

1 A Beginner's Guide to Conducting Reproducible  
2 Research

3 Jesse M. Alston<sup>1,2</sup> and Jessica A. Rick<sup>1,3</sup>

4 <sup>1</sup>*Program in Ecology, University of Wyoming, 1000 E University Dr., Laramie, WY 82072 USA*

5 <sup>2</sup>*Department of Zoology and Physiology, University of Wyoming, 1000 E University Dr., Laramie, WY*  
6 *82072 USA*

7 <sup>3</sup>*Department of Botany, University of Wyoming, 1000 E University Dr., Laramie, WY 82072 USA*

## 8 **Abstract**

9 Reproducible research is widely acknowledged as an important tool for improving science and  
10 reducing harm from the “replication crisis”, yet research in ecology and evolutionary biology  
11 remains largely irreproducible. In this article, we make the case for why all research should be  
12 reproducible, explain why research is often not reproducible, and offer a simple framework that  
13 researchers can use to make their research more reproducible. Researchers can increase the  
14 reproducibility of their work by improving data management practices, writing more readable  
15 code, and increasing use of the many available platforms for sharing data and code. While  
16 reproducible research is often associated with a set of advanced tools for sharing data and code,  
17 reproducibility is just as much about maintaining work habits that are already widely  
18 acknowledged as best practices for research. Increasing reproducibility will increase rigor,  
19 trustworthiness, and transparency while benefiting both practitioners of reproducible research and  
20 their fellow researchers.

21 *Key words: data management, data repository, software, open science, replication*

## 22 **Introduction**

23 Replication is a fundamental tenet of science, but there is increasing fear among scientists that too  
24 few scientific studies can be replicated. This has been termed the “replication crisis” (Ioannidis,  
25 2005; Schooler, 2014). Scientific papers often include inadequate detail to enable replication  
26 (Haddaway and Verhoeven, 2015; Archmiller et al., 2020), many attempted replications of  
27 well-known scientific studies have failed in a wide variety of disciplines (Bohannon, 2015;  
28 Hewitt, 2012; Moonesinghe et al., 2007; Open Science Collaboration, 2015), and rates of paper  
29 retractions are increasing (Cokol et al., 2008; Steen et al., 2013). Because of this, researchers are

30 working to develop new ways for researchers, research institutions, research funders, and journals  
31 to overcome this problem (Peng, 2011; Sandve et al., 2013; Stodden et al., 2013; Fiedler et al.,  
32 2012).

33       Because replicating studies with new independent data is expensive, rarely published in  
34 high-impact journals, and sometimes even methodologically impossible, computationally  
35 reproducible research (most often termed simply “reproducible research”) is often suggested as a  
36 pathway for increasing our ability to assess the validity and rigor of scientific results (Peng,  
37 2011). Research is reproducible when others can reproduce the results of a scientific study given  
38 only the original data, code, and documentation (Essawy et al., 2020). This approach focuses on  
39 the research process after data collection is complete, and it has many (though not all) of the  
40 advantages of replicating studies with independent data while minimizing the largest barrier (i.e.,  
41 the financial and time costs of collecting new data). Replicating studies remains the gold standard  
42 for rigorous scientific research, but reproducibility is increasingly viewed as a minimum standard  
43 that all scientists should strive toward (Peng, 2011; Sandve et al., 2013; Archmiller et al., 2020;  
44 Culina et al., 2020).

45       This commentary describes basic requirements for such reproducible research in the fields of  
46 ecology and evolutionary biology. In it, we make the case for why all research should be  
47 reproducible, explain why research is often not reproducible, and present a simple three-part  
48 framework all researchers can use to make their research more reproducible. These principles are  
49 applicable to researchers working in all sub-disciplines within ecology and evolutionary biology  
50 with data sets of all sizes and levels of complexity.

## 51 **Why Do Reproducible Research?**

### 52 **Reproducible research benefits those who do it**

53 Reproducible research is a by-product of careful attention to detail throughout the research  
54 process, and allows researchers to ensure that they can repeat the same analysis multiple times  
55 with the same results, at any point in that process. Because of this, researchers who conduct  
56 reproducible research are the primary beneficiaries of this practice.

57 First, reproducible research helps researchers remember how and why they performed  
58 specific analyses during the course of a project. This enables easier explanation of work to  
59 collaborators, supervisors, and reviewers, and it allows collaborators to conduct supplementary  
60 analyses more quickly and more efficiently.

61 Second, reproducible research enables researchers to quickly and simply modify analyses  
62 and figures. This is often requested by supervisors, collaborators, and reviewers across all stages  
63 of a research project, and expediting this process saves substantial amounts of time. When  
64 analyses are reproducible, creating a new figure may be as easy as changing one value in a line of  
65 code and re-running a script, rather than spending hours recreating a figure from scratch.

66 Third, reproducible research enables quick reconfiguration of previously conducted research  
67 tasks so that new projects that require similar tasks become much simpler and easier. Science is an  
68 iterative process, and many of the same tasks are performed over and over. Conducting research  
69 reproducibly enables researchers to re-use earlier materials (e.g., analysis code, file organization  
70 systems) to execute these common research tasks more efficiently in subsequent iterations.

71 Fourth, conducting reproducible research is a strong indicator to fellow researchers of rigor,  
72 trustworthiness, and transparency in scientific research. This can increase the quality and speed of  
73 peer review, because reviewers can directly access the analytical process described in a

74 manuscript. Peer reviewers' work becomes easier and they may be able to answer methodological  
75 questions without asking the authors. Reviewers can check whether code matches with methods  
76 described in the text of a manuscript to make sure that authors correctly performed the analyses as  
77 described, and it increases the probability that errors are caught during the peer-review process,  
78 decreasing the likelihood of corrections or retractions after publication. Finally, it also protects  
79 researchers from accusations of research misconduct due to analytical errors, because it is  
80 unlikely that researchers would openly share fraudulent code and data with the rest of the research  
81 community.

