

1 Data rescue: saving environmental data from extinction

2 Ellen K. Bledsoe^{1,2,†,*} (0000-0002-3629-7235), Joseph B. Burant^{1,3,4,†,*} (0000-0002-0713-3100),
3 Gracielle T. Higinio^{1,5,†} (0000-0003-2791-8383), Dominique G. Roche^{1,6} (0000-0002-3326-
4 864X), Sandra A. Binning^{1,4} (0000-0002-2804-9979), Kerri Finlay^{1,2} (0000-0001-6835-8832),
5 Jason Pither^{1,7} (0000-0002-7490-6839), Laura S. Pollock^{1,3} (0000-0002-6004-4027), Jennifer M.
6 Sunday^{1,3} (0000-0001-9372-040X), Diane S. Srivastava^{1,6,*} (0000-0003-4541-5595)

7 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

† 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

11 Abstract

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

27 Keywords

28 Data archiving, historical data, long-term ecological research, open data, open science,
29 reproducibility

Total word count (excluding title page and supplementary materials) = 7,821

Abstract = 199

Main text (excluding abstract, figure captions, boxes, and references) = 4,037

Figures: 2

Figure 1 caption = 80

Figure 2 caption = 83

Boxes: 3

Box 1 = 458 words

Box 2 = 408 words, 2 images

Box 3 = 1,202 words, 4 images

References = 49 (1,354 words)

30 Why Rescue Data?

31 Data are among the most valuable outputs of research and scholarship; beyond helping
32 answer important questions, they inform new lines of inquiry, new testable hypotheses, and
33 future data collection efforts. Observational and experimental data derived from ecology,
34 evolution, conservation and environmental sciences (hereafter, environmental data) are essential
35 to establishing historical trajectories of ecosystems (“baselines”; McClenachan et al., 2012),
36 understanding how species and communities respond to environmental change (Gatti et al.,
37 2015), and designing and evaluating the outcomes of management efforts (Hawkin et al., 2013;
38 Willis et al., 2007). While data collection is often targeted to particular populations,
39 communities, or locations, the reuse (i.e., aggregation, collation, and synthesis) of data from
40 different contexts is essential to establishing broader ecological knowledge and informing
41 conservation management (Renaut et al., 2018). Yet, despite their high value, data are often
42 misplaced, filed away, or otherwise rendered unusable, often through poor data management
43 practices (Vines et al., 2014). In their unusable and “at-risk” state, these data represent an
44 egregious waste of resources expended on their collection (Buxton et al., 2021; Box 1).
45 Languishing data, however, also offer an enormous opportunity. **Data rescue**—defined here as
46 the identification, preservation, and sharing of valuable data and associated metadata at risk of
47 loss—has the potential to realize substantial benefits for society, especially considering the
48 crucial roles that baseline data play in informing management and policy decisions. The ultimate
49 goal of data rescue is to make previously inaccessible or poorly preserved data available for
50 (re)use, ideally through archiving them in a permanent, publicly accessible, and reusable format.

51 In recent years, there has been a strong push from within the scientific and scholarly
52 communities for increased openness in the practice of science, including in ecology and
53 evolution (e.g., O’Dea et al., 2021). Calls for more transparency and accessibility in science are
54 not new (e.g., Eamon, 1985); the last decade, however, has seen a surge in general awareness and
55 promotion of open science practices (e.g., open access publishing and open data, code, software,
56 and peer-review) and their benefits (Powers & Hampton, 2019). These initiatives have not been
57 without criticism, with many researchers unsure about sharing their data due to real or perceived
58 concerns about data misuse and loss of control (Roche et al., 2014; Smith & Roberts, 2016;
59 Stieglitz et al., 2020). Others have acknowledged important caveats to the general appeal for

60 openness (e.g., valid considerations about security, confidentiality, equity, and Indigenous data
61 sovereignty and governance; Borgman, 2018; Walter & Suina, 2018; Lennox et al., 2020; Buck,
62 2021). Despite the legitimacy of (some of) these concerns, the benefits of data sharing are
63 apparent (Powers & Hampton; 2019; Soeharjono & Roche, 2021). Even so, large amounts of
64 data remain private and unavailable for reuse by other scientists. For example, in a sample of
65 more than 4,000 ecology and evolution papers, only one in five papers (21.5%) had a data
66 availability statement or associated open data (Roche et al., 2021), and less than half of archived
67 datasets in ecology and evolution are reusable (Roche et al., 2015; Roche et al., 2021).

68 Open science initiatives have developed rapidly, and the last few years have seen a rise in
69 the number of institutions, governments, funding agencies, and publishers who have
70 implemented policies that require the open, permanent, and accessible sharing of data (e.g.,
71 FAIR data principles [see *Data sharing*; Wilkinson et al., 2016], the Ecological Society of
72 America's new [Open Research policy](#), and the European Commission's [OpenAIRE open access
73 and open data policy](#)). These requirements, and participation by scientists, will enhance our
74 ability to evaluate, reuse, and synthesize increasingly rich and complex ecological data.
75 However, open data policies are not retroactive and, therefore, do little to address issues of
76 access to and preservation of previously-collected data (Vines et al., 2014). Arguably, data
77 collected prior to the adoption of widespread sharing practices remain a public good, funded by
78 taxpayers and governments, so rescuing datasets to ensure their longevity and accessibility is
79 imperative.

80 Here, we present general guidelines for implementing data rescue, with a focus on
81 environmental data. These recommendations are based on past and ongoing data rescue projects
82 by the Living Data Project, an initiative of the Canadian Institute of Ecology and Evolution
83 (CIEE), which aims to identify and secure vulnerable datasets and bring new life to them through
84 collaborative analysis and synthesis. We include examples using historical (Box 2) and recent
85 data rescue efforts (Box 3). We anticipate these guidelines will (a) focus attention on the current
86 threats to the usability and integrity of previously-collected data, (b) stimulate broader
87 consideration of the utility of previously-collected datasets for current research efforts, (c)
88 encourage people with access to or knowledge of unarchived data to work towards their
89 preservation, (d) provide a reference for those looking to apply data rescue techniques in the

90 context of their own work, and (e) help foster a strong culture of data stewardship such that data
91 rescue becomes unnecessary in the future.

92 Guidelines for data rescue

93 Imperiled data can be found nearly everywhere, such as non-profit organizations,
94 conservation councils, academic institutions, and government agencies (think: historical data
95 only available on paper records in basement filing cabinets, digitized data stored only on floppy
96 disks, etc.). Finding data to rescue is usually the easy part; implementing a successful data rescue
97 mission, however, requires a more strategic approach (Fig. 1). Some steps involved in data
98 rescue are closely aligned with recommended practices in research data management (see
99 *Metadata, Data Compilation, Validation, Archiving and Sharing* sections). Several resources
100 have already outlined “best” practices for data collection (Broman & Woo, 2018), management
101 (e.g., BES, 2018), and archiving (Cook et al., 2001; Renault et al., 2018; Whitlock, 2011; White
102 et al., 2013), yet these are written with current or future data collection in mind and do not
103 address historically-collected or unmanaged data. Below, we outline seven key steps for data
104 rescue, from identifying high-priority datasets to archiving and sharing them for (re)use.

105 1. Data prioritization

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

112 High-value environmental datasets have some common features. Scale is a key factor, as
113 datasets comprising long time series or covering a broad spatial extent are often important for
114 establishing temporal and spatial dynamics of change (e.g., population declines, range shifts,
115 etc.). The age of a dataset may be relevant, as older datasets can establish important baselines for
116 a species or system, and the value of such datasets increases with time. The subject of the data is
117 also critical, as the societal value of the data may be higher when it involves species or

118 ecosystems with conservation, cultural, or economic value. Additional considerations are the
119 rarity of the data (e.g., data from an undersampled region or ecosystem), their uniqueness or
120 irreplaceability (e.g., data from a historical event, such as a natural disaster), and the potential
121 costs of recollecting the data, if that is a possibility. Finally, how the data might be re-used in the
122 future is important, with the most high-value datasets having many, immediate potential use
123 scenarios. This is, perhaps, the most difficult (and subjective) factor to assess.

