graduate students Jenna Loesberg and
545 Amelia Hesketh, along with a handful of undergraduate data entry technicians, sorted, entered,
546 and digitized the data. They determined that there were data from two different locations—the
547 Hondo-Slave Lake region and the Athabasca Oil Sands region—both of which included data on
548 vascular plant cover, bryoid cover, and forest mensuration, among other datasets. Some data
549 were found only on printed-out scans of hand-written datasheets and needed to be entered into a
550 digital format. Other data, which had already been entered and digitized, were stored in hundreds
551 of text files which required extensive reformatting and cleaning before they could be compiled
552 into usable datasets. Metadata also needed to be written and consolidated into one document for
553 future reuse; while most of the data had clear documentation, some data were lost, as no
554 documentation about the meaning of the variable names or the values in the column could be
555 found or determined. With all of this work completed, the data and metadata of this rich and

556 expansive dataset will be archived and made publicly available through University of Alberta's
557 Dataverse repository and hopefully published as a data paper.

558

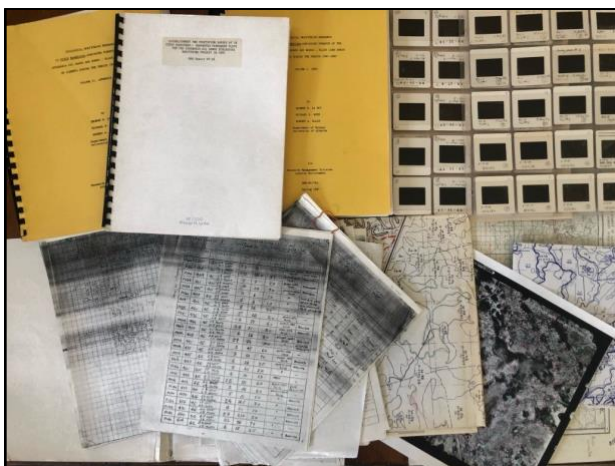
559 **Box 3.1.** *Photograph of researchers collecting data in the Athabasca Oil Sands region of*
560 *northern Alberta in 1982. This is one of 16 sites established by Dr. George La Roi in the region*
561 *in the 1980s to study seasonal and annual dynamics of boreal forests.* Image credit: unknown.



562

563

564 **Box 3.2.** *Photograph of loose data sheets, maps, reports, and picture slides; these items and*
565 *many more filled the boxes of research material left behind by Dr. George La Roi after his*
566 *passing in 2018.* Image credit: A. Hesketh.



567

568 *Out of the Archives and into the (Digital) Light of Day*

569 The archived theses and dissertations of former graduate students represent a rich, though
570 not fully realized, source of historical data. In particular, those prepared prior to the advent of
571 modern computer technologies and software, such as word processors and tools for statistical
572 analysis, contain troves of raw and summary data that have not been digitized and archived, and
573 so remain inaccessible to present-day researchers. As a result, the reuse of any raw or
574 summarized data from the thesis would first require data extraction and digitization.

575 Urban areas have expanded in size, number, and human population density in recent
576 decades, accompanied by changes in the abundance and diversity of bird populations that inhabit
577 these regions. Determining how biodiversity has changed in response to historical changes in
578 human activity and land use is central to understanding the impacts of these environmental
579 changes and predicting the potential for future declines. In a data rescue project proposed to us
580 by a then-doctoral student, Dr. Harold Eyster, LDP intern Andrea Brown worked to secure the
581 data contained in three University of British Columbia graduate theses (Weber, 1972; Lancaster,
582 1976; Melles, 2000), with a particular focus on data pertaining to surveys of bird abundances at
583 various locations around Greater Vancouver, British Columbia, Canada. While the specific
584 questions and research topics differed between these theses, the fact that all three surveyed the
585 same (or nearby) sites in Greater Vancouver over the span of several decades means that, in
586 combination, they present an opportunity to establish a baseline against which to compare
587 current and future trends (Box. 3.3 shows an example of the change in conditions at one of the
588 sites sampled by Weber (1972)). This project was identified as a priority for the LDP because the
589 data were both at-risk (much of the data existed only in non-digital formats and none of the

590 datasets are in active use) and high value (the data provide a valuable frame of reference for
591 studying changes in urban bird diversity).