82 Finally, reproducible research increases paper citation rates (Piwowar et al., 2007;  
83 McKiernan et al., 2016) and allows other researchers to cite code and data in addition to  
84 publications. This enables a given research project to have more impact than it would if the data  
85 or methods were hidden from the public. For example, researchers can re-use code from a paper  
86 with similar methods and organize their data in the same manner as the original paper, then cite  
87 code from the original paper in their manuscript. A third team of researchers may conduct a  
88 meta-analysis on the phenomenon described in these two research papers, and thus use and cite  
89 both of these papers and the data from those papers in their meta-analysis. Papers are more likely  
90 to be cited in these re-use cases if full information about data and analyses are available  
91 (Whitlock, 2011; Culina et al., 2018).

## 92 **Reproducible research benefits the research community**

93 Reproducible research also benefits others in the scientific community. Sharing data, code, and  
94 detailed research methods and results leads to faster progress in methodological development and  
95 innovation because research is more accessible to more scientists (Mislán et al., 2016; Parr and  
96 Cummings, 2005; Roche et al., 2015).

97 First, reproducible research allows others to learn from your work. Scientific research has a  
98 steep learning curve, and allowing others to access data and code gives them a head start on  
99 performing similar analyses. For example, researchers who are new to an analytical technique can  
100 use code shared with the research community by researchers with more experience with that  
101 technique to learn how to rigorously perform and validate these analyses. This allows researchers  
102 to conduct research that is more rigorous from the outset, rather than having to spend months or  
103 years trying to figure out current “best practices” through trial and error. Modifying existing  
104 resources can also save time and effort for experienced researchers—even experienced coders can  
105 modify existing code much faster than they can write code from scratch. Sharing code thus allows  
106 experienced researchers to perform similar analyses more quickly.

107 Second, reproducible research allows others to understand and reproduce a researcher’s  
108 work. Allowing others to access data and code makes it easier for other scientists to perform  
109 follow-up studies to increase the strength of evidence for the phenomenon of interest. It also  
110 increases the likelihood that similar studies are compatible with one another, and that a group of  
111 studies can together provide evidence in support of or in opposition to a concept. In addition,  
112 sharing data and code increases the utility of these studies for meta-analyses that are important for  
113 generalizing and contextualizing the findings of studies on a topic. Meta-analyses in ecology and  
114 evolutionary biology are often hindered by incompatibility of data between studies, or lack of  
115 documentation for how those data were obtained (Stewart, 2010; Culina et al., 2018).  
116 Well-documented, reproducible findings enhance the likelihood that data can be used in future  
117 meta-analyses (Gerstner et al., 2017).

118 Third, reproducible research allows others to protect themselves from your mistakes.  
119 Mistakes happen in science. Allowing others to access data and code gives them a better chance  
120 to critically analyze the work, which can lead to coauthors or reviewers discovering mistakes

121 during the revision process, or other scientists discovering mistakes after publication. This  
122 prevents mistakes from compounding over time and provides protection for collaborators,  
123 research institutions, funding organizations, journals, and others who may be affected when such  
124 mistakes happen.

## 125 **Barriers to Reproducible Research**

126 There are a number of reasons that most research is not reproducible. Rapidly developing  
127 technologies and analytical tools, novel interdisciplinary approaches, unique ecological study  
128 systems, and increasingly complex data sets and research questions hinder reproducibility, as does  
129 pressure on scientists to publish novel research quickly. This multitude of barriers can be  
130 simplified into four primary themes: (1) complexity, (2) technological change, (3) human error,  
131 and (4) concerns over intellectual property rights. Each of these concerns can contribute to  
132 making research less reproducible and can be valid in some scenarios. However, each of these  
133 factors can also be addressed easily via well-developed tools, protocols, and institutional norms  
134 concerning reproducible research.

135 **Complexity.** — Science is difficult, and scientific research requires specialized (and often  
136 proprietary) knowledge and tools that may not be available to everyone who would like to  
137 reproduce research. For example, studies in the fields of ecology and evolutionary biology often  
138 involve study systems, mathematical models, and statistical techniques that require a large  
139 amount of domain knowledge to understand, and these analyses can therefore be difficult to  
140 reproduce for those with limited understanding of any of the necessary underlying bases of  
141 knowledge. Some analyses may require high-performance computing clusters that use several  
142 different programming languages and software packages, or that are designed for specific

143 hardware configurations. Other analyses may be performed using proprietary software programs  
144 such as SAS statistical software (SAS Institute Inc., Cary, NC, USA) or ArcGIS (Esri, Redlands,  
145 CA, USA) that require expensive software licenses. Lack of knowledge, lack of institutional  
146 infrastructure, and lack of funding all make research less reproducible. However, most of these  
147 issues can be mitigated fairly easily. Researchers can cite primers on complex subjects or  
148 analyses to reduce knowledge barriers. They can also thoroughly annotate analytical code with  
149 comments explaining each step in an analysis, or provide extensive documentation on research  
150 software. Using open software (when possible) makes research more accessible for other  
151 researchers as well.

152 **Technological change.** — Hardware and software used in analyzing data both change over  
153 time, and they often change quickly. When old tools become obsolete, research becomes less  
154 reproducible. For example, reproducing research performed in 1960 using that era's  
155 computational tools would require a completely new set of tools today. Even research performed  
156 just a few years ago may have been conducted using software that is no longer available or is  
157 incompatible with other software that has since been updated. One minor update in a piece of  
158 software used in one minor analysis in an analytical workflow can render an entire project less  
159 reproducible. However, this too can be mitigated by using established tools in reproducible  
160 research. Careful documentation of versions of software used in analyses is a baseline  
161 requirement that anyone can meet. There are also more advanced tools that can help overcome  
162 such challenges in making research reproducible, including software containers, which are  
163 described in further detail below.

164 **Human error.** — Though fraudulent research is often cited as reason to make research more  
165 reproducible (e.g., Ioannidis 2005; Laine et al. 2007; Crocker and Cooper 2011), many more  
166 innocent reasons exist as to why research is often difficult to reproduce (e.g., Elliott 2014). People



167 forget small details of how they performed analyses. They fail to describe data collection  
168 protocols or analyses completely despite their best efforts and multiple reviewers checking their  
169 work. They fail to collect or thoroughly document data that seem unimportant during collection  
170 but later turn out to be vital for unforeseen reasons. Science is performed by fallible humans, and  
171 a wide variety of common events can render research less reproducible.

