

Nine simple ways to make it easier to (re)use your data

Ethan P. White, Elita Baldrige, Zachary T. Brym, Kenneth J. Locey, Daniel J. McGlinn, and Sarah R. Supp

Ethan P. White (ethan.white@usu.edu), Dept. of Biology and the Ecology Center, Utah State University, Logan, UT, USA, 84341

Elita Baldrige (elita.baldrige@usu.edu), Dept. of Biology and the Ecology Center, Utah State University, Logan, UT, USA, 84341

Zachary T. Brym (zachary.brym@usu.edu), Dept. of Biology and the Ecology Center, Utah State University, Logan, UT, USA, 84341

Kenneth J. Locey (kenneth.locey@usu.edu), Dept. of Biology, Utah State University, Logan, UT, USA, 84341

Daniel J. McGlinn (daniel.mcglinn@usu.edu), Dept. of Biology and the Ecology Center, Utah State University, Logan, UT, USA, 84341

Sarah R. Supp (sarah.supp@usu.edu), Dept. of Biology and the Ecology Center, Utah State University, Logan, UT, USA, 84341

Abstract

Sharing data is increasingly considered to be an important part of the scientific process. Making your data publicly available allows original results to be reproduced and new analyses to be conducted. While sharing your data is the first step in allowing reuse, it is also important that the data be easy to understand and use. We describe nine simple ways to make it easy to reuse the data that you share and also make it easier to work with it yourself. Our recommendations focus on making your data understandable, easy to analyze, and readily available to the wider community of scientists.

Introduction

Sharing data is increasingly recognized as an important component of the scientific process (Whitlock et al. 2010). The sharing of scientific data is beneficial because it allows replication of research results and reuse in meta-analyses and projects not originally intended by the data collectors (Parr and Cummings 2005, Poisot et al. 2013). In ecology and evolutionary biology, sharing occurs through a combination of formal data repositories like GenBank (<http://www.ncbi.nlm.nih.gov/genbank/>) and Dryad (<http://datadryad.org/>), and through individual and institutional websites.

While data sharing is increasingly common and straightforward, much of the shared data in ecology and evolutionary biology are not easily reused because they do not follow best practices in terms of data structure, metadata, and licensing (Jones et al. 2006). This makes it more difficult to work with existing data and therefore makes the data less useful than it could be (Jones et al. 2006, Reichman et al. 2011). Here we provide a list of 9 simple ways to make it easier to reuse the data that you share.

36 Our recommendations focus on making your data understandable, easy to work with, and available
37 to the wider community of scientists. They are designed to be simple and straightforward to
38 implement, and as such represent an introduction to good data practices rather than a comprehensive
39 treatment. We contextualize our recommendations with examples from ecology and evolutionary
40 biology, though many of the recommendations apply broadly across scientific disciplines. Following
41 these recommendations makes it easier for anyone to reuse your data including other members of
42 your lab and even yourself.

43 **1. Share your data**

44 The first and most important step in sharing your data is to share your data. The recommendations
45 below will help make your data more useful, but sharing it in any form is a big step forward. So,
46 why should you share your data?

47 Data sharing provides substantial benefits to the scientific community (Fienberg and Martin 1985)
48 and the researchers sharing the data. For the scientific community it allows 1) the results of existing
49 analyses to be reproduced and improved upon (Fienberg and Martin 1985, Poisot et al. 2013), 2)
50 data to be combined in meta-analyses to reach general conclusions (Fienberg and Martin 1985), 3)
51 new approaches to be applied to the data and new questions asked using it (Fienberg and Martin
52 1985), and 4) approaches to scientific inquiry that could not be considered without broad scale
53 data sharing (Hampton et al. 2013). As a result, data sharing is increasingly required by funding
54 agencies (Poisot et al. 2013, e.g., [NSF](#), [NIH](#), [NSERC](#), [FWF](#)), journals (Whitlock et al. 2010), and
55 potentially by law (e.g. [FASTR](#), [OSTP Policy](#)). For data collectors, data sharing provides credit
56 for publication of data products (Poisot et al. 2013) and can increase citation metrics (Piwowar et
57 al. 2007, Piwowar and Vision 2013). In addition, data that are well-documented and standardized
58 make future reuse easier for the original investigator.