124 The risks of data loss are similarly multifold. Data can be physically lost, and this risk is
125 highest for datasets for which there is only one copy (paper or digital). Data can also be
126 functionally lost when the datasets are unreadable because they are in older or defunct file
127 formats (e.g., Lotus 1-2-3™) or in obsolete storage media (e.g., floppy disks). Data can also be
128 functionally lost when vital knowledge about collection or meaning of the data is lost (e.g.,
129 because the collector/creator of the data is deceased, retired, or otherwise no longer active in
130 their field). Ultimately, balancing the data's value and risk of loss is essential for effective
131 prioritization of data rescue efforts.

132 2. Team creation

133 Data rescue takes a team, with different roles needed at different points in the rescue
134 process. We first consider those currently in possession of the data: *data creators* are typically
135 involved in generating the ideas that lead to the data's collection and retain the intellectual
136 property rights and responsibilities for the data, even if not directly involved in collecting or
137 managing the data products; *data collectors* generate or collect the original data and, therefore,
138 provide valuable input for documenting the data (see *Metadata creation*); and *data stewards* are
139 responsible for managing and maintaining the data (i.e., organizing and keeping data safely
140 archived, including instances where researchers have been bequeathed data or organizations act
141 as custodians of data collected by past employees). These roles are often played by the same
142 person, though not always. For example, in a mentee-mentor relationship such as that between a
143 graduate student and supervisor, the student may play all three roles as data creator, collector,
144 and (temporary) steward, while the advisor may retain the data long-term as the principal
145 investigator, thereby acting as data creator and (long-term) steward. Having at least one person
146 who is a data creator, collector, or steward, as part of the data rescue team is imperative for a
147 successful data rescue mission.

148 A *data management expert* is another key role. Usually, a data manager is the one that
149 plans the data lifecycle, but in a data rescue project this role is mainly focused on organizing and
150 documenting the digitized datasets. This person will have the skills to connect datasets, clean and
151 manage data, and compile previously unwritten information. Additionally, if any data have not
152 been entered into a digital format, a *data entry technician* will be an integral part of the team,
153 ensuring that all necessary data have been digitized in the appropriate format and validated
154 against the original records.

155 3. Metadata creation

156 *Metadata* are information about the data, typically contained in a file separate from the
157 dataset (Michener et al., 1997). Metadata generally describe the data collection process (e.g.,
158 types of data collected, methodology, and contributors), a description of the variables in the
159 dataset (e.g., column headings for tabular data; “data dictionary”), abbreviations, units of
160 measurement, and other relevant information necessary to understanding how the data were
161 generated and how to (re)use them (e.g., why some measurements are lacking; British Ecological
162 Society, 2018). We recommend early creation of the metadata, as this often informs the
163 remaining data rescue process and structure of the compiled dataset.

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

172 The metadata file format varies (often dependent on the type of data or chosen
173 repository), but one useful format is a text file written in Extensible Markup Language (XML;
174 see examples at <https://www.xmlfiles.com/xml/>). Tools like XML have been developed
175 specifically for writing and storing metadata in a format that is both human *and* machine
176 readable, not only ensuring that end users understand the data structure but also facilitating use
177 by other software/programming tools (e.g., search engines) that may rely on metadata being

178 available in a standardized form. Each variable is stored as a “tag,” and its description is stored
179 between tags. There is a variation of XML called Ecological Metadata Language (EML; Fegraus
180 et al., 2005; Jones et al., 2019; see <https://eml.ecoinformatics.org/>) which offers a set of
181 suggested tags specific to describing environmental data.

182 4. Data transfer and compilation

183 For the data rescue team to work effectively, all team members should have access to the
184 data and metadata files. However, this might only be possible if all files are already in a digital
185 format; any physical copies should either be photographed or scanned first or entrusted to the
186 team member responsible for data entry and validation. From there, discussion about how the
187 data should be compiled most effectively can ensue. While the details of data compilation will
188 need to be tailored to each dataset, the workflow should be as reproducible as possible. For
189 example, any edits made to the data should be done in a file separate from the original; a digital
190 file with untouched original data should always remain. Additionally, all major decisions should
191 be documented in the metadata.

192 In structuring the data, we recommend following Wickham’s (2014) tidy data principles,
193 which consist of 3 core concepts: (1) each variable has its own column, (2) each observation has
194 its own row, and (3) each type of observational unit is in its own data table (e.g., individual-level
195 measurements from a population, such as mass, in one table and population-level metrics, such
196 as abundance, in another). If there are multiple data tables, they should be connected to each
197 other by one or more variables that uniquely identify individual observations (i.e., primary keys
198 in a relational database; Codd, 1990). While we advocate for tidy data principles, as they are
199 most likely to generate a data structure that will be useful in subsequent analyses, sometimes
200 compromises will need to be made.

201 5. Data cleaning and validation

202 Following data entry and compilation, data cleaning can be one of the most time-
203 intensive steps of the data management process. Data cleaning is the process of identifying and
204 fixing issues, such as data entry errors or incomplete records. In addition to common steps like
205 correcting typographical or entry errors, data cleaning commonly includes checking for data

206 completeness (i.e., that the data from all records are fully and correctly transcribed) and
207 uniformity (i.e., that variables are recorded in a consistent way for all records, ensuring common
208 measurement units, etc.) and otherwise ensuring the data conform to expected standards. For
209 ecological or biodiversity data, other common data cleaning steps include checking for common
210 date formats (e.g., the International Organization for Standardization (ISO) 8601 standard
211 recommends date-time objects be recorded as YYYY-MM-DD hh:mm:ss + UTC offset),
212 ensuring geographic coordinates are complete and standardized (e.g., ISO 6709 applies to the
213 representation of spatial information), and correcting misspellings or synonyms in taxonomic
214 information. Many tools have been developed to help with specific aspects of data cleaning (e.g.,
215 the *taxize* package in R can be used to correct taxonomies; Chamberlin & Szocs, 2013).

216 Related to data cleaning, data validation involves the comparison of the dataset against a
217 set of assertions determined *a priori* (e.g., dry body mass of an organism should be less than its
218 wet mass) or *post hoc* (e.g., the ratio of dry to wet mass should be similar among replicates).
219 Data validation is important for ensuring data quality and integrity by evaluating the data against
220 a set of expectations to confirm the structure and content of the data are appropriate. In the case
221 of data rescue, unlike most recently or currently collected data, data validation may come with
222 the extra challenge that the original data creator or collector may be unreachable or deceased. As
223 such, having original members of the data team (Fig. 1, Step 2; see *Team creation*) is particularly
224 beneficial for effective data validation. Common data validation techniques include plotting the
225 data in various ways to assist with identifying incorrect or improbable values, checking that the
226 contents (e.g., number of unique values in a column) or dimensions of the data match
227 expectations, cross-checking data from different columns or tables for mutual compatibility, and
228 evaluating summary statistics or other outputs that characterize the data. In addition, many tools
229 exist to help with the data validation process, including open-source, “point-and-click” software
230 (e.g., OpenRefine) as well as a number of programming tools (e.g., the *assertr* and *validate*
231 packages in R; Fischetti, 2020; van der Loo & de Jonge, 2021).

232 Although the exact implementation of data cleaning and validation steps will vary by
233 dataset, many of the principles described in the *Data transfer and compilation* section are also
234 relevant here. Validation should be conducted in as reproducible a way as possible (e.g., in a
235 script file that can be run on the original or cleaned data files), and any errors identified should
236 be corrected without manipulating the original (raw) data files. Importantly, any changes should

237 be well documented (e.g., as comments in the script or as notes in the metadata), as should the
238 rationale behind the corrections.