172 While not all of these challenges can be avoided by performing research reproducibly, a  
173 well-documented research process can guard against small errors and sloppy analyses. For  
174 example, carefully recording details such as when and where data were collected, what decisions  
175 were made during data collection, and what labeling conventions were used can make a huge  
176 difference in making sure that those data can later be used appropriately or re-purposed.  
177 Unintentional errors often occur during the data wrangling stage of a project, and these can be  
178 mitigated by keeping multiple copies of data to prevent data loss, carefully documenting the  
179 process for converting raw data into clean data, and double-checking a small test set of data  
180 before manipulating the data set as a whole.

181 **Intellectual property rights.** — Researchers often hesitate to share data and code because  
182 doing so may allow other researchers to use data and code incorrectly or unethically. Other  
183 researchers may use publicly available data without notifying authors, leading to incorrect  
184 assumptions about the data that result in invalid analyses. Researchers may use publicly available  
185 data or code without citing the original data owners or code writers, who then do not receive  
186 proper credit for gathering expensive data or writing time-consuming code. Researchers may  
187 want to conceal data from others so that they can perform new analyses on those data in the future  
188 without worrying about others scooping them using the shared data. Rational self-interest can  
189 lead to hesitation to share data and code via many pathways, and we acknowledge that making  
190 data openly available is likely the most controversial aspect of reproducible research (e.g., Cassey

191 and Blackburn 2006; Hampton et al. 2013; Mills et al. 2015; Whitlock et al. 2016; Mills et al.  
192 2016). However, new tools for sharing data and code (outlined below and in Table 1) are making  
193 it easier for researchers to receive credit for doing so and to prevent others from using their data  
194 during an embargo period.

## 195 **A Three-Step Framework for Conducting Reproducible**

### 196 **Research**

197 Conducting reproducible research is not exceedingly difficult, nor does it require encyclopedic  
198 knowledge of esoteric research tools and protocols. Whether they know it or not, most researchers  
199 already perform much of the work required to make research reproducible. To clarify this point,  
200 we outline below some basic steps toward making research more reproducible in three stages of a  
201 research project: (1) before data analysis, (2) during analysis, and (3) after analysis. We discuss  
202 practical tips that anyone can use, as well as more advanced tools for those who would like to  
203 move beyond basic requirements (Table 1). Most readers will recognize that reproducible  
204 research largely consists of widely accepted best practices for scientific research, and that striving  
205 to meet a reasonable benchmark of reproducibility is both more valuable and more attainable than  
206 researchers may think.

### 207 **Before data analysis: data storage and organization**

208 Reproducibility starts in the planning stage, with sound data management practices. It does not  
209 arise simply from sharing data and code online after a project is done. It is difficult to reproduce  
210 research when data are disorganized or missing, or when it is impossible to determine where or

211 how data originated.

212 First, data should be backed up at every stage of the research process and stored in multiple  
213 locations. This includes raw data (e.g., physical data sheets or initial spreadsheets), clean  
214 analysis-ready data (i.e., final data sets), and steps in between. Because it is entirely possible that  
215 researchers unintentionally alter or corrupt data while cleaning it up, raw data should always be  
216 kept as a back up. It is good practice to scan and save data sheets or lab notebook pages  
217 associated with a data set to ensure that these are kept paired with the digital data set. Ideally,  
218 different copies should be stored in different locations and using different storage media (e.g.,  
219 paper copies *and* an external hard drive *and* cloud storage) to minimize risk of data loss from any  
220 single cause. Computers crash, hard drives are misplaced and stolen, and servers are  
221 hacked—researchers should not leave themselves vulnerable to those events.

222 Digital data files should be stored in useful, flexible, portable, non-proprietary formats.  
223 Storing data digitally in a “flat” file format is almost always a good idea. Flat file formats are  
224 those that store data as plain text with one record per line (e.g., .csv or .txt files) and are the  
225 most portable formats across platforms, as they can be opened by anyone without proprietary  
226 software programs. For more complex data types, multi-dimensional relational formats such as  
227 json, hdf5, or other discipline-specific formats (e.g., biom and EML) may be appropriate.  
228 However, the complexity of these formats makes them difficult for many researchers to access  
229 and use appropriately, so it is best to stick with simpler file formats when possible.

230 It is often useful to transform data into a ‘tidy’ format (Wickham, 2014) when cleaning up  
231 and standardizing raw data. Tidy data are in long format (i.e., variables in columns, observations  
232 in rows), have consistent data structure (e.g., character data are not mixed with numeric data for a  
233 single variable), and have informative and appropriately formatted headers (e.g., reasonably short  
234 variable names that do not include problematic characters like spaces, commas, and parentheses).

235 Data in this format are easy to manipulate, model, and visualize during analysis.

236 Metadata explaining what was done to clean up the data and what each of the variables  
237 means should be stored along with the data. Data are useless unless they can be interpreted  
238 (Roche et al., 2015); metadata is how we maximize data interpretability across potential users. At  
239 a minimum, all data sets should include informative metadata that explains how and why data  
240 were collected, what variable names mean, whether a variable consists of raw or transformed  
241 data, and how observations are coded. Metadata should be placed in a sensible location that pairs  
242 it with the data set it describes. A few rows of metadata above a table of observations within the  
243 same file may work in some cases, or a paired text file can be included in the same directory as  
244 the data if the metadata must be more detailed. In the latter case, it is best to stick with a simple  
245 `.txt` file for metadata to maximize portability.