59 Despite these potential benefits to the community, individual incentives have historically been
60 insufficient to encourage widespread data sharing. Reluctance to share data is largely due to
61 concerns about 1) competition for publications based on the shared data, 2) a lack of recognition
62 for sharing data, and 3) a perception that sharing data is technically difficult and time consuming
63 (Palmer et al. 2004, Parr and Cummings 2005, Hampton et al. 2013). However, changes in how
64 data is treated and shared have increasingly ameliorated these issues. First, many data sharing
65 initiatives allow for data embargoes or limitations on direct competition that allow authors to develop
66 their publications and thus avoid competition for deriving publications from the data. Second, as
67 mentioned above, datasets are now considered citable entities and data providers receive recognition
68 in the form of increased citation metrics and credit on CVs and grant applications (Piwowar et al.
69 2007, Piwowar and Vision 2013, Poisot et al. 2013). Finally, data archives have become increasingly
70 common and easy to use (Parr and Cummings 2005, Hampton et al. 2013), and in some cases
71 sharing data requires no more effort than uploading a file to a website. As a result, it is increasingly
72 beneficial to the individual researcher to share data.

73 2. Provide metadata

74 The first key to using data is understanding it. Metadata is information about the data including how
75 it was collected, what the units of measurement are, and descriptions of how to best use the data
76 (Michener and Jones 2012). Clear metadata makes it easier to figure out if a dataset is appropriate
77 for a project. It also makes data easier to use by both the original investigators and by other scientists
78 by making it easy to figure out how to work with the data. Without clear metadata, datasets can be
79 overlooked or go unused due to the difficulty of understanding the data (Fraser and Gluck 1999,
80 Zimmerman 2003). Undocumented data also becomes less useful over time as information about
81 the data is gradually lost (Michener et al. 1997).

82 Metadata can take several forms, including descriptive file and column names, a written description
83 of the data, images (i.e., maps, photographs), and specially structured information that can be read
84 by computers (i.e., machine readable metadata). Good metadata should provide the following
85 information (Michener et al. 1997, Zimmerman 2003, Strasser et al. 2012):

- 86 • The what, when, where, and how of data collection.
- 87 • How to find and access the data.
- 88 • Suggestions on the suitability of the data for answering specific questions.
- 89 • Warnings about known problems or inconsistencies in the data, e.g., general descriptions of
90 data limitations or a column in a table to indicate the quality of individual data points.
- 91 • Information to check that the data are properly imported, e.g., the number of rows and columns
92 in the dataset and the total sum of numerical columns.

93 Just like any other scientific publication, metadata should be logically organized, complete, and clear
94 enough to enable interpretation and use of the data (Zimmerman 2007). Specific metadata standards
95 exist (e.g., Ecological Metadata Language [EML](#), Directory Interchange Format [DIF](#), Darwin Core
96 [DWC](#) (Wieczorek et al. 2012), Dublin Core Metadata Initiative [DCMI](#), Federal Geographic Data
97 Committee [FGDC](#) (Reichman et al. 2011, Whitlock 2011, Michener and Jones 2012). These
98 standards are designed to provide consistency in metadata across different datasets and also to
99 allow computers to interpret the metadata automatically. This allows broader and more efficient
100 use of shared data because computers can be relied on to identify (and potentially combine) data
101 from many different datasets for synthetic analyses (Brunt et al. 2002, Jones et al. 2006). While
102 following these standards is valuable, the most important thing is having metadata regardless of the
103 specific form.

104 Writing good metadata does not necessarily require a lot of extra time. The easiest way to develop
105 metadata is to start describing your data during the planning and data collection stages. This will
106 help you stay organized, make it easier to work with your data after it has been collected, and make
107 eventual publication of the data easier. If you decide to take the extra step and follow metadata
108 standards, there are tools designed to make this easier including: [KNB Morpho](#), [USGS xtme](#), and
109 [FGDC workbook](#).

110 3. Provide an unprocessed form of the data

111 Often, the data used in scientific analyses are modified in some way from the original form in which
112 they were collected. Values are averaged, units are converted, or indices are calculated from direct
113 measurements or observations to address the focal research questions and to fix issues associated
114 with the raw data. However, the best way to process data depends on the question being asked
115 and corrections for common data limitations often change as better approaches are developed. It
116 can also be very difficult to combine data from multiple sources that have each been processed in
117 different ways. Therefore, to make your data as useful as possible it is best to share the data in as
118 raw a form as possible. That means providing your data in a form that is as close as possible to the
119 field measurements and observations from which your analysis started.

120 This is not to say that your data are best suited for analysis in the raw form, but providing it in the
121 raw form gives data users the most flexibility. Of course, your work to develop and process the
122 data is also very important and can be quite valuable for other scientists using your data. This is
123 particularly true when correcting data for common limitations. Providing both the raw and processed
124 forms of the data, and clearly explaining the differences between them in the metadata, is an easy
125 way to include the benefits of both data forms. An alternate approach is to share the unprocessed
126 data along with the code that process the data to the form you used for analysis. This allows other
127 scientists to assess and potentially modify the process by which you arrived at the values used in
128 your analysis.