239 6. Data archiving

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

247 There is a growing movement to archive data on public (and open) data repositories
248 rather than, or in addition to, private or institutional systems (e.g., a lab hard drive). Many
249 governments and funding agencies have recently implemented new data management protocols
250 that either encourage or mandate the archiving, though not necessarily sharing, of all data
251 generated using their resources (see below; e.g., Canada’s [Tri-agency Research Data](#)
252 [Management Policy](#)). With each year that passes after a publication, data that have not been
253 publicly archived are 17% less likely to be recoverable (Vines et al., 2014; see also Tedersoo et
254 al., 2021). As such, we consider public archiving to be an essential part of data rescue, since
255 private archiving does not mitigate the possibility that data will need to be “re-rescued” in the
256 future. Once the data and metadata are compiled and validated, they should be placed in a data
257 repository to maintain the data in a secure and retrievable format for the future. Importantly, the
258 push for public archiving does not contradict the need for privacy or sensitivity associated with
259 some datasets; it is possible to publicly archive data while maintaining restrictions on when and
260 how the data are accessed (see below). We suggest, however, that most environmental data
261 should be openly accessible upon archiving, with some clear exceptions (e.g., data pertaining to
262 threatened species or Indigenous data sovereignty).

263 There are many data repositories from which to choose (see [r3data.org](#) for a
264 comprehensive list), with some being very generalized (e.g., Dryad, Dataverse, Figshare,
265 Zenodo) and others catering to specific types of data (e.g., DataONE for environmental data,
266 GenBank for genetic sequences). Data repositories tend to use a distributed (i.e., decentralized)

267 approach to storing data and have contingency plans in place to ensure the longevity of archived
268 datasets. Which repository to choose will also be influenced by whether the data will remain
269 private or be made openly and publicly accessible upon upload, or sometime in the near future
270 (Roche et al., 2014). Some repositories allow for the long-term storage of datasets regardless of
271 whether they are made openly available (e.g., Dataverse); others require that the data be open
272 access if they are to be hosted by the repository (e.g., Dryad). Many archives also offer an option
273 to place an embargo, or delay, on the publication of data. Most data repositories will establish a
274 Digital Object Identifier (DOI), a unique identifier which will remain constant for the lifetime of
275 the object, even if the object or metadata change. If the data will be openly available, we suggest
276 explicitly stating the terms of use, such as noting that authors should be contacted if the data are
277 to be included in a publication or adding a copyright statement, such as those from Creative
278 Commons (e.g., CC0, CC-BY, etc.).

279 7. Data sharing

280 The final step in the data rescue workflow is to ensure that the data meet open science
281 standards and that their use can be tracked. Open science principles include transparency,
282 participation and accessibility (Bartling & Friesike, 2014). These values can be addressed in
283 different ways, sometimes making the process overwhelming for researchers who are not trained
284 in data management. The FAIR and CARE principles, the first of which focuses on how data can
285 be made useful and the second on how we can promote justice through responsibly sharing open
286 data, summarize ways these values can be met through a combination of actions.

287 The **FAIR** principles aim to improve **F**indability, **A**ccessibility, **I**nteroperability and
288 **R**eusability of datasets (Wilkinson et al., 2016). Providing human- and machine-readable
289 metadata improves both the findability and accessibility of a dataset. Combined with proper
290 archiving and identification, strong metadata helps increase the discoverability of datasets. As
291 mentioned in the *Data archiving* section, adding a DOI makes the data trackable and citable,
292 improving the reproducibility of analyses. A comprehensive metadata file also allows
293 interoperability, or the ability of the data to be combined with other datasets in different ways
294 and in different systems. Additionally, accessibility and reusability can be achieved through
295 licenses, which explicitly describe the usage and attribution rights of the data.

296 The **CARE** principles focus on datasets that used traditional knowledge or benefited
297 somehow from Indigenous lands, promoting transparency and participation of open data (Carroll
298 et al., 2020). They aim to address and encourage consideration of the Collective benefit for
299 Indigenous Peoples, **A**uthority to control (recognizing Indigenous data sovereignty),
300 **R**esponsibility to be respectful with Indigenous Peoples involved in the dataset collection, and
301 **E**thics (by assuring participation of Indigenous Peoples in the assessment of benefits, harms and
302 usability of the data; Carroll et al., 2020). These principles are meant to begin addressing the
303 larger, complicated history of colonialism in ecology, evolution, and related disciplines. While
304 these guidelines were written with current and future data collection in mind, they are equally
305 applicable to and important for previously collected data, and we recommend that all researchers
306 who are rescuing datasets take these principles into consideration.

307 Conclusion

308 Ultimately, we hope to reach a point where data rescue is no longer needed. This requires
309 researchers, funding agencies, and publishers to align their views around ethical and professional
310 obligations to archive data and make them publicly accessible where appropriate. It also requires
311 a culture change that sees best practices in data managing, archiving, and publicly sharing data
312 become the default in publicly funded research. To achieve this goal, data sharing and
313 accessibility need to be prioritized as critical components of the scientific enterprise. We believe
314 that the solution to shifting the culture around data sharing is two-fold. First, there must be
315 continued, long-term investment in data management (Mons, 2020). Such investment includes
316 not only infrastructure but also training and support for students and personnel (Renaut et al.,
317 2018; Soeharjono & Roche, 2021). Additionally, publishers, employers, and funding agencies
318 must require some level of accountability from researchers to preserve data in accessible formats
319 and, if appropriate, make the data openly available to anyone interested (Mons, 2020). Until
320 these institutional-level paradigm shifts occur, however, smaller-scale and innovative data rescue
321 is an integral part of environmental data curation.

322 Currently, training in data management and shifting regulations regarding data
323 availability have, rightfully, focused on present and future data and data practices. With such a
324 strong eye to the future, however, much of the data of the past is being left behind. Data rescue

325 presents an opportunity to mitigate this loss of historical data while also providing additional,
326 less tangible benefits. In the CIEE Living Data Project, our mission of breathing life into
327 languishing data is concomitant with training the next generations of scientists in data
328 management best practices and forging connections amongst researchers across a wide variety of
329 career stages and trajectories, thus ensuring the longevity of scientific knowledge and preparing
330 students for a data-rich future.

331 Acknowledgements

332 The Living Data Project (LDP) is a collaborative initiative by researchers at institutions
333 across Canada: University of British Columbia-Vancouver, University of British Columbia-
334 Okanagan, University of Regina, McGill University, and Université de Montréal. The authors
335 recognize that we live and work on the traditional, ancestral, treaty, and unceded territories of
336 many Indigenous peoples, including the Coast Salish Peoples, xwməθkwəy̓əm (Musqueam),
337 Syilx (Okanagan), nēhiyawak (Cree), anihšīnāpēk (Saulteaux), Dakota, Lakota, Nakoda,
338 Attawanderon, Mississaugas, kanien'kehà:ka (Mohawk), and Haudenosaunee, and the homeland
339 of the Métis/Michif Nation.

340 We thank LDP members and partner organizations for their thoughtful discussions on
341 data rescue, best practices in data management, and fostering more open and transparent science.
342 We are grateful to our project partners and the data rescue interns, particularly those whose work
343 we reference: Jenna Loesberg and Amelia Hesketh, who worked in partnership with Drs. Ellen
344 Macdonald and Justine Karst, and Andrea Brown, who worked with Dr. Harold Eyster.