592 During her internship, Brown first worked to transcribe the data from the earlier two of
593 the theses, Weber (1972) and Lancaster (1976), which were archived as scans of typewritten
594 documents and did not have data available in digital form. Among other challenges, digitization
595 required the conversion of non-standard data types (see, e.g., Box 3.4) into “tidy” forms that
596 could be used and interpreted programmatically. Data from the third thesis, Melles (2000), were
597 made available by the original author in a Microsoft Excel spreadsheet, and so only required
598 cleaning and manipulation, and conversion to a non-proprietary format. Later work included
599 efforts to rationalize the datasets so that they might be used in combination with each other (e.g.,
600 standardizing column names and other formatting, or combining similar or related tables into a
601 single file). Given the extensive data manipulation required, clear metadata were developed to
602 document the various steps taken to generate the final dataset and document other details from
603 the theses that were not captured during the digitization process. The data have been archived on
604 the UBC Dataverse repository (Brown, Eyster, & Lancaster, 2021; Brown, Eyster, & Melles,
605 2021; Brown, Eyster, & Weber 2021) and linked with the original theses.

606

607 **Box 3.3.** Comparison of the historical and current appearance of one of the sampling locations
 608 for urban bird surveys conducted in Vancouver, British Columbia, Canada. Photographs show
 609 the view looking west from the intersection of 24th Avenue West at Wallace Street (49.251°N,
 610 123.191°W). The historical reference is reproduced from Weber (1972); the contemporary image
 611 is shared with permission from the photographer.



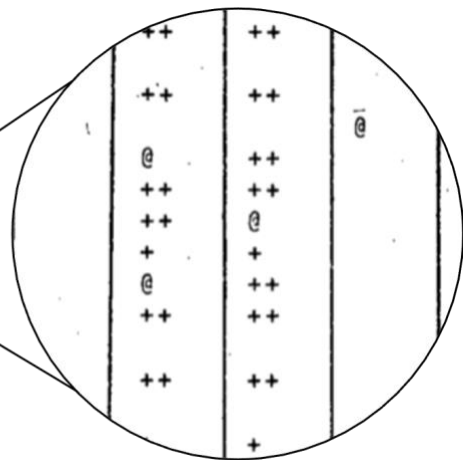
612
 613 April 1970 (Image credit: W.C. Weber) October 2021 (Image credit: © C.N. Nemeth)

614 **Box 3.4.** Example of non-standard (untidy) data to be rationalized and digitized. This example
 615 table contains symbolic data representing the significance of correlations between habitat
 616 features. These symbols were converted to numeric factors during digitization. Reproduced with
 617 modification from Lancaster (1976; see: Appendix 4, p. 103-104 therein).

Habitat Features	C>1.5	E<1.5	E>1.5	HERB	WEED	RHD1	RHD2	RHD3	HFD	TDD	TDC	TDE	FOOD	TOTVEG
SLANT														
FLAT														
ROAD														
LANE														
PVT+S														
LAWN														
D<7.5														
D>7.5														
C<7.5														
C>7.5														
BDEVG														
D<1.5														
D>1.5														
C<1.5														
C>1.5														
E<1.5														
E>1.5														
HERB														
WEED														
RHD1														
RHD2														
RHD3														
HFD														
TDD														
TDC														
TDE														
FOOD														
TOTVEG														

Explanation of Symbols (see also Table 1)
 + = Positive correlation *, - = significant at = .05
 - = Negative correlation **, -- = significant at = .01
 @ = correlation coefficient greater than .9800