129 4. Use standard data formats

130 Everyone has their own favorite tools for storing and analyzing data. To make it easy to use your
131 data it is best to store it in a standard format that can be used by many different kinds of software.
132 Good standard formats include the type of file, the overall structure of the data, and the specific
133 contents of the file.

134 Use standard file formats

135 You should use file formats that are readable by most software and, when possible, are non-
136 proprietary (Borer et al. 2009, Strasser et al. 2011, 2012). Certain kinds of data in ecology and
137 evolution have well established standard formats such as FASTA files for nucleotide or peptide
138 sequences (<http://zhanglab.ccmb.med.umich.edu/FASTA/>) and the Newick files for phylogenetic
139 trees (<http://evolution.genetics.washington.edu/phylip/newicktree.html>). Use these well-defined
140 formats when they exist, because that is what other scientists and most existing software will be
141 able to work with most easily.

142 Data that does not have a well-defined standard format is often stored in tables. To increase
143 reuseability, tabular data should be stored in a format that can be opened by any type of software,
144 i.e. text files. These text files use delimiters to indicate different columns. Commas are the most
145 commonly used delimiter (i.e., comma-delimited text files with the .csv extension). Tabs can also
146 be used as a delimiter, although problems can occur in displaying the data correctly when importing
147 data from one program to another. In contrast to plain text files, proprietary formats such as those

148 used by Microsoft Excel (e.g. .xls, .xlsx) can be difficult to load into other programs. In addition,
149 these types of files can become obsolete, eventually making it difficult to open the data files at all if
150 the newer versions of the software no longer support the original format (Borer et al. 2009, Strasser
151 et al. 2011, 2012).

152 When naming files you should use descriptive names so that it is easy to keep track of what data
153 they contain (Borer et al. 2009, Strasser et al. 2011, 2012). If there are multiple files in a dataset,
154 name them in a consistent manner to make it easier to automate working with them. You should
155 also avoid spaces in file names, which can cause problems for some software (Borer et al. 2009).
156 Spaces in file names can be avoided by using camel case (e.g., RainAvg) or by separating the words
157 with underscores (e.g., rain_avg).

158 Use standard table formats

159 Data tables are ubiquitous in ecology and evolution. Tabular data provides a great deal of flexibility
160 in how data can be structured. However, this flexibility also makes it easy to structure your data
161 in a way that is difficult to (re)use. We provide three simple recommendations to help ensure that
162 tabular data are properly structured to allow the data to be easily imported and analyzed by most
163 data management systems and common analysis software, such as R and Python.

- 164 • Each row should represent a single observation (i.e., record) and each column should represent
165 a single variable or type of measurement (i.e., field) (Borer et al. 2009, Strasser et al. 2011,
166 2012). This is the standard formatting for tables in the most commonly used database
167 management systems and analysis packages, and makes the data easy to work with in the
168 most general way.
- 169 • Every cell should contain only a single value (Strasser et al. 2012). For example, do not
170 include units in the cell with the values (Figure 1) or include multiple measurements in a
171 single cell, and break taxonomic information up into single components with one column
172 each for family, genus, species, subspecies, etc. Violating this rule makes it difficult to process
173 or analyze your data using standard tools, because there is no easy way for the software to
174 treat the items within a cell as separate pieces of information.
- 175 • There should only be one column for each type of information (Borer et al. 2009, Strasser et
176 al. 2011, 2012). The most common violation of this rule is cross-tab structured data (http://en.wikipedia.org/wiki/Cross_tabulation), where different columns contain measurements of
177 the same variable (e.g., in different sites, treatments, etc.; Figure 1).

178 While cross-tab data can be easier to read and may be appropriate for data collection, this format
179 makes it difficult to link the records with additional data (e.g., the location and environmental
180 conditions at a site) and it cannot be properly used by most common database management and
181 analysis tools (e.g., relational databases, dataframes in R and Python, etc.). If tabular data are
182 currently in a cross-tab structure, there are tools to help restructure the data including functions in
183 Excel, R (e.g., melt() function in the R package reshape; Wickham 2007), and Python (e.g., melt()
184 function in the Pandas Python module <http://pandas.pydata.org/>).