246 Finally, researchers should organize files in a sensible, user-friendly structure and make sure  
247 that all files have informative names. It should be easy to tell what is in a file or directory from its  
248 name, and a consistent naming protocol (e.g., ending the filename with the date created or version  
249 number) provides even more information when searching through files in a directory. A consistent  
250 naming protocol for both directories and files also makes coding simpler by placing data,  
251 analyses, and products in logical locations with logical names. It is often more useful to organize  
252 files in small blocks of similar files, rather than having one large directory full of hundreds of  
253 files. For example, Noble (2009) suggests organizing computational projects within a main  
254 directory for each project, with sub-directories for the manuscript (`doc/`), data files (`data/`),  
255 analyses (`scripts/` or `src/`), and analysis products (`results/`) within that directory. While this  
256 specific organization scheme may differ for other types of research, keeping all of the research  
257 products and documentation for a given project organized in this way makes it much easier to find  
258 everything at all stages of the research process, and to archive it or share it with others once the

259 project is finished.

260 Throughout the research process, from data acquisition to publication, version control can be  
261 used to record a project's history and provide a log of changes that have occurred over the life of a  
262 project or research group. Version control systems record changes to a file or set of files over time  
263 so that you can recall specific versions later, compare differences between versions of files, and  
264 even revert files back to previous states in the event of mistakes. Many researchers use version  
265 control systems to track changes in code and documents over time. The most popular version  
266 control system is Git, which is often used via hosting services such as GitHub, GitLab, and  
267 BitBucket (Table 1). These systems are relatively easy to set up and use, and they systematically  
268 store snapshots of data, code, and accompanying files throughout the duration of a project.  
269 Version control also enables a specific snapshot of data or code to be easily shared, so that code  
270 used for analyses at a specific point in time (e.g., when a manuscript is submitted) can be  
271 documented, even if that code is later updated.

## 272 **During analysis: best coding practices**

273 When possible, all data wrangling and analysis should be performed using coding scripts—as  
274 opposed to using interactive or point-and-click tools—so that every step is documented and  
275 repeatable by yourself and others. Code both performs operations on data and serves as a log of  
276 analytical activities. Because of this second function, code (unlike point-and-click programs) is  
277 inherently reproducible. Most errors are unintentional mistakes made during data wrangling or  
278 analysis, so having a record of these steps ensures that analyses can be checked for errors and are  
279 repeatable on future data sets. If operations are not possible to script, then they should be  
280 well-documented in a log file that is kept in the appropriate directory.

281 Analytical code should be thoroughly annotated with comments. Comments embedded

282 within code serve as metadata for that code, substantially increasing its usefulness. Comments  
283 should contain enough information for an informed stranger to easily understand what the code  
284 does, but not so much that sorting through comments is a chore. Code comments can be tested for  
285 this balance by a friend who is knowledgeable about the general area of research but is not a  
286 project collaborator. In most scripting languages, the first few lines of a script should include a  
287 description of what the script does and who wrote it, followed by small blocks that import data,  
288 packages, and external functions. Data cleaning and analytical code then follows those sections,  
289 and sections are demarcated using a consistent protocol and sufficient comments to explain what  
290 function each section of code performs.

291       Following a clean, consistent coding style makes code easier to read. Many well-known  
292 organizations (e.g., RStudio, Google) offer style guidelines for software code that were developed  
293 by many expert coders. Researchers should take advantage of these while keeping in mind that all  
294 style guides are subjective to some extent. Researchers should work to develop a style that works  
295 for them. This includes using a consistent naming convention (e.g., camelCase or snake\_case)  
296 to name objects and embedding meaningful information in object names (e.g., using “\_mat” as a  
297 suffix for objects to denote matrices or “\_df” to denote data frames). Code should also be written  
298 in relatively short lines and grouped into blocks, as our brains process narrow columns of data  
299 more easily than longer ones (Martin, 2009). Blocks of code also keep related tasks together and  
300 can function like paragraphs to make code more comprehensible.

301       There are several ways to prevent coding mistakes and make code easier to use. First,  
302 researchers should automate repetitive tasks. For example, if a set of analysis steps are being used  
303 repeatedly, those steps can be saved as a function and loaded at the top of the script. This reduces  
304 the size of a script and eliminates the possibility of accidentally altering some part of a function  
305 so that it works differently in different locations within a script. Similarly, researchers can use

306 loops to make code more efficient by performing the same task on multiple values or objects in  
307 series (though it is also important to note that nesting too many loops inside one another can  
308 quickly make code incomprehensible). A third way to reduce mistakes is to reduce the number of  
309 hard-coded values that must be changed to replicate analyses on an updated or new data set. It is  
310 often best to read in the data file(s) and assign parameter values at the beginning of a script, so  
311 that those variables can then be used throughout the rest of the script. When operating on new  
312 data, these variables can then be changed once at the beginning of a script rather than multiple  
313 times in locations littered throughout the script.

314       Because incompatibility between operating systems or program versions can inhibit the  
315 reproducibility of research, the current gold standard for ensuring that analyses can be used in the  
316 future is to create a software container, such as a Docker (Merkel, 2014) or Singularity  
317 (Kurtzer et al., 2017) image (Table 1). Containers are standalone, portable environments that  
318 contain the entire computing environment used in an analysis: software, all of its dependencies,  
319 libraries, binaries, and configuration files, all bundled into one package. Containers can then be  
320 archived or shared, allowing them to be used in the future, even as packages, functions, or  
321 libraries change over time. If creating a software container is infeasible or a larger step than  
322 researchers are willing to take, it is important to thoroughly report all software packages used,  
323 including version numbers.

## 324 **After data analysis: finalizing results and sharing**

325 After the steps above have been followed, it is time for the step most people associate with  
326 reproducible research: sharing research with others. As should be clear by now, sharing the data  
327 and code is far from the only component of reproducible research; however, once Steps 1 and 2  
328 above are followed, it becomes the easiest step. All input data, scripts, program versions,

329 parameters, and important intermediate results should be made publicly and easily accessible.  
330 Various solutions are now available to make data sharing convenient, standardized, and accessible  
331 in a variety of research areas. There are many ways to do this, several of which are described  
332 below.