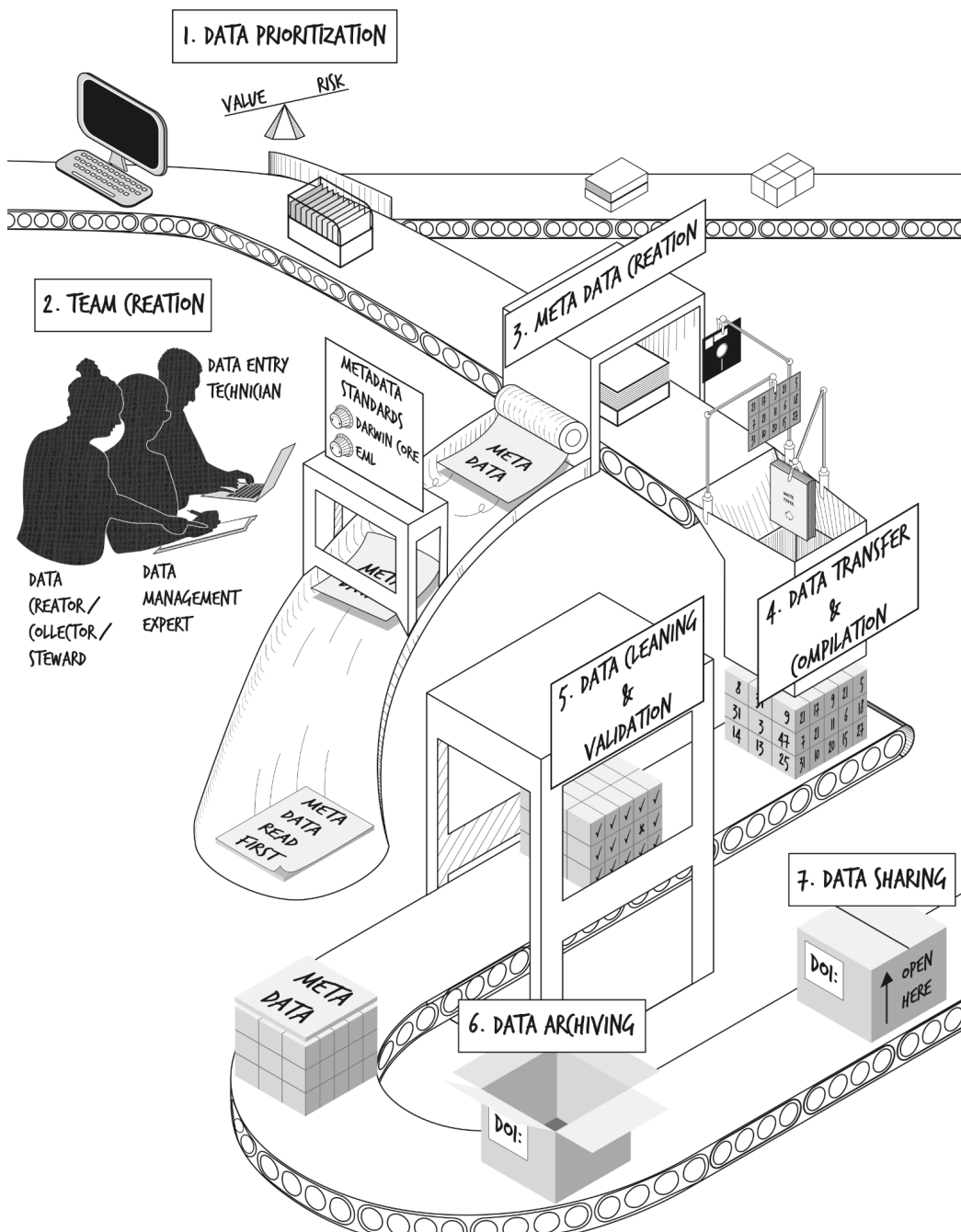
345 Funding for the LDP is provided through a Collaborative Research and Training
346 Experience (CREATE) grant to the Canadian Institute of Ecology and Evolution from the
347 Natural Sciences and Engineering Research Council of Canada. EKB, JBB, and GTH were
348 supported by the CREATE grant; EKB was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

351

352 Author contributions

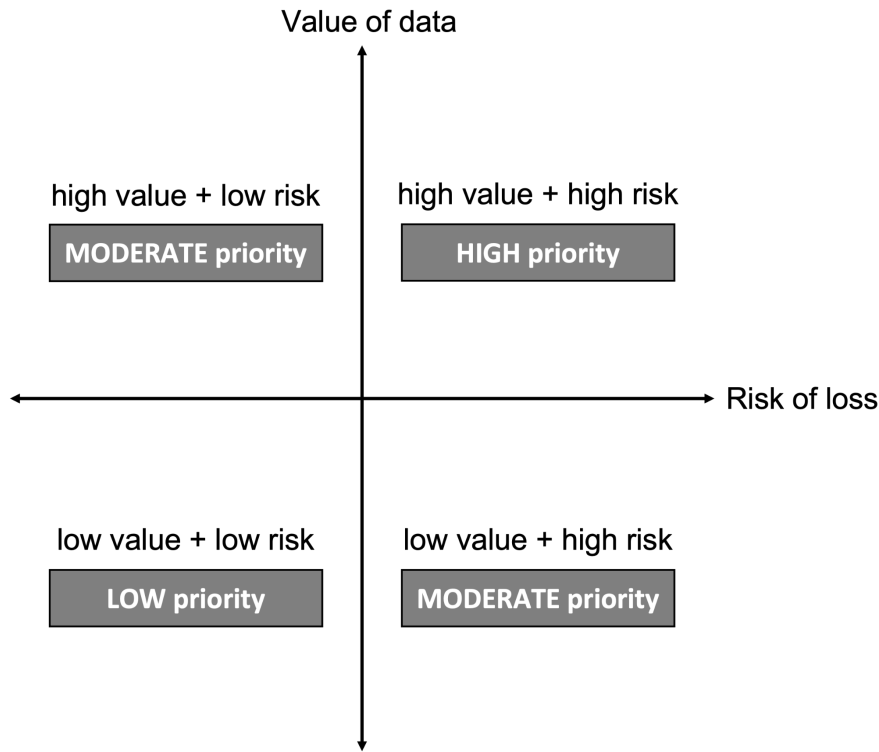
353 EKB, JBB, DGR, and DSS proposed the initial idea for the manuscript; all authors
354 contributed to developing the methods of data rescue we describe and subsequent discussions
355 about the paper. EKB, JBB, and GTH wrote the first draft. DSS created the first draft of the
356 figure. All authors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

358 **Figure 1.** Steps in the data rescue assembly line.



359
 360 Figure 1. Steps in the data rescue assembly line. First, data must be prioritized for rescue (Step
 361 1). After team creation (Step 2) and metadata creation (Step 3), the data must be transferred and
 362 compiled into a logical format (Step 4). After data cleaning and validation (Step 5) is complete,
 363 the finalized data and metadata should be archived on a long-term data repository (Step 6). The
 364 ultimate goal is to have the rescued data openly available for reuse (Step 7).

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



367
368 Figure 2. Prioritizing data for rescue: balancing the value of the data and its risk of loss. With
369 many datasets in need of preservation and limited resources, the first step in the data rescue
370 process requires developing a list of priorities for consideration and identifying relevant datasets
371 (Fig. 1). We consider data prioritization to be a balance between the assessed value of a dataset
372 in question and the potential risk of its loss in the absence of intervention (see *Data prioritization*
373 under *Guidelines*).

374 **Box 1. Spilt oil, spent money, and lost data**

375 In 1989, the oil tanker *Exxon Valdez* struck the Bligh Reef in Prince William Sound, less than 2.5
376 km from the Alaskan shore. As a result, approximately 37,000 tonnes of crude oil spilled into the
377 sound, leading to catastrophic short- and long-term ecological consequences. The *Exxon Valdez*
378 Oil Spill Trustee Council (EVOSTC) was established in 1991 to oversee the spending of funds
379 from a civil settlement in 1991 between *Exxon*, the United States federal government and the
380 state government of Alaska. A large portion of funds were directed towards determining and
381 monitoring the impacts of the oil spill on oceanographic, environmental, and ecological
382 conditions. Prior to 2003, there was no requirement for data preservation or availability;
383 afterwards, all projects were awarded under explicit conditions from EVOSTC that data be
384 preserved and made publicly available (Jones et al., 2018). In their annual report from 2010, the
385 EVOSTC notes that some \$151.2 million USD were spent on “research, monitoring, and general
386 restoration” during 1992-2010 fiscal years (EVOSTC, 2012). Most funding went to state and
387 federal agencies, though a few projects were awarded to universities, professional societies,
388 consultants, and other private entities (EVOSTC, 2018).