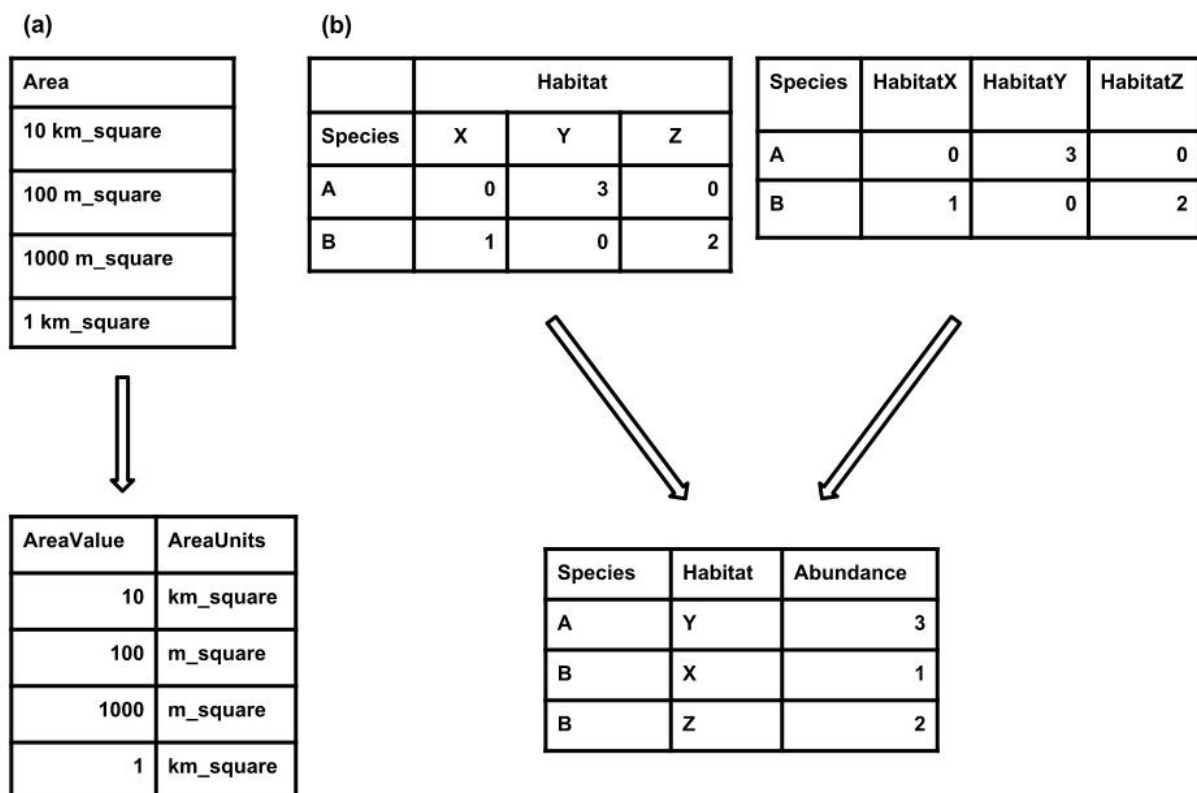


Figure 1: Examples of how to restructure two common issues with tabular data. (a) Each cell should only contain a single value. If more than one value is present then the data should be split into multiple columns. (b) There should be only one column for each type of information. If there are multiple columns then the column header should be stored in one column and the values from each column should be stored in a single column.

186 In addition to following these basic rules you should also make sure to use descriptive column
187 names (Borer et al. 2009). Descriptive column names make the data easier to understand and
188 therefore make data interpretation errors less likely. As with file names, spaces can cause problems
189 for some software and should be avoided.

190 **Use standard formats within cells**

191 In addition to using standard table structures it is also important to ensure that the contents of each
192 cell do not cause problems for data management and analysis software. Specifically, we recommend
193 that you:

- 194 • Be consistent. For example, be consistent in your capitalization of words, choice of delimiters,
195 and naming conventions for variables.
- 196 • Avoid special characters. Most software for storing and analyzing data works best on plain
197 text, and accents and other special characters can make it difficult to import your data (Borer
198 et al. 2009, Strasser et al. 2012).
- 199 • Avoid using your delimiter in the data itself (e.g., commas in the notes filed of a comma-
200 delimited file). This can make it difficult to import your data properly. This means that if you
201 are using commas as the decimal separator (as is often done in continental Europe) then you
202 should use a non-comma delimiter (e.g., a tab).
- 203 • When working with dates use the YYYY-MM-DD format (i.e., follow the [ISO 8601](#) data
204 standard).

205 While these standard approaches make it easier to use your data, the most important thing is to
206 document the approach that you have taken in your metadata (e.g., specify the date format) so that
207 data users can understand how to work with the data.

208 **5. Use good null values**

209 Most ecological and evolutionary datasets contain missing or empty data values. Working with this
210 kind of “null” data can be difficult, especially when the null values are indicated in problematic ways.
211 There are many ways to indicate a missing/empty value and little agreement on which approach to
212 use. We recommend choosing a null value that is both compatible with most software and unlikely
213 to cause errors in analyses (Table 1).