333 Just as it is better to use scripts than interactive tools in analysis, it is better to produce tables  
334 and figures directly from code than to manipulate these using Adobe Illustrator, Microsoft  
335 Powerpoint, or other image editing programs. A large number of errors in finished manuscripts  
336 come from not remembering to change *all* relevant numbers or figures when a part of an analysis  
337 changes, and this task can be incredibly time-consuming when revising a manuscript. Truly  
338 reproducible figures and tables are created directly with code and integrated into documents in a  
339 way that allows automatic updating when analyses are re-run, creating a “dynamic” document.  
340 For example, documents written in  $\text{\LaTeX}$  and markdown incorporate figures directly from a  
341 directory, so a figure will be updated in the document when the figure is updated in the directory  
342 (see Xie 2015 for a much lengthier discussion of dynamic documents). Both  $\text{\LaTeX}$  and markdown  
343 can also be used to create presentations that can incorporate live-updated figures when code or  
344 data change, so that presentations can be reproducible as well. If using one of these tools is too  
345 large a leap, then simply producing figures directly from code—instead of adding annotations and  
346 arranging panels post-hoc—can make a substantial difference in increasing the reproducibility of  
347 these products.

348 Beyond creating dynamic documents, it is possible to make data wrangling, analysis, and  
349 creation of figures, tables, and manuscripts a “one-button” process using GNU Make  
350 (<https://www.gnu.org/software/make/>). GNU Make is a simple, yet powerful tool that can be used  
351 to coordinate and automate command-line processes, such as a series of independent scripts. For  
352 example, a Makefile can be written that will take the input data, clean and manipulate it, analyze



353 it, produce figures and tables with results, and update a `LATEX` or `markdown` manuscript document  
354 with those figures, tables, and any numbers included in the results. Setting up research projects to  
355 run in this way takes some time, but it can substantially expedite re-analyses and reduce  
356 copy-paste errors in manuscripts.

357       Currently, code and data that can be used to replicate research are often found in the  
358 supplementary material of journal articles. Some journals (e.g., *eLife*) are even experimenting  
359 with embedding data and code in articles themselves. However, this is not a fail-safe method of  
360 archiving data and analyses: supplementary materials can be lost if a journal switches publishers  
361 or when a publisher changes its website. In addition, research is only reproducible if it can be  
362 accessed, and many papers are published in journals that are locked behind paywalls that make  
363 them inaccessible to many researchers (Desjardins-Proulx et al., 2013; McKiernan et al., 2016;  
364 Alston, 2019). To increase access to publications, authors can post pre-prints of final (but  
365 pre-acceptance) versions of manuscripts on a pre-print server, or post-prints of manuscripts on  
366 post-print servers. There are several widely used pre-print servers (see Table 1 for three  
367 examples), and libraries at many research institutions host post-print servers.

368       Similarly, data and code shared on personal websites are only available as long as websites  
369 are maintained, and can be difficult to transfer when researchers migrate to another domain or  
370 website provider. Materials archived on personal websites are also often difficult for other  
371 scientists to find, as they are not usually linked to the published research and lack a permanent  
372 digital object identifier (DOI). To make research accessible to everyone, it is therefore better to  
373 use tools like data and code repositories than personal websites.

374       Data archiving in online repositories has become more popular in recent years, a trend  
375 resulting from a combination of improvements in technology for sharing data, an increase in  
376 -omics-scale data sets, and an increasing number of publisher and funding organizations who

377 encourage or mandate data archiving (Whitlock et al., 2010; Whitlock, 2011; Nosek et al., 2015).  
378 Data repositories are large databases that collect, manage, and store data sets for analysis, sharing,  
379 and reporting. Repositories may be either subject- or data-specific, or cross-disciplinary general  
380 repositories that accept multiple data types. Some are free and others require a fee for depositing  
381 data. Journals often recommend appropriate repositories on their websites, and these  
382 recommendations should be consulted when submitting a manuscript. Three commonly used  
383 general purpose repositories are Dryad, Zenodo, and Figshare; each of these creates a DOI that  
384 allows data and code to be citable by others. Before choosing a repository, researchers should  
385 explore commonly used options in their specific fields of research.

386       When data, code, software, and products of a research project are archived together, these  
387 are termed a “research compendium” (Gentleman and Lang, 2007). Research compendia are  
388 increasingly common, although standards for what is included in research compendia differ  
389 between scientific fields. They provide a standardized and easily recognisable way to organize the  
390 digital materials of a research project, which enables other researchers to inspect, reproduce, and  
391 extend research (Marwick et al., 2018).

392       In particular, the Open Science Framework (OSF; <http://osf.io/>) is a project management  
393 repository that goes beyond the repository features of Dryad, Zenodo, and Figshare to integrate  
394 and share components of a research project using collaborative tools. The goal of the OSF is to  
395 enable research to be shared at every step of the scientific process—from developing a research  
396 idea and designing a study, to storing and analyzing collected data and writing and publishing  
397 reports or papers (Sullivan et al., 2019). OSF is integrated with many other reproducible research  
398 tools, including widely used pre-print servers, version control software, and publishers.

## 399 **Conclusions**

400 While many researchers associate reproducible research primarily with a set of advanced tools for  
401 sharing research, reproducibility is just as much about simple work habits as the tools used to  
402 share data and code. We ourselves are not perfect reproducible researchers—we do not use all the  
403 tools mentioned in this commentary all the time and often fail to follow our own advice (almost  
404 always to our regret). Nevertheless, we recognize that reproducible research is a process rather  
405 than a destination and work hard to consistently increase the reproducibility of our work. We  
406 encourage others to do the same. Researchers can make strides toward a more reproducible  
407 research process by simply thinking carefully about data management and organization, coding  
408 practices, and processes for making figures and tables (e.g., Fig. 1). Time and expertise must be  
409 invested in learning and adopting these tools and tips, and this investment can be substantial.  
410 Nevertheless, we encourage our fellow researchers to work toward more open and reproducible  
411 research practices so we can all enjoy the resulting improvements in work habits, collaboration,  
412 scientific rigor, and trust in science.

## 413 **Acknowledgements**

414 Many thanks to J.G. Harrison, B.J. Rick, A.L. Lewanski, E.A. Johnson, and F.S. Dobson for  
415 providing helpful comments on pre-publication versions of this manuscript, and to C.A. Buerkle  
416 for inspiring this project during his Computational Biology course at the University of Wyoming.