389 From 2012-2014, a group of researchers from the National Center for Ecological
390 Analysis & Synthesis (NCEAS) worked to recover the historical datasets funded by EVOSTC,
391 focusing specifically on data collected between 1989-2010 (Jones et al., 2018). Of the 419
392 projects funded by EVOSTC during this time, only 27% of the datasets were able to be
393 recovered; after a total of 5 years hunting down datasets, this grew to 30% (Jones et al., 2018).
394 Using these numbers, we can roughly estimate the money spent on research for which the data
395 are unrecoverable (70% of datasets): **~\$105 million USD was spent collecting data that are no**
396 **longer recoverable and, therefore, effectively nonexistent to science.** While we do not know
397 the distribution of years from which data were recovered or how money was allocated by year,
398 this is likely a conservative estimate given that the original cost does not include the first 3 years
399 following the spill, when extensive ecological assessments would have been completed.

400 The NCEAS group also noted the reasons for their inability to recover the data. Instances
401 in which data collectors specifically stated that the data were lost or unrecoverable were rare

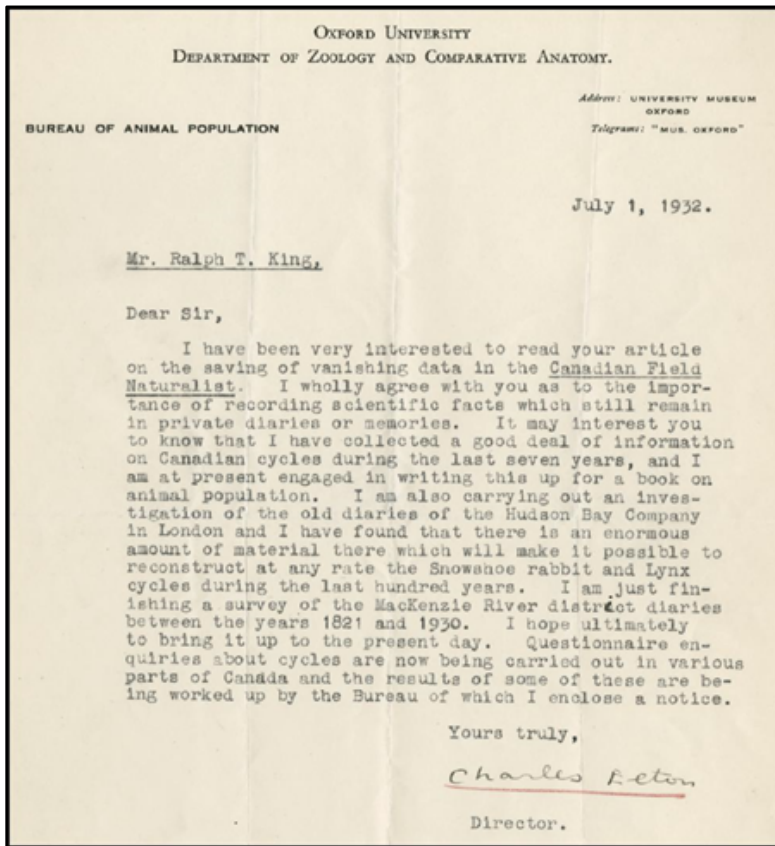
402 (Jones et al., 2018). Instead, over 80% of reasons for unrecovered data due to a lack or failure of
403 communication (~50% categorized as “communication lost”); the authors of the final report,
404 however, interpret this lack of communication as an unwillingness or inability by the data
405 owners to share data (Jones et al., 2018), highlighting the importance of proper documentation
406 and public archiving of data for longevity.
407

408 **Box 2. From fur trappers to fundamental ecological theory**

409 Charles Sutherland Elton (1900-1991) was a British ecologist, whose major contributions include
410 work on population cycling and community dynamics. In 1932, Elton established and became
411 the first director of the Bureau of Animal Population at Oxford University and the inaugural
412 editor of the *Journal of Animal Ecology* later the same year. As part of their work on population
413 cycling, Elton and his colleague Mary Nicholson endeavoured to recover historical records on
414 the number of Canada lynx (*Lynx canadensis*) and snowshoe hare (*Lepus americanus*) furs
415 collected by Hudson's Bay Company (HBC) trappers in Canada (Box 2.S1; Elton & Nicholson,
416 1942). In an effort spanning more than 15 years, Elton and Nicholson used these and other
417 records to collate information on trapping activities across the whole of Canada from 1886 to
418 1940 (and, for some areas, as far back as 1821). Much of this work was akin to data rescue,
419 including correspondence between Elton and the original data owners (see Elton's description of
420 the process in a letter to Ralph King in Box. 2.1; Elton & Nicholson, 1942), collation of data
421 from different sources, as well as data cleaning and validation. This work was not only central to
422 compiling among the longest time series of animal populations and revealing the now classical
423 example of ~12-year population cycles in snowshoe hare and lynx abundances (Box 2.2), but has
424 spurred an entire field of ecology (population/community cycling) and many decades of
425 ecological research in the Canadian Arctic. This is just one, elegant example of the immense
426 value of historical data, even those from unconventional places (like the ledgers of a colonial fur-
427 trading company), and the importance of working to identify and preserve them.

428

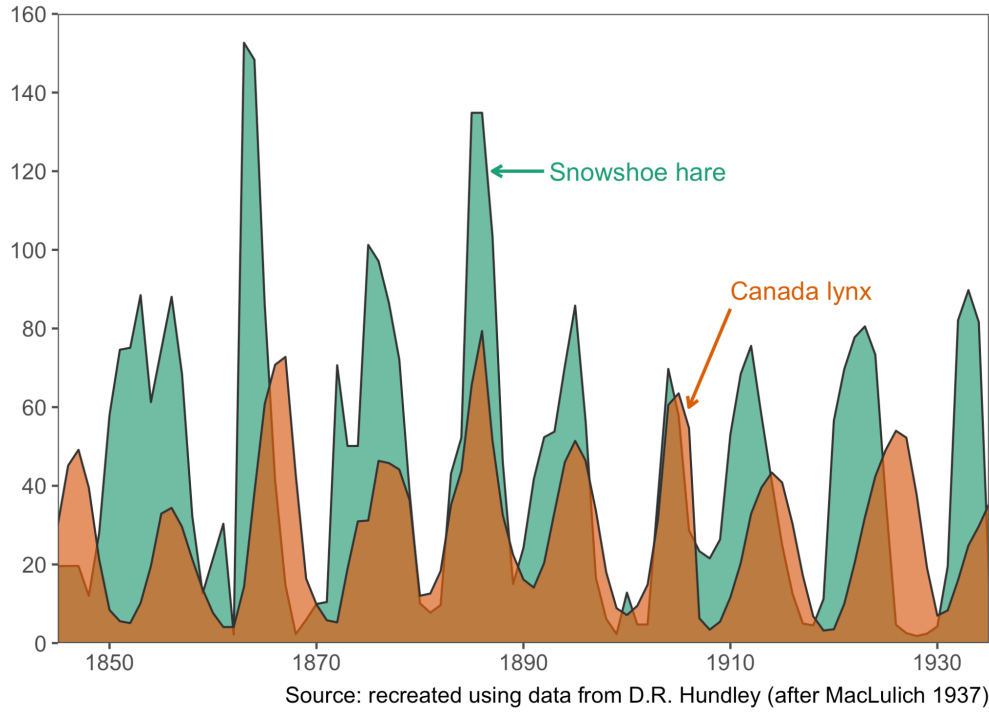
429 **Box 2.1.** *Letter from Charles S. Elton to Ralph T. King, dated 01 July 1932.*
430 In this letter, Elton expresses his interest in King's recent article on "saving vanishing data"
431 (King, 1932), which regarded many aspects of what we are calling "data rescue" and was itself
432 based on a paper of the same name written some three decades earlier (Haddon, 1903). This
433 letter was provided to us courtesy of Dr. Adam T. Ford and is available through the Elton
434 Archive at Oxford (Elton, 1932). A transcription of the letter's text is available in the Supporting
435 Information (Box 2.S2).