618



619 References

- 620 Bartling, S., & Friesike, S. (2014). *Opening Science: The Evolving Guide on How the Internet is*
621 *Changing Research, Collaboration and Scholarly Publishing*. Springer Open. ISBN: 978-
622 2-319-00025-1.
- 623 Borgman, C. L. (2018). Open data, grey data, and stewardship: universities at the privacy
624 frontier. *Berkeley Technology Law Journal*, 33, 365-412.
625 <https://doi.org/10.15779/Z38B56D489>
- 626 British Ecological Society. (2018). A guide to data management in ecology and evolution. *BES*
627 *Guides to Better Science*. British Ecological Society, London, UK.
628 [https://www.britishecologicalsociety.org/wp-content/uploads/2019/06/BES-Guide-Data-](https://www.britishecologicalsociety.org/wp-content/uploads/2019/06/BES-Guide-Data-Management-2019.pdf)
629 [Management-2019.pdf](https://www.britishecologicalsociety.org/wp-content/uploads/2019/06/BES-Guide-Data-Management-2019.pdf)
- 630 Broman, K. W., & Woo, K. H. (2018). Data organization in spreadsheets. *American Statistician*,
631 72, 2-10. <https://doi.org/10.1080/00031305.2017.1375989>
- 632 Brown, A., Eyster, H., & Lancaster, R. K. (2021). Data for: Bird communities in relation to the
633 structure of urban habitats. *Scholars Portal Dataverse*.
634 <https://doi.org/10.5683/SP2/YD6N7C>
- 635 Brown, A., Eyster, H., & Melles, S. J. (2021). Data for: Effects of landscape and local habitat
636 features on bird communities: a study of an urban gradient in greater Vancouver.
637 *Scholars Portal Dataverse*. <https://doi.org/10.5683/SP2/BPLPAP>
- 638 Brown, A., Eyster, H., & Weber, W. C. (2021). Data for: Birds in cities: a study of populations,
639 foraging ecology and nest-sites of urban birds. *Scholars Dataverse Portal*.
640 <https://doi.org/10.5683/SP2/K5LMLA>

641 Buck, S. (2021). Beware performative reproducibility. *Nature*, 595, 151.
642 <https://doi.org/10.1038/d41586-021-01824-z>

643 Buxton, R. T., Nyboer, E. A., Pigeon, K. E., Raby, G. D., Rytwinski, T., Gallagher, A. J.,
644 Schuster, R., Lin, H.-Y., Fahrig, L., Bennett, J.R, Cooke, S.J., & Roche, D.G. (2021).
645 Avoiding wasted research resources in conservation science. *Conservation Science and*
646 *Practice*, 3(2), e329. <https://doi.org/10.1111/csp2.329>

647 Carroll, S. R., Garba, I., Figueroa-Rodríguez, O. L., Holbrook, J., Lovett, R., Materechera, S.,
648 Parsons, M., Raseroka, K., Rodriguez-Lonebear, D., Rowe, R., Sara, R., Walker, J. D.,
649 Anderson, J., & Hudson, M. (2020). The CARE principles for Indigenous data
650 governance. *Data Science Journal*, 19, 43. <https://doi.org/10.5334/dsj-2020-043>

651 Chamberlain, S., & Szocs, E. (2013). taxize – taxonomic search and retrieval in R. *F1000*
652 *Research*, 2, 191. <https://doi.org/10.12688/f1000research.2-191.v2>

653 Codd, E. F. (1990). *The Relational Model for Database Management: Version 2*. Addison-
654 Wesley Longman Publishing.

655 Cook, R. B., Olson, R. J., Kanciruk, P., & Hook, L. A. (2001). Best practices for preparing
656 ecological data sets to share and archive. *Bulletin of the Ecological Society of America*,
657 82, 138-141. <https://www.jstor.org/stable/20168543>

658 Eamon, W. (1985). From the secrets of nature to public knowledge: the origins of the concept of
659 openness in science. *Minerva*, 23, 321-347. <https://doi.org/10.1007/BF01096442>

660 Elton, C. S. (1932). Letter to Ralph T. King, 01 July. MS. Eng. c3328 A72, Elton Archives,
661 Weston Library, University of Oxford.

662 Elton, C. S., & Nicholson, M. (1942). The ten-year cycle in numbers of the lynx in Canada.
663 *Journal of Animal Ecology*, 11, 215-244. <https://www.jstor.org/stable/1358>

664 EVOSTC (*Exxon Valdez Oil Spill Trustee Council*) (2012). 2010 Annual Report.
665 <https://evostc.state.ak.us/media/4411/2010annualreport.pdf>
666 EVOSTC (*Exxon Valdez Oil Spill Trustee Council*) (2018). *Exxon Valdez Oil Spill Final*
667 *and Annual Reports*. <https://evostc.state.ak.us/media/4291/finalandannualreports.pdf>
668 Fegraus, E. H., Andelman, S., Jones, M. B., & Schildhauer, M. (2005). Maximizing the
669 Value of Ecological Data with Structured Metadata: An Introduction to Ecological
670 Metadata Language (EML) and Principles for Metadata Creation. *Bulletin of the*
671 *Ecological Society of America*. 86(3), 158-68.
672 <http://www.jstor.org/stable/bullecosociamer.86.3.158>.
673 Fischetti, T. (2020). *assertr: assertive programming for R analysis pipelines*. R package version
674 2.7. <https://CRAN.R-project.org/package=assertr>
675 Gatti, G., Bianchi, C. N., Parravicini, v., Rovere, A., Peirano, A., Montefalcone, M., Massa, F.,
676 & Morri, C. (2015). Ecological change, sliding baselines and the importance of historical
677 data: lessons from combining observational and quantitative data on a temperate reef over
678 70 years. *PLoS One*, 10, e0123268. <https://doi.org/10.1371/journal.pone.0118581>
679 Haddon, A. C. (1903). The saving of vanishing data. *Popular Science Monthly*, 63, 222-229.
680 [https://en.wikisource.org/wiki/Popular_Science_Monthly/Volume_62/January_1903/The](https://en.wikisource.org/wiki/Popular_Science_Monthly/Volume_62/January_1903/The_Saving_of_Vanishing_Data)
681 [_Saving_of_Vanishing_Data](https://en.wikisource.org/wiki/Popular_Science_Monthly/Volume_62/January_1903/The_Saving_of_Vanishing_Data)
682 Hawking, S. J., Firth, L. B., McHugh, M., Poloczanska, E. S., Herbert, R. J. H., Burrows, M. T.,
683 Kendall, M. A., Moore, P. J., Thompson, R. C., Jenkins, S. R., Sims, D. W., Genner, M.
684 J., & Mieszkowska, N. (2013). Data rescue and re-use: recycling old information to
685 inform new policy concerns. *Marine Policy*, 42, 91-98.
686 <https://doi.org/10.1016/j.marpol.2013.02.001>

687 Jones, M. B., Blake, R., Couture, J., & Ward, C. (2018). Collaborative data management and
688 holistic synthesis of impacts and recovery status associated with the *Exxon Valdez* oil
689 spill. *Exxon Valdez Oil Spill Long-Term Monitoring Program (Gulf Watch Alaska) Final*
690 *Report* (project 16120120). Exxon Valdez Oil Spill Trustee Council, Anchorage, Alaska.
691 [http://www.gulfwatchalaska.org/wp-content/uploads/2018/08/16120120-Jones-et-al.-](http://www.gulfwatchalaska.org/wp-content/uploads/2018/08/16120120-Jones-et-al.-2018-Final-Report.pdf)
692 [2018-Final-Report.pdf](http://www.gulfwatchalaska.org/wp-content/uploads/2018/08/16120120-Jones-et-al.-2018-Final-Report.pdf)

693 Jones, M. B., O'Brien, M., Mecum, B., Boettiger, C., Schildhauer, M., Maier, M., Whiteaker, T.,
694 Earl, S., & Chong, S. (2019). Ecological Metadata Language version 2.2.0. KNB Data
695 Repository. <https://doi.org/10.5063/F11834T2>

696 Kelly, G., Easterday, K., Rapucciuolo, G., Koo, M. S., McIntyre, P., & Thorne, J. (2016).
697 Rescuing and sharing historical vegetation data for ecological analysis: the California
698 Vegetation Type Mapping Project. *Biodiversity Informatics*, 11, 40-62.
699 <https://core.ac.uk/download/pdf/162636907.pdf>

700 King, R. T. (1932). The saving of vanishing data. *Canadian Field Naturalist*, 46, 108-111.
701 <https://www.biodiversitylibrary.org/ia/canadianfieldnat1932otta/#page/134/mode/1up>

702 Knockaert, C., Tyberghein, L., Goffin, A., Vanhaecke, D., Ong'anda, H., Wakwabi, E. O., &
703 Mees, J. (2019). Biodiversity data rescue in the framework of a long-term Kenya-
704 Belgium cooperation in marine sciences. *Scientific Data* 6(85).
705 <https://doi.org/10.1038/s41597-019-0092-8>

706 Lancaster, R.K. (1976). Bird communities in relation to the structure of urban habitats. Thesis.
707 Department of Zoology, University of British Columbia.
708 <https://dx.doi.org/10.14288/1.0093863>

709 Lennox, R.J., Harcourt, R., Bennett, J.R., Davies, A., Ford, A.T., Frey, R.M., ..., & Cooke, S. J.
710 (2020). A novel framework to protect animal data in a world of biosurveillance.
711 *BioScience*, 70, 468-476. <https://doi.org/10.1093/biosci/biaa035>

712 van der Loo, M.P.J., & de Jonge, E. (2021). Data validation infrastructure for R. *Journal of*
713 *Statistical Software*, 97, 1–31. <https://doi.org/10.18637/jss.v097.i10>

714 McClenachan, L., Ferretti, F., & Baum, J. K. (2012). From archives to conservation: why
715 historical data are needed to set baselines for marine animals and ecosystems.
716 *Conservation Letters*, 5, 349-359. <https://doi.org/10.1111/j.1755-263X.2012.00253.x>

717 Melles, S. J. (2000). Effects of landscape and local habitat features on bird communities: a study
718 of an urban gradient in Greater Vancouver. Thesis. Department of Forest Sciences,
719 University of British Columbia. <https://dx.doi.org/10.14288/1.0099590>

720 Michener, W. K., Brunt, J. W., Helly, J. J., Kirchner, T. B. & Stafford, S. G. (1997).
721 Nongeospatial metadata for the ecological sciences. *Ecological Applications*, 7, 330-342.

722 Mills, J. A., Teplitsky, C., Arroyo, B., Charmantier, A., Becker, P. H., Birkhead, T. R., Bize, P.,
723 Blumstein, D. T., Bonenfant, C., Boutin, S., Bushuev, A., Cam, E., Cockburn, A., Côté S.
724 D., Coulson, J. C., Daunt, F., Dingemans, N. J., Doligez, B., Drummond H., Espie, R. H.
725 M., et al. (2015). Archiving primary data: solutions for long-term studies. *Trends in*
726 *Ecology and Evolution*, 30, 581-589. <https://doi.org/10.1016/j.tree.2015.07.006>

727 Mons, B. (2020). Invest 5% of research funds in ensuring data are reusable. *Nature*, 578(7796),
728 491. <https://doi.org/10.1038/d41586-020-00505-7>

729 Norton, D. C., Assel, R. A., Meyers, D., Hibner, B. A., Morse, N., Trimble, P. J., Cronk, K., &
730 Rubens, M. (2000). Great Lakes ice data rescue project (Technical memorandum
731 GLERL-117). Great Lakes Environmental Research Laboratory, National Oceanographic

732 and Atmospheric Administration (NOAA).
733 <https://repository.library.noaa.gov/view/noaa/11024>

734 O’Dea, R. E., Parker, T. H., Chee, Y. E., Culina, A., Drobniak, S. M., Duncan, D. H., ... &
735 Nakagawa, S. (2021). Towards open, reliable, & transparent ecology and evolutionary
736 biology. *BMC Biology*, 19(1), 1-5. <https://doi.org/10.1186/s12915-021-01006-3>

737 Powers, S. M., & Hampton, S. E. (2019). Open science, reproducibility, and transparency in
738 ecology. *Ecological Applications*, 29, e01822. <https://doi.org/10.1002/eap.1822>

739 Ritchie, H. (2021). COVID’s lessons for climate, sustainability and more from Our World in
740 Data. *Nature*, 598:9. <https://doi.org/10.1038/d41586-021-02691-4>

741 Roche, D. G., Berberi, I., Dhane, F., Lauzon, F., Soeharjono, S., Dakin, R., & Binning, S. A.
742 (2021). The quality of open datasets shared by researchers in ecology and evolution is
743 moderately repeatable and slow to change. *EcoEvoRxiv*.
744 <https://doi.org/10.32942/osf.io/d63js>

745 Roche, D. G., Kruuk, L. E. B., Lanfear, R., & Binning, S. A. (2015). Public data archiving in
746 ecology and evolution: how well are we doing? *PLoS Biology*, 13, e1002295.
747 <https://doi.org/10.1371/journal.pbio.1002295>

748 Roche, D. G., Lanfear, R., Binning, S. A., Haff, T. M., Schwanz, L. E., Cain, K. E., Kokko, H.,
749 Jennions, M. D., & Kruuk, L. E. (2014). Troubleshooting public data archiving:
750 suggestions to increase participation. *PLoS Biology*, 12(1), e1001779.
751 <https://doi.org/10.1371/journal.pbio.1001779>

752 Smith, R., & Roberts, I. (2016). Time for sharing data to become routine: the seven excuses for
753 not doing so are all invalid. *F1000 Research*, 5, 781.
754 <https://doi.org/10.12688/f1000research.8422.1>

755 Specht, A., Bolton, M. P., Kingsford, B., Specht, R. L., & Belbin, L. (2018). A story of data won,
756 data lost and data re-found: the realities of ecological data preservation. *Biodiversity Data*
757 *Journal*, 6, e29073. <https://doi.org/10.3897/BDJ.6.e28073>

758 Stieglitz, S. Wilms, K., Mirbabaie, M., Hofeditz, L., Brenger, B., López, A., & Rehwald, S.
759 (2020). When are researchers willing to share their data? - Impacts of values and
760 uncertainty on open data in academia. *PLoS One*, 15, e0234172.
761 <https://doi.org/10.1371/journal.pone.0234172>

762 Soeharjono, S., & Roche, D. R. (2021). Reported individual costs and benefits of sharing open
763 data among Canadian academic faculty in ecology and evolution. *BioScience*, biab024.
764 <https://doi.org/10.1093/biosci/biab024>

765 Tedersoo, L., Küngas, R., Oras, E., Köster, K., Eenmaa, H., Leijen, Ä., Pedaste, M., Raju, M.,
766 Astapova, A., Lukner, H., Korgerman, K., & Sepp, T. (2021). Data sharing practices and
767 data availability upon request differ across scientific disciplines. *Scientific Data*, 8, 192.
768 <https://doi.org/10.1093/biosci/biab024>

769 Vines, T. H., Albert, A. Y. K., Andrew, R. L., Débarre, F., Bock, D. G., Franklin, M. T., Gilbert,
770 K. J., Moore, J., Renault, S., & Rennison, D. J. (2014). The availability of research data
771 declines rapidly with article age. *Current Biology*, 24, 94-97.
772 <https://doi.org/10.1016/j.cub.2013.11.014>

773 Walter, M., & Suina, M. (2018). Indigenous data, indigenous methodologies and indigenous data
774 sovereignty. *International Journal of Social Research Methodology*, 22, 233-243.
775 <https://doi.org/10.1080/13645579.2018.1531228>

776 Weber, W. C. (1972). Birds in cities: a study of populations, foraging ecology and nest-sites of
777 urban birds. Thesis. Department of Zoology, University of British Columbia.
778 <https://dx.doi.org/10.14288/1.0101293>

779 Wickham, H. (2014). Tidy Data. *Journal of Statistical Software*, 59(10).
780 <http://dx.doi.org/10.18637/jss.v059.i10>

781 Wilkinson, M. D., Dumontier, M., Aalbersberg, IJ. J., Appleton, G., Axton, M., Baak, A.,
782 Blomberg, N., Boiten, J., Bonino de Silva Santos, L., Bourne, P. E., Bouwman, J.,
783 Brooke, A. J., Clark, T., Crosas, M., Dillo, I., Dumon, O., Edmunds, S., Evelo, C. T.,
784 Finkers, R., Gonzalez-Beltran, A., et al. (2016). The FAIR Guiding Principles for
785 scientific data management and stewardship. *Scientific Data*, 3, 160018.
786 <https://doi.org/10.1038/sdata.2016.18>

787 Willis, K. J., Araùjo, M. B., Bennett, K. D., Figueroa-Rangel, B., Freud, C. A., & Myers, N.
788 (2007). How can a knowledge of the past help to conserve the future? Biodiversity
789 conservation and the relevance of long-term ecological data. *Philosophical Transactions*
790 *of the Royal Society B*, 362, 175-187. <https://doi.org/10.1098/rstb.2006.1977>

791 White, E. P., Baldrige, E., Brym, Z. T., Locey, K. J., McGlenn, D. J. & Supp, S. R. (2013). Nine
792 simple ways to make it easier to (re) use your data. *Ideas in Ecology and Evolution*, 6, 1–
793 10. <https://doi.org/10.4033/iee.2013.6b.6.f>

794 Whitlock, M. C. (2011). Data archiving in ecology and evolution: best practices. *Trends in*
795 *Ecology and Evolution*, 26, 61-65. <https://doi.org/10.1016/j.tree.2010.11.006>

796 Yenni, G. M., Christensen, E. M., Bledsoe, E. K., Supp, S. R., Diaz, R. M., White, E. P., &
797 Ernest, S. M. (2019). Developing a modern data workflow for regularly updated data.
798 *PLoS Biology*, 17(1), e3000125. <https://doi.org/10.1371/journal.pbio.3000125>