## 417 **References**

- 418 Alston, J. M. (2019). Open access principles and practices benefit conservation. *Conservation*  
419 *Letters*, 12(6):e12672.
- 420 Archmiller, A. A., Johnson, A. D., Nolan, J., Edwards, M., Elliott, L. H., Ferguson, J. M.,  
421 Iannarilli, F., Vélez, J., Vitense, K., Johnson, D. H., and Fieberg, J. (2020). Computational  
422 reproducibility in The Wildlife Society’s flagship journals. *Journal of Wildlife Management*,  
423 84(5):1012–1017.
- 424 Bohannon, J. (2015). Many psychology papers fail replication test. *Science*, 349(6251):910–911.
- 425 Cassey, P. and Blackburn, T. M. (2006). Reproducibility and repeatability in ecology. *BioScience*,  
426 56(12):958–959. Publisher: Oxford Academic.
- 427 Cokol, M., Ozbay, F., and Rodriguez-Esteban, R. (2008). Retraction rates are on the rise. *EMBO*  
428 *reports*, 9(1):2–2.
- 429 Crocker, J. and Cooper, M. L. (2011). Addressing scientific fraud. *Science*,  
430 334(6060):1182–1182.
- 431 Culina, A., Berg, I. v. d., Evans, S., and Sánchez-Tójar, A. (2020). Low availability of code in  
432 ecology: a call for urgent action. *PLOS Biology*, 18(7):e3000763. Publisher: Public Library of  
433 Science.
- 434 Culina, A., Crowther, T. W., Ramakers, J. J. C., Gienapp, P., and Visser, M. E. (2018). How to do  
435 meta-analysis of open datasets. *Nature Ecology & Evolution*, 2(7):1053–1056.
- 436 Desjardins-Proulx, P., White, E. P., Adamson, J. J., Ram, K., Poisot, T., and Gravel, D. (2013).  
437 The case for open preprints in biology. *PLOS Biology*, 11(5):e1001563.

438 Elliott, D. B. (2014). The impact factor: a useful indicator of journal quality or fatally flawed?  
439 *Ophthalmic and Physiological Optics*, 34(1):4–7.

440 Essawy, B. T., Goodall, J. L., Voce, D., Morsy, M. M., Sadler, J. M., Choi, Y. D., Tarboton, D. G.,  
441 and Malik, T. (2020). A taxonomy for reproducible and replicable research in environmental  
442 modelling. *Environmental Modelling & Software*, page 104753.

443 Fiedler, A. K., Landis, D. A., and Arduser, M. (2012). Rapid shift in pollinator communities  
444 following invasive species removal. *Restoration Ecology*, 20(5):593–602.

445 Gentleman, R. and Lang, D. T. (2007). Statistical analyses and reproducible research. *Journal of*  
446 *Computational and Graphical Statistics*, 16(1):1–23.

447 Gerstner, K., Moreno-Mateos, D., Gurevitch, J., Beckmann, M., Kambach, S., Jones, H. P., and  
448 Seppelt, R. (2017). Will your paper be used in a meta-analysis? Make the reach of your  
449 research broader and longer lasting. *Methods in Ecology and Evolution*, 8(6):777–784.

450 Haddaway, N. R. and Verhoeven, J. T. A. (2015). Poor methodological detail precludes  
451 experimental repeatability and hampers synthesis in ecology. *Ecology and Evolution*,  
452 5(19):4451–4454.

453 Hampton, S. E., Strasser, C. A., Tewksbury, J. J., Gram, W. K., Budden, A. E., Batcheller, A. L.,  
454 Duke, C. S., and Porter, J. H. (2013). Big data and the future of ecology. *Frontiers in Ecology*  
455 *and the Environment*, 11(3):156–162. \_eprint:  
456 <https://esajournals.onlinelibrary.wiley.com/doi/pdf/10.1890/120103>.

457 Hewitt, J. K. (2012). Editorial policy on candidate gene association and candidate  
458 gene-by-environment interaction studies of complex traits. *Behavior Genetics*, 42(1):1–2.

459 Ioannidis, J. P. A. (2005). Why most published research findings are false. *PLOS Medicine*,  
460 2(8):e124.

461 Kurtzer, G. M., Sochat, V., and Bauer, M. W. (2017). Singularity: scientific containers for  
462 mobility of compute. *PLOS ONE*, 12(5):e0177459.

463 Laine, C., Goodman, S. N., Griswold, M. E., and Sox, H. C. (2007). Reproducible research:  
464 moving toward research the public can really trust. *Annals of Internal Medicine*, 146(6):450.

465 Martin, R. C. (2009). *Clean code: a handbook of agile software craftsmanship*. Prentice Hall,  
466 Upper Saddle River, NJ, USA.

467 Marwick, B., Boettiger, C., and Mullen, L. (2018). Packaging data analytical work reproducibly  
468 using R (and friends). *The American Statistician*, 72(1):80–88.

469 McKiernan, E. C., Bourne, P. E., Brown, C. T., Buck, S., Kenall, A., Lin, J., McDougall, D.,  
470 Nosek, B. A., Ram, K., Soderberg, C. K., Spies, J. R., Thaney, K., Updegrave, A., Woo, K. H.,  
471 and Yarkoni, T. (2016). How open science helps researchers succeed. *eLife*, 5:e16800.

472 Merkel, D. (2014). Docker: lightweight Linux containers for consistent development and  
473 deployment. *Linux Journal*, 2014(239):2:2.