436

437

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



440
441

442 **Box 3. Recent data rescue examples from the Living Data** 443 **Project.**

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

456 *Seeing the Forest Data for the Trees*

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

462 Dr. George H. La Roi was a professor of forest ecology at the University of Alberta for
463 35 years. In a 2016 email to colleagues, he implored for assistance archiving his extensive long-
464 term survey data from the boreal forests of northern Alberta. Before this could be accomplished,
465 however, Dr. La Roi passed away in 2018. Upon his passing, La Roi's children bequeathed his
466 legacy of highly valuable data to his former colleague, Dr. Ellen Macdonald, who had earlier
467 taken over sampling some of his long-term plots. With no living data creator and much of the
468 data in unorganized boxes containing unsorted datasheets, various documents, CD-ROMs, and
469 picture slides, the data was at high risk of being lost. Macdonald determined she would be unable

470 to tackle the boxes of materials alone and joined forces with her colleague, Dr. Justine Karst,
471 who had also come into possession of some of La Roi's boxes of data by way of University's
472 Botanic Garden. Together, they applied for an LDP data rescue internship. The value and
473 precarious circumstances of the dataset made it a high priority for rescue.

474 Over the course of two data rescue internships, Jenna Loesberg and Amelia Hesketh,
475 along with several undergraduate data entry technicians, sorted, entered, and digitized the data.
476 They recovered data recorded at two different locations (Hondo-Slave Lake and Athabasca Oil
477 Sands regions), both of which included measures of vascular plant cover, bryoid cover, and
478 forest mensuration, among other datasets. Some data were stored as printed scans of hand-filled
479 datasheets, and thus required digitization. Other data, which had already been entered and
480 digitized, were stored in hundreds of text files which required extensive reformatting and
481 cleaning before they could be compiled into usable datasets. Metadata also needed to be written
482 and consolidated into one document for future reuse; while most of the data had clear
483 documentation, some data were lost, since no documentation about the meaning of some variable
484 names or values in a column was recovered. With this work completed, the data and metadata of
485 this rich and expansive dataset will be archived and made publicly available through University
486 of Alberta's Dataverse repository and eventually accompanied by a data paper.

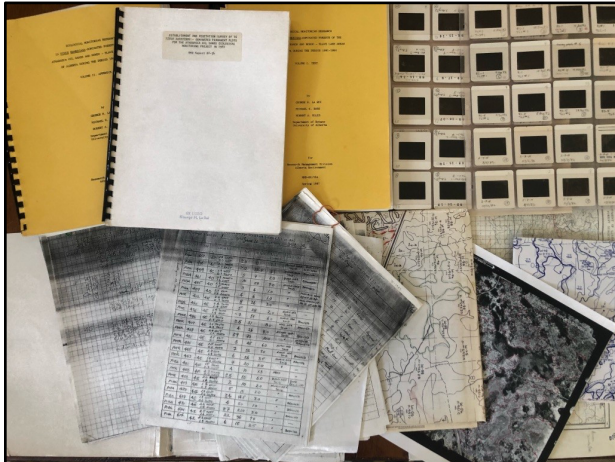
487

488 **Box 3.1.** Photograph of researchers collecting data in the Athabasca Oil Sands region of
489 northern Alberta in 1982 at one of 16 sites established by Dr. George La Roi. Image credit:
490 unknown.



491
492

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



496
497

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

499 The archived theses and dissertations of former graduate students represent a rich source
500 of historical data. In particular, those prepared prior to the advent of modern computer
501 technologies and software, such as word processors and tools for statistical analysis, may contain
502 troves of raw and summary data that remain un-digitized. As a result, the reuse of any raw or
503 summarized data from the thesis would first require data extraction and digitization.

504 Determining how biodiversity has changed in response to human activity and land use is
505 central to understanding the impacts of these environmental changes and predicting the potential
506 for future declines. In a data rescue project proposed to us by a then-doctoral student, Dr. Harold
507 Eyster, Andrea Brown worked to secure the data contained in three University of British
508 Columbia graduate theses (Weber, 1972; Lancaster, 1976; Melles, 2000). While the specific
509 questions and research topics differed between these theses, all three surveyed bird abundances
510 in the same (or nearby) sites in Greater Vancouver, British Columbia, over the span of several
511 decades, and in combination present an opportunity to establish a baseline against which to
512 compare current and future trends (Box. 3.3). This project was prioritized by the LDP because
513 the data were both at-risk (much of the data existed only in non-digital formats and none of the
514 datasets are in active use) and high value (the data provide a valuable frame of reference for
515 studying changes in urban bird diversity).

516 During the internship, Brown first transcribed the data from the earlier two of the theses,
517 Weber (1972) and Lancaster (1976), which were archived as scans of typewritten documents.
518 Among other challenges, digitization required the conversion of non-standard data types (Box
519 3.4) into “tidy” forms that could be used and interpreted programmatically. Data from the third
520 thesis, Melles (2000), were made available by the original author in a Microsoft Excel®
521 spreadsheet and so only required cleaning, manipulation, and conversion to a non-proprietary
522 format. Later work included efforts to rationalize the datasets so that they might be used in
523 combination with each other (e.g., standardizing column names and combining similar tables
524 into a single file). Given the extensive data manipulation required, clear metadata were
525 developed to document the various steps taken to generate the final datasets and document other
526 details from the theses that were not captured during the digitization process. The data have been
527 archived on the UBC Dataverse repository (Brown et al., 2021a, 2021b, 2021c) and linked with
528 the original theses.

529 **Box 3.3.** Comparison of the historical and current appearance of one of the sampling locations
 530 for urban bird surveys conducted in Vancouver, British Columbia, Canada. Photographs show
 531 the view looking west from the intersection of 24th Avenue West at Wallace Street (49.251°N,
 532 123.191°W).



533
 534 April 1970 (Image credit: W.C. Weber) October 2021 (Image credit: C.N. Nemeth)
 535

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

Habitat Features	C>1.5	E<1.5	E>1.5	HERB	WEED	FRD1	FRD2	FRD3	HRD	TDD	TDC	TDE	FOOD	TOTVEG
SLANT					++								++	
FLAT														
ROAD														
LANE														
PVT+S														
LAWN														
D>7.5	+		++		++	++	+	++		++	++		++	
D>7.5		++			++	++		++		++	++		++	
C>7.5	++		@		++	++	++	++		++	++		++	
C>7.5	++				++	++	++	++		++	++		++	
RBEVG					++	++	++	++		++	++		++	
D<1.5	++	@	++		++	++	++	++		++	++		++	
D>1.5	++	++	++		++	++	++	++		++	++		++	
C<1.5	++	++	++		++	++	++	++		++	++		++	
C>1.5	++	++	++		++	++	++	++		++	++		++	
E<1.5	++	++	++		++	++	++	++		++	++		++	
E>1.5	++	++	++		++	++	++	++		++	++		++	
HERB				*										
WEED					*									
FRD1						*	@	++		++	++		++	
FRD2							@	++		++	++		++	
FRD3								++		++	++		++	
HRD								++		++	++		++	
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