214 The null value that is most compatible with the software commonly used by biologists is the blank
215 (i.e., nothing; Table 1). Blanks are automatically treated as null values by R, Python, SQL, and
216 Excel. They are also easily spotted in a visual examination of the data. Note that a blank involves
217 entering nothing, it is not a space, so if you use this option make sure there are no hidden spaces.
218 There are two potential issues with blanks that should be considered:

- 219 1. It can be difficult to know if a value is missing or was overlooked during data entry.
- 220 2. Blanks can be confusing when spaces or tabs are used as delimiters in text files.

221 “NA” and “NULL” are reasonable null values, but they are only handled automatically by a subset of
 222 commonly used software (Table 1). “NA” can also be problematic if it is also used as an abbreviation
 223 (e.g., North America, Namibia, *Neotoma albigula*, sodium, etc.). We recommend against using
 224 numerical values to indicate nulls (e.g., 999, -999, etc.) because they typically require an extra step
 225 to remove from analyses and can be accidentally included in calculations. We also recommend
 226 against using non-standard text indications (e.g., No data, ND, missing, —) because they can cause
 227 issues with software that requires consistent data types within columns). Whichever null value
 228 that you use, only use one, use it consistently throughout the data set, and indicate it clearly in the
 229 metadata.

Null values	Problems	Compatibility	Recommendation
0	Indistinguishable from a true zero		Never use
blank	Hard to distinguish values that are missing from those overlooked on entry. Hard to distinguish blanks from spaces, which behave differently.	R, Python, SQL	Best option
999, -999	Not recognized as null by many programs without user input. Can be inadvertently entered into calculations.		Avoid
NA, na	Can also be an abbreviation (e.g., North America), can cause problems with data type (turn a numerical column into a text column). NA is more commonly recognized than na.	R	Good option
N/A	An alternate form of NA, but often not compatible with software		Avoid
NULL	Can cause problems with data type	SQL	Good option
None	Can cause problems with data type	Python	Avoid
No data	Can cause problems with data type, contains a space		Avoid

Missing	Can cause problems with data type	Avoid
-,+,. ,	Can cause problems with data type	Avoid

Table 1: Commonly used null values, limitations, compatibility with common software and a recommendation regarding whether or not it is a good option. Null values are indicated as being a null value for specific software if they work consistently and correctly with that software. For example, the null value “NULL” works correctly for certain applications in R, but does not work in others, so it is not presented in the table as R compatible.

6. Make it easy to combine your data with other datasets

Ecological and evolutionary data are often combined with other kinds of data. You can make it easier to combine your data with other data sources by including contextual data that appears across similar data sources. Two of the most common kinds of contextual data in ecology and evolution are taxonomy and geographic location. While this type of data is known and recorded in most studies (e.g., in field notebooks, on maps) it is frequently not included with the data. In general, if you have collected additional data or notes about a study organism or field site, there is a good chance that it will be useful to someone else, so including it with your data when you share it is a good idea. This kind of information can be included either as part of the data itself (e.g., in a new column or an additional table) or can be included in the metadata (e.g., the geographic location of the study site). For geographic data it is also important to include the datum (e.g., WGS-84) and sufficient precision (e.g., 4 decimal places if using decimal degrees) to allow the data to be combined with other geographic datasets.

When this data is included in a dataset it is often included as codes or abbreviations (e.g., DS instead of *Dipodomys spectabilis*, or site names instead of geographic coordinates). This can be useful for the data collector because it reduces data entry (e.g., typing a 1 into a plot column instead of entering both the latitude and longitude) and redundancy (e.g., a single column for a species ID rather than separate columns for family, genus, and species). However, without clear definitions these codes can be difficult to understand and make it more difficult to combine your data with external sources. One easy way to link your data to other datasets is to include additional tables that contain a column for the code and additional columns that describe the item in the standard way. For taxonomy, you might include a table with the species codes followed by their most current family, genus, and specific epithet. For site location, you could include a table with the site name or code followed by latitude and longitude, and other site information such as spatial extent, and temporal duration of sampling.

255 **7. Perform basic quality control**

256 Data, just like any other scientific product, should undergo some level of quality control (Reichman
257 et al. 2011). This is true regardless of whether you plan to share the data because quality control
258 will make it easier to analyze your own data and decrease the chance of making mistakes. However,
259 it is particularly important for data that will be shared because scientists using the data will not be
260 familiar with quirks in the data and how to work around them.

261 At its most basic, quality control can consist of a few quick sanity checks. More advanced quality
262 control can include automated checks on data as it is entered and double-entry of data (Lampe
263 and Weiler 1998, Michener and Jones 2012, Paulsen et al. 2012). This additional effort can be
264 time consuming but is valuable because it increases data accuracy by catching typographical errors,
265 reader/recorder error, out-of-range values, and questionable data in general (Lampe and Weiler
266 1998, Paulsen et al. 2012).

267 Before sharing your data we recommend performing a quick review. Start by performing a few
268 basic sanity checks. For example:

- 269 • If a column should contain numeric values, check that there are no non-numeric values in the
270 data.
- 271 • Check that empty cells actually represent missing data, and not mistakes in data entry, and
272 indicate that they are empty using the appropriate null values (see recommendation 6).
- 273 • Check for consistency in unit of measurement, data type (e.g., numeric, character), naming
274 scheme (e.g., taxonomy, location), etc.

275 These checks can be performed by carefully looking at the data or can be automated using common
276 programming and analysis tools like R or Python.

277 Then, ask someone else to look over your metadata and data and provide you with feedback about
278 anything they did not understand. In the same way that friendly reviews of papers can help catch
279 mistakes and identify confusing sections of papers, a friendly review of data can help identify
280 problems and things that are unclear in the data and metadata.

281 **8. Use an established repository**

282 For data sharing to be effective, data should be easy to find, accessible, and stored where it will
283 be preserved for a long time (Kowalczyk and Shankar 2011). To make your data (and associated
284 code) visible and easily accessible, and to ensure a permanent link to a well maintained website, we
285 suggest depositing your data in one of the major well-established repositories. This guarantees that
286 the data will be available in the same location for a long time, in contrast to personal and institutional
287 websites that do not guarantee long-term persistence. There are repositories available for sharing
288 almost any type of biological or environmental data. Repositories that host specific data types, such
289 as molecular sequences (e.g., DDBJ, GenBank, MG-RAST), are often highly standardized in data
290 type, format, and quality control. Other repositories host a wide array of data types and are less
291 standardized (e.g., Dryad, KNB, PANGAEA). In addition to the repositories focused on the natural
292 sciences there are also all-purpose repositories where data of any kind can be shared (e.g., figshare).

293 When choosing a repository you should consider where other researchers in your discipline are
 294 sharing their data. This helps to quickly identify the community's standard approach to sharing
 295 and increases the likelihood that other scientists will discover your data. In particular, if there is a
 296 centralized repository for a specific kind of data (e.g., GenBank for sequence data) then it should be
 297 used.

298 In cases where there is no *de facto* standard, it is worth considering differences among repositories
 299 in terms of use, data rights, and licensing (Table 3) and whether your funding agency or journal
 300 has explicit requirements or restrictions related to repositories. We also recommend that you use a
 301 repository that allows your dataset to be easily cited. Most repositories will describe how this works,
 302 but an easy way to guarantee that your data are citable is to confirm that the repository associates it
 303 with a persistent identifier, the most popular of which is the digital object identifier (DOI). DOIs are
 304 permanent unique identifiers that are independent of physical location and site ownership. There
 305 are also online tools for finding good repositories for your data including <http://databib.org> and
 306 <http://re3data.org>.

Repository	License	DOI	Metadata	Access	Notes
Dryad	CC0	Yes	Suggested	Open	Ecology & evolution data associated with publications
Ecological Archives	No	Yes	Required	Open	Publishes supplemental data for ESA journals and stand alone data papers
Knowledge Network for Biocomplexity	No	Yes	Required	Variable	Partners with ESA, NCEAS, DataONE
Paleobiology Database	Various CC	No	Optional	Variable	Paleontology specific
Data Basin	Various CC	No	Optional	Open	GIS data in ESRI files, limited free space
Pangaea	Various CC	Yes	Required	Variable	Editors participate in QA/QC
figshare	CC0	Yes	Optional	Open	Also allows deposition of other research outputs and private datasets

Table 2: Popular repositories for scientific datasets. This table does not include well-known molecular repositories (e.g. GenBank, EMBL, MG-RAST) that have become *de facto* standards in molecular and evolutionary biology. Consequently, several of these primarily serve the ecological community. These repositories are not exclusively used by members of specific institutions or museums, but accept data from the general scientific community.

9. Use an established and open license

Including an explicit license with your data is the best way to let others know exactly what they can and cannot do with the data you shared. Following the Panton Principles <http://pantonprinciples.org> we recommend:

1. Using well established licenses (or waivers) in order to clearly communicate the rights and responsibilities of both the people providing the data and the people using it.
2. Using the most open license (or waiver) possible, because even minor restrictions on data use can have unintended consequences for the reuse of the data (Schofield et al. 2009, Poisot et al. 2013).

The Creative Commons Zero (CC0) public domain dedication places no restrictions on data use and is considered by many to be one of the best ways to share data (e.g., (Schofield et al. 2009, Poisot et al. 2013), <http://blog.datadryad.org/2011/10/05/why-does-dryad-use-cc0/>). Several other licenses and waivers also accomplish these same goals <http://opendefinition.org/licenses/#Data>. Having a clear and open license (or waiver) will increase the chance that other scientists will be comfortable using your data.

Concluding remarks

Data sharing has the potential to transform the way we conduct ecological and evolutionary research (Fienberg and Martin 1985, Whitlock et al. 2010, Poisot et al. 2013). As a result, there are an increasing number of initiatives at the federal, funding agency, and journal levels to encourage or require the sharing of the data associated with scientific research (Piwowar and Chapman 2008, Whitlock et al. 2010, Poisot et al. 2013). However, making your data available is only the first step. To make data sharing as useful as possible it is necessary to make the data (re)usable with as little effort as possible (Jones et al. 2006, Reichman et al. 2011). This allows scientists to spend their time doing science rather than deciphering and cleaning up data.

We have provided a list of 9 practices that require only a small additional time investment but substantially improve the usability of data. These practices can be broken down into three major groups.

- 334 1. Well documented data are easier to understand.
- 335 2. Properly formatted data are easier to use in a variety of software.
- 336 3. Data that is shared in established repositories with open licenses is easier for others to find
- 337 and use.

338 Most of these recommendations are simply good practice for working with data regardless of
339 whether that data are shared or not. This means that following these recommendations (2-7) make
340 the data easier to work with for anyone, including you. This is particularly true when returning
341 to your own data for further analysis months or years after you originally collected or analyzed it.
342 In addition, data sharing often occurs within a lab or research group. Good data sharing practices
343 make these in-house collaborations faster, easier, and less dependent on lab members who may
344 have graduated or moved on to other endeavors. Following the other recommendations (1, 8, and 9)
345 provides broader benefits including academic credit in the form of published datasets and increased
346 citation metrics (Piwowar et al. 2007, Piwowar and Vision 2013, Poisot et al. 2013).

347 Many of these recommendations can be implemented at any point during a project, but the best
348 time to think about how to handle your data is before the project even starts (Michener and Jones
349 2012). A few hours of thought about how the data will be documented, structured, and shared at
350 the beginning of a project can prevent the need to restructure data or recall old information. This
351 will make it faster and easier to share your data when you are ready. By following these practices
352 we can assure that the data collected in ecology and evolution can be used to its full potential to
353 improve our understanding of biological systems.

354 Acknowledgments

355 Thanks to Karthik Ram for organizing this special section and inviting us to contribute. Carly
356 Strasser and Kara Woo recommended important references and David Harris and Carly Strasser
357 provided valuable feedback on null values, all via Twitter. Carl Boettiger, Matt Davis, Daniel
358 Hocking, Hilmar Lapp, Heinz Pampel, Karthik Ram, Thiago Silva, Carly Strasser, Tom Webb, and
359 @beroe (Twitter handle) provided valuable comments on the manuscript. Many of these comments
360 were part of the informal review process facilitated by posting this manuscript as a preprint to PeerJ
361 Preprints. The writing of this paper was supported by a CAREER grant from the U.S. National
362 Science Foundation (DEB 0953694) to EPW.

363 References

- 364 Borer, E. T., E. W. Seabloom, M. B. Jones, and M. Schildhauer. 2009. Some simple guidelines for
365 effective data management. *Bulletin of the Ecological Society of America* 90:205–214.
- 366 Brunt, J. W., P. McCartney, K. Baker, and S. G. Stafford. 2002. The future of ecoinformatics in
367 long term ecological research. Pages 14–18 *in* Proceedings of the 6th World Multiconference on
368 Systemics, Cybernetics and Informatics: SCI.
- 369 Fienberg, S. E., and M. E. Martin. 1985. Sharing research data. *Natl Academy Pr.*

- 370 Fraser, B., and M. Gluck. 1999. Usability of Geospatial Metadata or Space-Time Matters. *Bulletin*
371 *of the American Society for Information Science and Technology* 25:24–28.
- 372 Hampton, S. E., C. A. Strasser, J. J. Tewksbury, W. K. Gram, A. E. Budden, A. L. Batcheller, C.
373 S. Duke, and J. H. Porter. 2013. Big data and the future of ecology. *Frontiers in Ecology and the*
374 *Environment*.
- 375 Jones, M. B., M. P. Schildhauer, O. J. Reichman, and S. Bowers. 2006. The new bioinformatics:
376 integrating ecological data from the gene to the biosphere. *Annual Review of Ecology, Evolution,*
377 *and Systematics*:519–544.
- 378 Kowalczyk, S., and K. Shankar. 2011. Data sharing in the sciences. *Annual Review of Information*
379 *Science and Technology* 45:247–294.
- 380 Lampe, A. J., and J. M. Weiler. 1998. Data capture from the sponsors' and investigators' perspec-
381 tives: Balancing quality, speed, and cost. *Drug information journal* 32:871–886.
- 382 Michener, W. K., J. W. Brunt, J. J. Helly, T. B. Kirchner, and S. G. Stafford. 1997. Nongeospatial
383 metadata for the ecological sciences. *Ecological Applications* 7:330–342.
- 384 Michener, W. K., and M. B. Jones. 2012. Ecoinformatics: supporting ecology as a data-intensive
385 science.. *Trends in ecology & evolution* 27:85.
- 386 Palmer, M. A., E. S. Bernhardt, E. A. Chornesky, S. L. Collins, A. P. Dobson, C. S. Duke, B.
387 D. Gold, R. Jacobson, S. Kingsland, R. Kranz, M. J. Mappin, M. L. Martinez, F. Micheli, J. L.
388 Morse, M. L. Pace, M. Pascual, S. Palumbi, O. J. Reichman, A. Townsend, and M. G. Turner. 2004.
389 *Ecological Science and Sustainability for a Crowded Planet*.
- 390 Parr, C. S., and M. P. Cummings. 2005. Data sharing in ecology and evolution. *Trends in Ecology*
391 *and Evolution* 20:362–362.
- 392 Paulsen, A., S. Overgaard, and J. M. Lauritsen. 2012. Quality of Data Entry Using Single Entry,
393 Double Entry and Automated Forms Processing—An Example Based on a Study of Patient-Reported
394 Outcomes. *PloS one* 7:35087.
- 395 Piwowar, H. A., R. S. Day, and D. B. Fridsma. 2007. Sharing detailed research data is associated
396 with increased citation rate. *PLoS One* 2:308.
- 397 Piwowar, H. A., and W. W. Chapman. 2008. A review of journal policies for sharing research data.
398 *in* *ELPUB2008*.
- 399 Piwowar, H. A., and T. J. Vision. 2013. Data reuse and the open data citation advantage. *PeerJ*
400 *PrePrints* 1:1.
- 401 Poisot, T., R. Mounce, and D. Gravel. 2013. Moving toward a sustainable ecological science: don't
402 let data go to waste!
- 403 Reichman, O. J., M. B. Jones, and M. P. Schildhauer. 2011. Challenges and opportunities of open
404 data in ecology. *Science(Washington)* 331:703–705.
- 405 Schofield, P. N., T. Bubela, T. Weaver, L. Portilla, S. D. Brown, J. M. Hancock, D. Einhorn, G.
406 Tocchini-Valentini, M. H. de Angelis, and N. Rosenthal. 2009. Post-publication sharing of data and
407 tools. *Nature* 461:171–173.

- 408 Strasser, C. A., R. B. Cook, W. K. Michener, A. Budden, and R. Koskela. 2011. Promoting Data
409 Stewardship Through Best Practices. *in* Proceedings of the Environmental Information Management
410 Conference 2011 (EIM 2011). Oak Ridge National Laboratory (ORNL).
- 411 Strasser, C. A., R. Cook, W. K. Michener, and A. Budden. 2012. Primer on Data Management:
412 What you always wanted to know.
- 413 Whitlock, M. C. 2011. Data archiving in ecology and evolution: best practices. *Trends in ecology &*
414 *evolution* 26:61–65.
- 415 Whitlock, M. C., M. A. McPeck, M. D. Rausher, L. Rieseberg, and A. J. Moore. 2010. Data
416 archiving. *The American Naturalist* 175:145–146.
- 417 Wickham, H. 2007. Reshaping data with the reshape package. *Journal of Statistical Software* 21.
- 418 Wieczorek, J., D. Bloom, R. Guralnick, S. Blum, M. Döring, R. Giovanni, T. Robertson, and D.
419 Vieglais. 2012. Darwin Core: An evolving community-developed biodiversity data standard. *PloS*
420 *one* 7:29715.
- 421 Zimmerman, A. S. 2003. Data sharing and secondary use of scientific data: Experiences of
422 ecologists. The University of Michigan.
- 423 Zimmerman, A. S. 2007. Not by metadata alone: the use of diverse forms of knowledge to locate
424 data for reuse. *International Journal on Digital Libraries* 7:5–16.