474 Mills, J. A., Teplitsky, C., Arroyo, B., Charmantier, A., Becker, P. H., Birkhead, T. R., Bize, P.,  
475 Blumstein, D. T., Bonenfant, C., Boutin, S., Bushuev, A., Cam, E., Cockburn, A., Côté, S. D.,  
476 Coulson, J. C., Daunt, F., Dingemans, N. J., Doligez, B., Drummond, H., Espie, R. H. M.,  
477 Festa-Bianchet, M., Frentiu, F., Fitzpatrick, J. W., Furness, R. W., Garant, D., Gauthier, G.,  
478 Grant, P. R., Griesser, M., Gustafsson, L., Hansson, B., Harris, M. P., Jiguet, F., Kjellander, P.,  
479 Korpimäki, E., Krebs, C. J., Lens, L., Linnell, J. D. C., Low, M., McAdam, A., Margalida, A.,

480 Merilä, J., Møller, A. P., Nakagawa, S., Nilsson, J.-, Nisbet, I. C. T., Noordwijk, A. J. v., Oro,  
481 D., Pärt, T., Pelletier, F., Potti, J., Pujol, B., Réale, D., Rockwell, R. F., Ropert-Coudert, Y.,  
482 Roulin, A., Sedinger, J. S., Swenson, J. E., Thébaud, C., Visser, M. E., Wanless, S., Westneat,  
483 D. F., Wilson, A. J., and Zedrosser, A. (2015). Archiving Primary Data: Solutions for  
484 Long-Term Studies. *Trends in Ecology & Evolution*, 30(10):581–589. Publisher: Elsevier.

485 Mills, J. A., Teplitsky, C., Arroyo, B., Charmantier, A., Becker, P. H., Birkhead, T. R., Bize, P.,  
486 Blumstein, D. T., Bonenfant, C., Boutin, S., Bushuev, A., Cam, E., Cockburn, A., Côté, S. D.,  
487 Coulson, J. C., Daunt, F., Dingemanse, N. J., Doligez, B., Drummond, H., Espie, R. H. M.,  
488 Festa-Bianchet, M., Frentiu, F. D., Fitzpatrick, J. W., Furness, R. W., Gauthier, G., Grant, P. R.,  
489 Griesser, M., Gustafsson, L., Hansson, B., Harris, M. P., Jiguet, F., Kjellander, P., Korpimäki,  
490 E., Krebs, C. J., Lens, L., Linnell, J. D. C., Low, M., McAdam, A., Margalida, A., Merilä, J.,  
491 Møller, A. P., Nakagawa, S., Nilsson, J.-, Nisbet, I. C. T., Noordwijk, A. J. v., Oro, D., Pärt, T.,  
492 Pelletier, F., Potti, J., Pujol, B., Réale, D., Rockwell, R. F., Ropert-Coudert, Y., Roulin, A.,  
493 Thébaud, C., Sedinger, J. S., Swenson, J. E., Visser, M. E., Wanless, S., Westneat, D. F.,  
494 Wilson, A. J., and Zedrosser, A. (2016). Solutions for Archiving Data in Long-Term Studies: A  
495 Reply to Whitlock et al. *Trends in Ecology & Evolution*, 31(2):85–87. Publisher: Elsevier.

496 Mislan, K. A. S., Heer, J. M., and White, E. P. (2016). Elevating the status of code in ecology.  
497 *Trends in Ecology & Evolution*, 31(1):4–7.

498 Moonesinghe, R., Khoury, M. J., and Janssens, A. C. J. W. (2007). Most published research  
499 findings are false—but a little replication goes a long way. *PLOS Medicine*, 4(2):e28.

500 Noble, W. S. (2009). A quick guide to organizing computational biology projects. *PLOS*  
501 *Computational Biology*, 5(7):e1000424.

502 Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D., Breckler, S. J., Buck, S.,  
503 Chambers, C. D., Chin, G., Christensen, G., Contestabile, M., Dafoe, A., Eich, E., Freese, J.,  
504 Glennerster, R., Goroff, D., Green, D. P., Hesse, B., Humphreys, M., Ishiyama, J., Karlan, D.,  
505 Kraut, A., Lupia, A., Mabry, P., Madon, T., Malhotra, N., Mayo-Wilson, E., McNutt, M.,  
506 Miguel, E., Paluck, E. L., Simonsohn, U., Soderberg, C., Spellman, B. A., Turitto, J.,  
507 VandenBos, G., Vazire, S., Wagenmakers, E. J., Wilson, R., and Yarkoni, T. (2015). Promoting  
508 an open research culture. *Science*, 348(6242):1422–1425.

509 Open Science Collaboration (2015). Estimating the reproducibility of psychological science.  
510 *Science*, 349(6251):aac4716.

511 Parr, C. S. and Cummings, M. P. (2005). Data sharing in ecology and evolution. *Trends in*  
512 *Ecology & Evolution*, 20(7):362–363.

513 Peng, R. D. (2011). Reproducible research in computational science. *Science*,  
514 334(6060):1226–1227.

515 Piwowar, H. A., Day, R. S., and Fridsma, D. B. (2007). Sharing detailed research data is  
516 associated with increased citation rate. *PLOS ONE*, 2(3):e308.

517 Roche, D. G., Kruuk, L. E. B., Lanfear, R., and Binning, S. A. (2015). Public data archiving in  
518 ecology and evolution: how well are we doing? *PLOS Biology*, 13(11):e1002295.

519 Sandve, G. K., Nekrutenko, A., Taylor, J., and Hovig, E. (2013). Ten simple rules for  
520 reproducible computational research. *PLOS Computational Biology*, 9(10):e1003285.

521 Schooler, J. W. (2014). Metascience could rescue the ‘replication crisis’. *Nature*, 515(7525):9–9.



522 Steen, R. G., Casadevall, A., and Fang, F. C. (2013). Why has the number of scientific retractions  
523 increased? *PLOS ONE*, 8(7):e68397.

524 Stewart, G. (2010). Meta-analysis in applied ecology. *Biology Letters*, 6(1):78–81.

525 Stodden, V., Guo, P., and Ma, Z. (2013). Toward reproducible computational research: an  
526 empirical analysis of data and code policy adoption by journals. *PLOS ONE*, 8(6):e67111.

527 Sullivan, I., DeHaven, A., and Mellor, D. (2019). Open and reproducible research on Open  
528 Science Framework. *Current Protocols Essential Laboratory Techniques*, 18(1):e32.

529 Whitlock, M., McPeck, M., Rausher, M., Rieseberg, L., and Moore, A. (2010). Data archiving.  
530 *American Naturalist*, 175(2):145–146.

531 Whitlock, M. C. (2011). Data archiving in ecology and evolution: best practices. *Trends in*  
532 *Ecology & Evolution*, 26(2):61–65.