540
 541

542 References

- 543 Bartling S, Friesike S. 2014 *Opening Science: The Evolving Guide on How the Internet is*
544 *Changing Research, Collaboration and Scholarly Publishing*. Springer Open.
545 (ISBN:978-2-319-00025-1)
- 546 BES. 2018 A guide to data management in ecology and evolution. BES Guides to Better Science.
547 British Ecological Society, London, UK. ([https://www.britishecologicalsociety.org/wp-](https://www.britishecologicalsociety.org/wp-content/uploads/2019/06/BES-Guide-Data-Management-2019.pdf)
548 [content/uploads/2019/06/BES-Guide-Data-Management-2019.pdf](https://www.britishecologicalsociety.org/wp-content/uploads/2019/06/BES-Guide-Data-Management-2019.pdf))
- 549 Borgman CL. 2018 Open data, grey data, and stewardship: universities at the privacy frontier.
550 *Berkeley Tech. Law J.* **33**, 365-412. (doi:[10.15779/Z38B56D489](https://doi.org/10.15779/Z38B56D489))
- 551 Broman KW, Woo KH. 2018 Data organization in spreadsheets. *Am. Stat.* **72**, 2-10.
552 (doi:[10.1080/00031305.2017.1375989](https://doi.org/10.1080/00031305.2017.1375989))
- 553 Brown A, Eyster H, Lancaster RK. 2021a Data for: Bird communities in relation to the structure
554 of urban habitats. *Scholars Portal Dataverse*. (doi:[10.5683/SP2/YD6N7C](https://doi.org/10.5683/SP2/YD6N7C))
- 555 Brown A, Eyster H, Melles SJ. 2021b Data for: Effects of landscape and local habitat features on
556 bird communities: a study of an urban gradient in greater Vancouver. *Scholars Portal*
557 *Dataverse*. (doi:[10.5683/SP2/BPLPAP](https://doi.org/10.5683/SP2/BPLPAP))
- 558 Brown A, Eyster H, Weber WC. 2021c Data for: Birds in cities: a study of populations, foraging
559 ecology and nest-sites of urban birds. *Scholars Dataverse Portal*.
560 (doi:[10.5683/SP2/K5LMLA](https://doi.org/10.5683/SP2/K5LMLA))
- 561 Buck S. 2021 Beware performative reproducibility. *Nature* **595**, 151. (doi:[10.1038/d41586-021-](https://doi.org/10.1038/d41586-021-01824-z)
562 [01824-z](https://doi.org/10.1038/d41586-021-01824-z))
- 563 Buxton RT, Nyboer EA, Pigeon KE, Raby GD, Rytwinski T, Gallagher AJ, Schuster R, Lin H-
564 Y, Fahrig L, Bennett JR, et al. 2021 Avoiding wasted research resources in conservation
565 science. *Converv. Sci. Pract.* **3**, e329. (doi:[10.1111/csp2.329](https://doi.org/10.1111/csp2.329))
- 566 Carroll SR, Garba I, Figueroa-Rodríguez OL, Holbrook J, Lovett R, Materechera S, Parsons M,
567 Raseroka K, Rodriguez-Lonebear D, Rowe R, et al. 2020 The CARE principles for
568 Indigenous data governance. *Data Sci. J.* **19**, 43. (doi:[10.5334/dsj-2020-043](https://doi.org/10.5334/dsj-2020-043))
- 569 Chamberlain S, Szocs E. 2013 taxize – taxonomic search and retrieval in R. *F1000 Res.* **2**, 191.
570 (doi:[10.12688/f1000research.2-191.v2](https://doi.org/10.12688/f1000research.2-191.v2))

571 Codd EF. 1990 *The Relational Model for Database Management: Version 2*. Addison-Wesley
572 Longman Publishing. (ISBN:978-0-201-14192-4)

573 Cook RB, Olson RJ, Kanciruk P, Hook LA. 2001 Best practices for preparing ecological data
574 sets to share and archive. *Bull. Ecol. Soc. Am.* **82**, 138-141. (
575 <https://www.jstor.org/stable/20168543>)

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

578 Elton CS. 1932 Letter to Ralph T. King, 01 July. MS. Eng. c3328 A72, Elton Archives, Weston
579 Library, University of Oxford.

580 Elton CS, Nicholson M. 1942 The ten-year cycle in numbers of the lynx in Canada. *J. Anim.*
581 *Ecol.* **11**, 215-244. (<https://www.jstor.org/stable/1358>)

582 EVOSTC. 2012 2010 Annual Report. *Exxon Valdez Oil Spill Trustee Council*.
583 (<https://evostc.state.ak.us/media/4411/2010annualreport.pdf>)

584 EVOSTC. 2018 *Exxon Valdez Oil Spill Final and Annual Reports*. *Exxon Valdez Oil Spill*
585 *Trustee Council*. (<https://evostc.state.ak.us/media/4291/finalandannualreports.pdf>)

586 Fegraus EH, Andelman S, Jones MB, Schildhauer M. 2005 Maximizing the value of
587 ecological data with structured metadata: an introduction to Ecological Metadata
588 Language (EML) and principles for metadata creation. *Bull. Ecol. Soc. Am.* **86**, 158-68.
589 (<http://www.jstor.org/stable/bullecosociamer.86.3.158>)

590 Fischetti T. 2020 assertr: assertive programming for R analysis pipelines. R package version 2.7.
591 (<https://CRAN.R-project.org/package=assertr>)

592 Gatti G, Bianchi CN, Parravicini V, Rovere A, Peirano A, Montefalcone M, Massa F, Morri C.
593 2015 Ecological change, sliding baselines and the importance of historical data: lessons
594 from combining observational and quantitative data on a temperate reef over 70 years.
595 *PLoS One* **10**, e0123268. (doi:[10.1371/journal.pone.0118581](https://doi.org/10.1371/journal.pone.0118581))

596 Haddon AC. 1903 The saving of vanishing data. *Pop. Sci. Mon.* **63**, 222-229.
597 ([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)
598 [_Saving_of_Vanishing_Data](https://en.wikisource.org/wiki/Popular_Science_Monthly/Volume_62/January_1903/The_Saving_of_Vanishing_Data))

599 Jones MB, Blake R, Couture J, Ward C. 2018 Collaborative data management and holistic
600 synthesis of impacts and recovery status associated with the *Exxon Valdez* oil spill. *Exxon*
601 *Valdez Oil Spill Long-Term Monitoring Program (Gulf Watch Alaska) Final Report*

602 (project 16120120). *Exxon Valdez* Oil Spill Trustee Council, Anchorage, Alaska.
603 ([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)
604 [2018-Final-Report.pdf](http://www.gulfwatchalaska.org/wp-content/uploads/2018/08/16120120-Jones-et-al.-2018-Final-Report.pdf))
605 Jones MB, O'Brien M, Mecum B, Boettiger C, Schildhauer M, Maier M, Whiteaker T, Earl S,
606 Chong S. 2019 Ecological Metadata Language version 2.2.0. KNB Data Repository.
607 (doi:[10.5063/F11834T2](https://doi.org/10.5063/F11834T2))
608 King RT. 1932 The saving of vanishing data. *Can. Field Nat.* **46**, 108-111.
609 (<https://www.biodiversitylibrary.org/ia/canadianfieldnat1932otta/#page/134/mode/1up>)
610 Lancaster RK. 1976 Bird communities in relation to the structure of urban habitats. Thesis.
611 Department of Zoology, University of British Columbia. (doi:[10.14288/1.0093863](https://doi.org/10.14288/1.0093863))
612 Lennox RJ, Harcourt R, Bennett JR, Davies A, Ford AT, Frey RM, Hayward MW, Hussey ME,
613 Iverson SJ, Kays R, et al.. (2020). A novel framework to protect animal data in a world of
614 biosurveillance. *BioScience* **70**, 468-476. (doi:[10.1093/biosci/biaa035](https://doi.org/10.1093/biosci/biaa035))
615 van der Loo MPJ, de Jonge E. 2021 Data validation infrastructure for R. *J. Stat. Softw.* **97**, 1–31.
616 (doi:[10.18637/jss.v097.i10](https://doi.org/10.18637/jss.v097.i10))
617 McClenachan L, Ferretti F, Baum JK. 2012 From archives to conservation: why historical data
618 are needed to set baselines for marine animals and ecosystems. *Conserv. Lett.* **5**, 349-359.
619 (doi:[10.1111/j.1755-263X.2012.00253.x](https://doi.org/10.1111/j.1755-263X.2012.00253.x))
620 Melles SJ. 2000 Effects of landscape and local habitat features on bird communities: a study of
621 an urban gradient in Greater Vancouver. Thesis. Department of Forest Sciences,
622 University of British Columbia. (doi:[10.14288/1.0099590](https://doi.org/10.14288/1.0099590))
623 Michener WK, Brunt JW, Helly JJ, Kirchner TB, Stafford SG. 1997 Non-geospatial metadata for
624 the ecological sciences. *Ecol. Appl.* **7**, 330-342. (doi:[10.1890/1051-](https://doi.org/10.1890/1051-0761(1997)007[0330:NMFTES]2.0.CO;2)
625 [0761\(1997\)007\[0330:NMFTES\]2.0.CO;2](https://doi.org/10.1890/1051-0761(1997)007[0330:NMFTES]2.0.CO;2))
626 Mons B. 2020 Invest 5% of research funds in ensuring data are reusable. *Nature* **578**, 491.
627 (doi:[10.1038/d41586-020-00505-7](https://doi.org/10.1038/d41586-020-00505-7))
628 O'Dea RE, Parker TH, Chee YE, Culina A, Drobniak SM, Duncan DH, Fidler F, Gould E, Ihle
629 M, Kelly CD. 2021 Towards open, reliable, & transparent ecology and evolutionary
630 biology. *BMC Biol.* **19**, 1-5. (doi:[10.1186/s12915-021-01006-3](https://doi.org/10.1186/s12915-021-01006-3))
631 Powers SM, Hampton SE. 2019 Open science, reproducibility, and transparency in ecology.
632 *Ecol. Appl.* **29**, e01822. (doi:[10.1002/eap.1822](https://doi.org/10.1002/eap.1822))

633 Renaut S, Budden AE, Gravel D, Poisot T, Peres-Neto P. 2018 Management, archiving, and
634 sharing for biologists and the role of research institutions in the technology-oriented age.
635 *Bioscience* 68, 400-411. ([10.1093/biosci/biy038](https://doi.org/10.1093/biosci/biy038))

636 Roche DG, Berberi I, Dhane F, Lauzon F, Soeharjono S, Dakin R, Binning SA. 2021 The quality
637 of open datasets shared by researchers in ecology and evolution is moderately repeatable
638 and slow to change. *EcoEvoRxiv*. (doi:[10.32942/osf.io/d63js](https://doi.org/10.32942/osf.io/d63js))

639 Roche DG, Kruuk LEB, Lanfear R, Binning SA. 2015 Public data archiving in ecology and
640 evolution: how well are we doing? *PLoS Biol.* **13**, e1002295.
641 (doi:[10.1371/journal.pbio.1002295](https://doi.org/10.1371/journal.pbio.1002295))

642 Roche DG, Lanfear R, Binning SA, Haff TM, Schwanz LE, Cain KE, Kokko H, Jennions MD,
643 Kruuk LE. 2014 Troubleshooting public data archiving: suggestions to increase
644 participation. *PLoS Biol.* **12**, e1001779. (doi:[10.1371/journal.pbio.1001779](https://doi.org/10.1371/journal.pbio.1001779))

645 Smith R, Roberts I. 2016 Time for sharing data to become routine: the seven excuses for not
646 doing so are all invalid. *F1000 Res.* **5**, 781. (doi:[10.12688/f1000research.8422.1](https://doi.org/10.12688/f1000research.8422.1))

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

650 Soeharjono S, Roche DR. 2021 Reported individual costs and benefits of sharing open data
651 among Canadian academic faculty in ecology and evolution. *BioScience* **71**, biab024.
652 (doi:[10.1093/biosci/biab024](https://doi.org/10.1093/biosci/biab024))

653 Tedersoo L, Kungas R, Oras E, Köster K, Eenmaa H, Leijen Ä, Pedaste M, Raju M, Astapova A,
654 Lukner H, Korgerman K, Sepp T. 2021 Data sharing practices and data availability upon
655 request differ across scientific disciplines. *Sci. Data* **8**, 192. (doi:[10.1093/biosci/biab024](https://doi.org/10.1093/biosci/biab024))

656 Vines TH, Albert AYK, Andrew RL, Débarre F, Bock DG, Franklin MT, Gilbert KJ, Moore J,
657 Renault S, Rennison DJ. (2014). The availability of research data declines rapidly with
658 article age. *Curr. Biol.* **24**, 94-97. (doi:[10.1016/j.cub.2013.11.014](https://doi.org/10.1016/j.cub.2013.11.014))

659 Walter M, Suina M. 2018 Indigenous data, indigenous methodologies and indigenous data
660 sovereignty. *Int. J. Soc. Res. Methodol.* **22**, 233-243.
661 (doi:[10.1080/13645579.2018.1531228](https://doi.org/10.1080/13645579.2018.1531228))

662 Weber WC. 1972 Birds in cities: a study of populations, foraging ecology and nest-sites of urban
663 birds. Thesis. Department of Zoology, University of British Columbia.
664 (doi:[10.14288/1.0101293](https://doi.org/10.14288/1.0101293))

665 Wickham H. 2014 Tidy Data. *J. Stat. Softw.* **59**, 1-23. (doi:[10.18637/jss.v059.i10](https://doi.org/10.18637/jss.v059.i10))

666 Wilkinson MD, Dumontier M, Aalbersberg IJJ, Appleton, G., Axton, M., Baak, A., Blomberg N,
667 Boiten J-W, Bonino da Silva Santos L, Bourne PE, et al. 2016 The FAIR Guiding
668 Principles for scientific data management and stewardship. *Sci. Data* **3**, 160018.
669 (doi:[10.1038/sdata.2016.18](https://doi.org/10.1038/sdata.2016.18))

670 Willis KJ, Araújo MB, Bennett KD, Figueroa-Rangel B, Freud CA, Myers N. 2007 How can a
671 knowledge of the past help to conserve the future? Biodiversity conservation and the
672 relevance of long-term ecological data. *Phil. Trans. R. Soc. Lond. B* **362**, 175-187.
673 (doi:[10.1098/rstb.2006.1977](https://doi.org/10.1098/rstb.2006.1977))

674 White EP, Baldrige E, Brym ZT, Locey KJ, McGlinn DJ, Supp SR. (2013). Nine simple ways
675 to make it easier to (re) use your data. *Ideas in Ecology and Evolution*, 6, 1–10.
676 (doi:[10.4033/iee.2013.6b.6.f](https://doi.org/10.4033/iee.2013.6b.6.f))

677 Whitlock MC. 2011 Data archiving in ecology and evolution: best practices. *Trends Ecol. Evol.*
678 **26**, 61-65. (doi:[10.1016/j.tree.2010.11.006](https://doi.org/10.1016/j.tree.2010.11.006))

679 Yenni GM, Christensen EM, Bledsoe EK, Supp SR, Diaz RM, White EP, Ernest SM. 2019
680 Developing a modern data workflow for regularly updated data. *PLoS Biol.* **17**,
681 e3000125. (doi:[10.1371/journal.pbio.3000125](https://doi.org/10.1371/journal.pbio.3000125))

682