533 Whitlock, M. C., Bronstein, J. L., Bruna, E. M., Ellison, A. M., Fox, C. W., McPeck, M. A.,  
534 Moore, A. J., Noor, M. A. F., Rausher, M. D., Rieseberg, L. H., Ritchie, M. G., and Shaw, R. G.  
535 (2016). A Balanced Data Archiving Policy for Long-Term Studies. *Trends in Ecology &*  
536 *Evolution*, 31(2):84–85. Publisher: Elsevier.

537 Wickham, H. (2014). Tidy data. *Journal of Statistical Software*, 59(i10).

538 Xie, Y. (2015). *Dynamic documents with R and knitr*. CRC Press.

## Tables

Table 1: A list of advanced tools commonly used for reproducible research, aggregated by function. This list is not intended to be comprehensive, but should serve as a good starting point for those interested in moving beyond basic requirements.

	Free	Open Source	Website
<b>Data and Code Management</b>			
Version control			
GitHub	Y <sup>a</sup>	N	<a href="https://github.com">https://github.com</a>
BitBucket	Y <sup>a</sup>	N	<a href="https://bitbucket.com">https://bitbucket.com</a>
GitLab	Y <sup>a</sup>	Y	<a href="https://www.gitlab.com">https://www.gitlab.com</a>
Make			
GNU Make	Y	Y	<a href="https://www.gnu.org/software/make/">https://www.gnu.org/software/make/</a>
Software containers and virtual machines			
Docker	Y	Y	<a href="https://docker.com">https://docker.com</a>
Singularity	Y <sup>a</sup>	Y	<a href="https://syslabs.io">https://syslabs.io</a>
Oracle VM VirtualBox	Y	Y	<a href="https://virtualbox.org">https://virtualbox.org</a>
<b>Sharing Research</b>			
Preprint Servers			
ArXiv	Y		<a href="https://arxiv.org/">https://arxiv.org/</a>
bioRxiv	Y		<a href="https://www.biorxiv.org/">https://www.biorxiv.org/</a>
EcoEvoRxiv	Y		<a href="https://ecoevorxiv.org/">https://ecoevorxiv.org/</a>
Manuscript creation			
Overleaf	Y <sup>a</sup>	Y	<a href="https://overleaf.com">https://overleaf.com</a>
TeXstudio	Y	Y	<a href="https://www.texstudio.org/">https://www.texstudio.org/</a>
Rstudio	Y	Y	<a href="https://rstudio.org">https://rstudio.org</a>
Data Repositories			
Dryad	N		<a href="https://datadryad.org/">https://datadryad.org/</a>
Figshare	Y <sup>a</sup>		<a href="https://figshare.com/">https://figshare.com/</a>
Zenodo	Y		<a href="https://zenodo.org/">https://zenodo.org/</a>
Open Science Framework	Y		<a href="https://osf.io/">https://osf.io/</a>

<sup>a</sup> free to use, but paid premium options with more features are available

540 **Figure Captions**

541 **Figure 1.** A ten-point checklist to guide researchers toward greater reproducibility in their  
542 research. Researchers should give careful thought before, during, and after analysis to ensure  
543 reproducibility of their work.

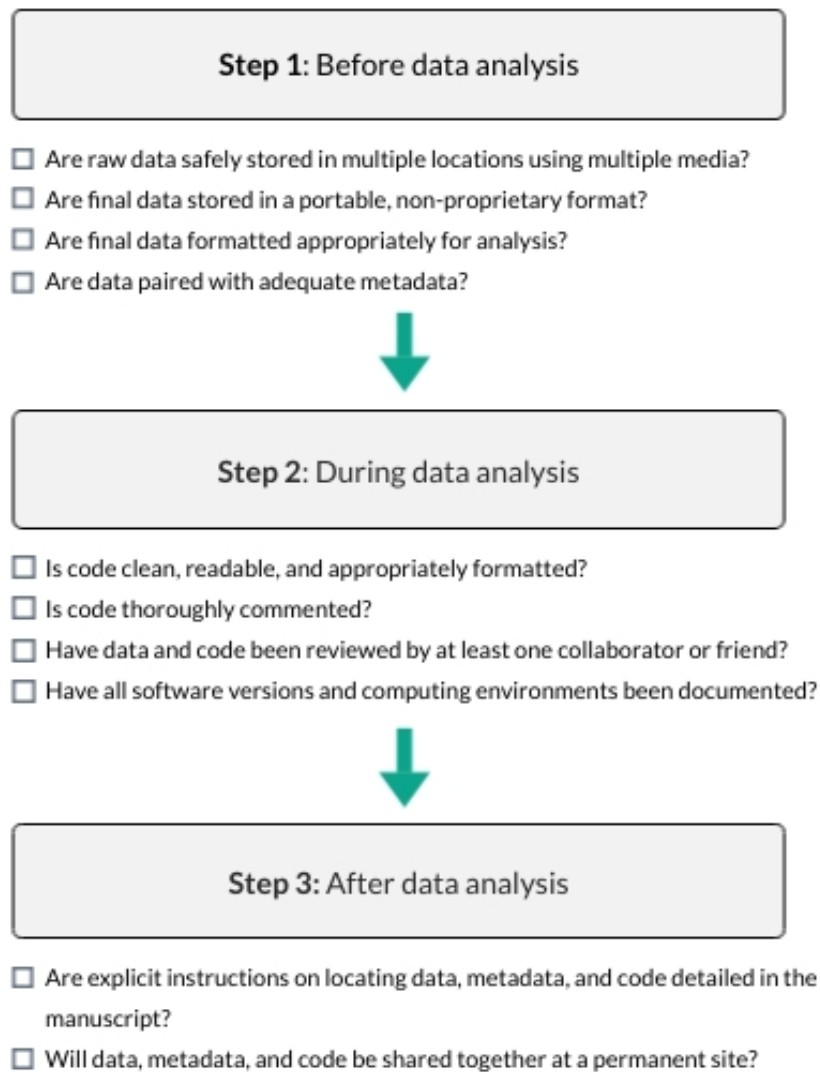


Figure 1: