

1 A Beginner's Guide to Conducting Reproducible
2 Research

3 Jesse M. Alston^{1,2} and Jessica A. Rick^{1,3}

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

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

7 ³*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 most fields within biology remains
11 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 reproduction
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, “reproducible research”
35 is often suggested as a method for increasing our ability to assess the validity and rigor of
36 scientific results (Peng, 2011). Research is reproducible when others can reproduce scientific
37 results given only the original data, code, and documentation (Essawy et al., 2020). This
38 commentary describes basic requirements for such reproducibility in biological research. In it, we
39 make the case for why all research should be reproducible, explain why research is often not
40 reproducible, and present a simple three-part framework all researchers can use to make their
41 research more reproducible. These principles are applicable to researchers working in all types of
42 biological research with data sets of all sizes and levels of complexity.

43 **Why Do Reproducible Research?**

44 **Reproducible research benefits those who do it**

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

49 First, reproducible research helps researchers remember how and why they performed
50 specific analyses during the course of a project. This enables easier explanation of work to
51 collaborators, supervisors, and reviewers, and it allows collaborators to conduct supplementary

52 analyses more quickly and more efficiently.

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

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

63 Fourth, conducting reproducible research is a strong indicator to fellow researchers of rigor,
64 trustworthiness, and transparency in scientific research. This can increase the quality and speed of
65 peer review, because reviewers can directly access the analytical process described in a
66 manuscript. Peer reviewers' work becomes easier and they may be able to answer methodological
67 questions without asking the authors. It also protects researchers from accusations of research
68 misconduct due to analytical errors, because it is unlikely that researchers would openly share
69 fraudulent code and data with the rest of the research community. In addition, reviewers can
70 check whether code matches with methods described in the text of a manuscript, to make sure that
71 authors correctly performed the analyses as described. Finally, it increases the probability that
72 errors are caught during the peer-review process, decreasing the likelihood of corrections or
73 retractions after publication.

74 Finally, reproducible research increases paper citation rates (Piwowar et al., 2007;
75 McKiernan et al., 2016) and allows other researchers to cite code and data in addition to

76 publications. This enables a given research project to have more impact than it would if the data
77 or methods were hidden from the public. For example, researchers can re-use code from a paper
78 with similar methods and organize their data in the same manner as the original paper, then cite
79 code from the original paper in their manuscript. Another researcher may conduct a meta-analysis
80 on the phenomenon described in the two research papers, and thus use and cite both the two
81 papers and the data from those papers in their meta-analysis. Papers are more likely to be cited in
82 these re-use cases if full information about data and analyses are available (Whitlock, 2011;
83 Culina et al., 2018).

84 **Reproducible research benefits the research community**

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

89 First, reproducible research allows others to learn from your work. Scientific research has a
90 steep learning curve, and allowing others to access data and code gives them a head start on
91 performing similar analyses. For example, junior researchers can use code shared with the
92 research community by more senior researchers to learn how to perform advanced analyses. This
93 allows junior researchers to conduct research that is more rigorous from the outset, rather than
94 having to spend months or years trying to figure out “best practices” through trial and error.
95 Modifying existing resources can also save time and effort for experienced researchers—even
96 experienced coders can modify existing code much faster than they can write code from scratch.
97 Sharing code thus allows experienced researchers to perform similar analyses more quickly.

98 Second, reproducible research allows others to understand and reproduce a researcher’s

99 work. Allowing others to access data and code makes it easier for other scientists to perform
100 follow-up studies to increase the strength of evidence for the phenomenon of interest. It also
101 increases the likelihood that similar studies are compatible with one another, and that all of these
102 studies can provide evidence in support of or in opposition to a concept. In addition, sharing data
103 and code increases the utility of these studies for meta-analyses that are important for
104 generalizing and contextualizing the findings of studies on a topic. Meta-analyses in ecology and
105 evolutionary biology are often hindered by incompatibility of data between studies, or lack of
106 documentation for how those data were obtained (Stewart, 2010; Culina et al., 2018).
107 Well-documented, reproducible findings enhance the likelihood that data can be used in future
108 meta-analyses (Gerstner et al., 2017).

109 Third, reproducible research allows others to protect themselves from your mistakes.
110 Mistakes happen in science. Allowing others to access data and code gives them a better chance
111 to critically analyze the work, which can lead to coauthors or reviewers discovering mistakes
112 during the revision process, or other scientists discovering mistakes after publication. This
113 prevents mistakes from compounding over time and provides protection for collaborators,
114 research institutions, funding organizations, journals, and others who may be affected when such
115 mistakes happen.

116 **Barriers to Reproducible Research**

117 There are a number of reasons that most research is not reproducible. Rapidly developing
118 technologies and analytical tools, novel interdisciplinary approaches, unique ecological study
119 systems, and increasingly complex data sets and research questions hinder reproducibility, as does
120 pressure on scientists to publish novel research quickly. This multitude of barriers can be

121 simplified into four primary themes: (1) complexity, (2) technological change, (3) human error,
122 and (4) concerns over intellectual property rights. Each of these concerns can contribute to
123 making research less reproducible and can be valid in some scenarios. However, each of these
124 factors can also be addressed easily via well-developed tools, protocols, and institutional norms
125 concerning reproducible research.

126 **Complexity.** — Science is difficult, and scientific research requires specialized (and often
127 proprietary) knowledge and tools that may not be available to everyone who would like to
128 reproduce research. For example, analyses of genomic data require researchers to possess a vast
129 base of knowledge about statistical methodologies and the molecular architecture of DNA, and
130 genomic analyses are therefore difficult to reproduce for those with limited knowledge of the
131 subject. Some analyses may require high-performance computing clusters that use several
132 different programming languages and software packages, or that are designed for specific
133 hardware configurations. Other analyses may be performed using proprietary software programs
134 such as SAS statistical software (SAS Institute Inc., Cary, NC, USA) or ArcGIS (Esri, Redlands,
135 CA, USA) that require expensive software licenses. Lack of knowledge, lack of institutional
136 infrastructure, and lack of funding all make research less reproducible. However, most of these
137 issues can be mitigated fairly easily. Researchers can cite primers on complex subjects or
138 analyses to reduce knowledge barriers. They can also thoroughly annotate analytical code with
139 comments explaining each step in an analysis, or provide extensive documentation on research
140 software. Using open software (when possible) makes research more accessible for other
141 researchers as well.

142 **Technological change.** — Hardware and software both change over time, and they often
143 change quickly. When old tools become obsolete, research becomes less reproducible. For
144 example, reproducing research performed in 1960 using that era's computational tools would

145 require a completely new set of tools today. Even research performed just a few years ago may
146 have been conducted using software that is no longer available or is incompatible with other
147 software that has since been updated. One minor update in a piece of software used in one minor
148 analysis in an analytical workflow can render an entire project less reproducible. However, this
149 too can be mitigated by using established tools in reproducible research. Careful documentation
150 of versions of software used in analyses is a baseline requirement that anyone can meet. There are
151 also more advanced tools that can help overcome such challenges in making research
152 reproducible, including software containers, which are described in further detail below.

153 **Human error.** — Though fraudulent research is often cited as reason to make research more
154 reproducible (e.g., Ioannidis 2005; Laine et al. 2007; Crocker and Cooper 2011), many more
155 innocent reasons exist as to why research is often difficult to reproduce (e.g., Elliott 2014). People
156 forget small details of how they performed analyses. They fail to describe data collection
157 protocols or analyses completely despite their best efforts and multiple reviewers checking their
158 work. They perform sloppy analyses because they just want to be done with a project that feels
159 like it is taking forever to complete. Science is performed by fallible humans, and a wide variety
160 of common events can render research less reproducible.

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

169 before manipulating the data set as a whole.

170 **Intellectual property rights.** — Researchers often hesitate to share data and code because
171 doing so may allow other researchers to use data and code incorrectly or unethically. Other
172 researchers may use publicly available data without notifying authors, leading to incorrect
173 assumptions about the data that result in invalid analyses. Researchers may use publicly available
174 data or code without citing the original data owners or code writers, who then do not receive
175 proper credit for gathering expensive data or writing time-consuming code. Researchers may
176 want to conceal data from others so that they can perform new analyses on those data in the future
177 without worrying about others scooping them using the shared data. Rational self-interest can
178 lead to hesitation to share data and code via many pathways. However, new tools for sharing data
179 and code are making it easier for researchers to receive credit for doing so and to prevent others
180 from using their data during an embargo period.

181 **A Three-Step Framework for Conducting Reproducible** 182 **Research**

183 Conducting reproducible research is not exceedingly difficult, nor does it require encyclopedic
184 knowledge of esoteric research tools and protocols. Whether they know it or not, most researchers
185 already perform much of the work required to make research reproducible. To clarify this point,
186 we outline below some basic steps toward making research more reproducible in three stages of a
187 research project: (1) before data analysis, (2) during analysis, and (3) after analysis. We discuss
188 practical tips that anyone can use, as well as more advanced tools for those who would like to
189 move beyond basic requirements (Table 1). Most readers will recognize that reproducible
190 research largely consists of widely accepted best practices for scientific research, and that striving

191 to meet a reasonable benchmark of reproducibility is both more valuable and more attainable than
192 researchers may think.

193 **Before data analysis: data storage and organization**

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

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

208 Digital data files should be stored in useful, flexible, portable, non-proprietary formats.
209 Storing data digitally in a “flat” file format is almost always a good idea. Flat file formats are
210 those that store data as plain text with one record per line (e.g., .csv or .txt files) and are the
211 most portable formats across platforms, as they can be opened by anyone without proprietary
212 software programs. For more complex data types, multi-dimensional relational formats such as
213 json, hdf5, or other discipline-specific formats (e.g., biom and EML) may be appropriate.

214 However, the complexity of these formats makes them difficult for many researchers to access
215 and use appropriately, so it is best to stick with simpler file formats when possible.

216 It is often useful to transform data into a ‘tidy’ format (Wickham, 2014) when cleaning up
217 and standardizing raw data. Tidy data are in long format (i.e., variables in columns, observations
218 in rows), have consistent data structure (e.g., character data are not mixed with numeric data for a
219 single variable), and have informative and appropriately formatted headers (e.g., reasonably short
220 variable names that do not include problematic characters like spaces, commas, and parentheses).
221 Data in this format are easy to manipulate, model, and visualize during analysis.

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

232 Finally, researchers should organize files in a sensible, user-friendly structure and make sure
233 that all files have informative names. It should be easy to tell what is in a file or directory from its
234 name, and a consistent naming protocol (e.g., ending the filename with the date created or version
235 number) provides even more information when searching through files in a directory. A consistent
236 naming protocol for both directories and files also makes coding simpler by placing data,
237 analyses, and products in logical locations with logical names. It is often more useful to organize

238 files in small blocks of similar files, rather than having one large directory full of hundreds of
239 files. For example, Noble (2009) suggests organizing computational projects within a main
240 directory for each project, with sub-directories for the manuscript (`doc/`), data files (`data/`),
241 analyses (`scripts/` or `src/`), and analysis products (`results/`) within that directory. While this
242 specific organization scheme may differ for other types of research, keeping all of the research
243 products and documentation for a given project organized in this way makes it much easier to find
244 everything at all stages of the research process, and to archive it or share it with others once the
245 project is finished.

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

258 **During analysis: best coding practices**

259 When possible, all data wrangling and analysis should be performed using coding scripts—as
260 opposed to using interactive or point-and-click tools—so that every step is documented and

261 repeatable by yourself and others. Code both performs operations on data and serves as a log of
262 analytical activities. Because of this second function, code (unlike point-and-click programs) is
263 inherently reproducible. Most errors are unintentional mistakes made during data wrangling or
264 analysis, so having a record of these steps ensures that analyses can be checked for errors and are
265 repeatable on future data sets. If operations are not possible to script, then they should be
266 well-documented in a log file that is kept in the appropriate directory.

267 Analytical code should be thoroughly annotated with comments. Comments embedded
268 within code serve as metadata for that code, substantially increasing its usefulness. Comments
269 should contain enough information for an informed stranger to easily understand what the code
270 does, but not so much that sorting through comments is a chore. Code comments can be tested for
271 this balance by a friend who is knowledgeable about the general area of research but is not a
272 project collaborator. In most scripting languages, the first few lines of a script should include a
273 description of what the script does and who wrote it, followed by small blocks that import data,
274 packages, and external functions. Data cleaning and analytical code then follows those sections,
275 and sections are demarcated using a consistent protocol and sufficient comments to explain what
276 function each section of code performs.

277 Following a clean, consistent coding style makes code easier to read. Many well-known
278 organizations (e.g., RStudio, Google) offer style guidelines for software code that were developed
279 by many expert coders. Researchers should take advantage of these while keeping in mind that all
280 style guides are subjective to some extent. Researchers should work to develop a style that works
281 for them. This includes using a consistent naming convention (e.g., camelCase or snake_case)
282 to name objects and embedding meaningful information in object names (e.g., using “_mat” as a
283 suffix for objects to denote matrices or “_df” to denote data frames). Code should also be written
284 in relatively short lines and grouped into blocks, as our brains process narrow columns of data

285 more easily than longer ones (Martin, 2009). Blocks of code also keep related tasks together and
286 can function like paragraphs to make code more comprehensible.

287 There are several ways to prevent coding mistakes and make code easier to use. First,
288 researchers should automate repetitive tasks. For example, if a set of analysis steps are being used
289 repeatedly, those steps can be saved as a function and loaded at the top of the script. This reduces
290 the size of a script and eliminates the possibility of accidentally altering some part of a function
291 so that it works differently in different locations within a script. Similarly, researchers can use
292 loops to make code more efficient by performing the same task on multiple values or objects in
293 series (though it is also important to note that nesting too many loops inside one another can
294 quickly make code incomprehensible). A third way to reduce mistakes is to reduce the number of
295 hard-coded values that must be changed to replicate analyses on an updated or new data set. It is
296 often best to read in the data file(s) and assign parameter values at the beginning of a script, so
297 that those variables can then be used throughout the rest of the script. When operating on new
298 data, these variables can then be changed once at the beginning of a script rather than multiple
299 times in locations littered throughout the script.

300 Because incompatibility between operating systems or program versions can inhibit the
301 reproducibility of research, the current gold standard for ensuring that analyses can be used in the
302 future is to create a software container, such as a Docker (Merkel, 2014) or Singularity
303 (Kurtzer et al., 2017) image (Table 1). Containers are lightweight, standalone, portable
304 environments that contain the entire computing environment used in an analysis: software, all of
305 its dependencies, libraries, binaries, and configuration files, all bundled into one package.
306 Containers can then be archived or shared, allowing them to be used in the future, even as
307 packages, functions, or libraries change over time. If creating a software container is infeasible or
308 a larger step than readers are willing to take, it is important to thoroughly report all software

309 packages used, including version numbers.

310 **After data analysis: finalizing results and sharing**

311 After the steps above have been followed, it is time for the step most people associate with
312 reproducible research: sharing research with others. As should be clear by now, sharing the data
313 and code is far from the only component of reproducible research; however, once Steps 1 and 2
314 above are followed, it becomes the easiest step. All input data, scripts, program versions,
315 parameters, and important intermediate results should be made publicly and easily accessible.
316 Various solutions are now available to make data sharing convenient, standardized, and accessible
317 in a variety of research areas. There are many ways to do this, several of which are described
318 below.

319 Just as it is better to use scripts than interactive tools in analysis, it is better to produce tables
320 and figures directly from code than to manipulate these using Adobe Illustrator, Microsoft
321 Powerpoint, or other image editing programs. A large number of errors in finished manuscripts
322 come from not remembering to change *all* relevant numbers or figures when a part of an analysis
323 changes, and this task can be incredibly time-consuming when revising a manuscript. Truly
324 reproducible figures and tables are created directly with code and integrated into documents in a
325 way that allows automatic updating when analyses are re-run, creating a “dynamic” document.
326 For example, documents written in \LaTeX and markdown incorporate figures directly from a
327 directory, so a figure will be updated in the document when the figure is updated in the directory
328 (see Xie 2015 for a much lengthier discussion of dynamic documents). Both \LaTeX and markdown
329 can also be used to create presentations that can incorporate live-updated figures when code or
330 data change, so that presentations can be reproducible as well. If using one of these tools is too
331 large a leap, then simply producing figures directly from code—instead of adding annotations and

332 arranging panels post-hoc—can make a substantial difference in increasing the reproducibility of
333 these products.

334 Beyond creating dynamic documents, it is possible to make data wrangling, analysis, and
335 creation of figures, tables, and manuscripts a “one-button” process using GNU Make
336 (<https://www.gnu.org/software/make/>). GNU Make is a simple, yet powerful tool that can be used
337 to coordinate and automate command-line processes, such as a series of independent scripts. For
338 example, a `Makefile` can be written that will take the input data, clean and manipulate it, analyze
339 it, produce figures and tables with results, and update a `LATEX` or `markdown` manuscript document
340 with those figures, tables, and any numbers included in the results. Setting up research projects to
341 run in this way takes some time, but it can substantially expedite re-analyses and reduce
342 copy-paste errors in manuscripts.

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

354 Similarly, data and code shared on personal websites are only available as long as websites
355 are maintained, and can be difficult to transfer when researchers migrate to another domain or

356 website provider. Materials archived on personal websites are also often difficult for other
357 scientists to find, as they are not usually linked to the published research and lack a permanent
358 digital object identifier (DOI). To make research accessible to everyone, it is therefore better to
359 use tools like data and code repositories than personal websites.

360 Data archiving in online repositories has become more popular in recent years, a trend
361 resulting from a combination of improvements in technology for sharing data, an increase in
362 omics-scale data sets, and an increasing number of publisher and funding organizations who
363 encourage or mandate data archiving (Whitlock et al., 2010; Whitlock, 2011; Nosek et al., 2015).
364 Data repositories are large databases that collect, manage, and store data sets for analysis, sharing,
365 and reporting. Repositories may be either subject- or data-specific, or cross-disciplinary general
366 repositories that accept multiple data types. Some are free and others require a fee for depositing
367 data. Journals often recommend appropriate repositories on their websites, and these
368 recommendations should be consulted when submitting a manuscript. Three commonly used
369 general purpose repositories are Dryad, Zenodo, and Figshare; each of these creates a DOI that
370 allows data and code to be citable by others. Before choosing a repository, researchers should
371 explore commonly used options in their specific fields of research.

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

378 In particular, the Open Science Framework (OSF; <http://osf.io/>) is a project management
379 repository that goes beyond the repository features of Dryad, Zenodo, and Figshare to integrate

380 and share components of a research project using collaborative tools. The goal of the OSF is to
381 enable research to be shared at every step of the scientific process—from developing a research
382 idea and designing a study, to storing and analyzing collected data and writing and publishing
383 reports or papers (Sullivan et al., 2019). OSF is integrated with many other reproducible research
384 tools, including widely used pre-print servers, version control software, and publishers.

385 **Conclusions**

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

399 **Acknowledgements**

400 Many thanks to J.G. Harrison, B.J. Rick, and A.L. Lewanski for providing helpful comments on
401 early versions of this manuscript, and to C.A. Buerkle for inspiring this project during his
402 Computational Biology course at the University of Wyoming.

403 **References**

- 404 Alston, J. M. (2019). Open access principles and practices benefit conservation. *Conservation*
405 *Letters*, 12(6):e12672.
- 406 Archmiller, A. A., Johnson, A. D., Nolan, J., Edwards, M., Elliott, L. H., Ferguson, J. M.,
407 Iannarilli, F., Vélez, J., Vitense, K., Johnson, D. H., and Fieberg, J. (2020). Computational
408 reproducibility in The Wildlife Society’s flagship journals. *Journal of Wildlife Management*,
409 84(5):1012–1017.
- 410 Bohannon, J. (2015). Many psychology papers fail replication test. *Science*, 349(6251):910–911.
- 411 Cokol, M., Ozbay, F., and Rodriguez-Esteban, R. (2008). Retraction rates are on the rise. *EMBO*
412 *reports*, 9(1):2–2.
- 413 Crocker, J. and Cooper, M. L. (2011). Addressing scientific fraud. *Science*,
414 334(6060):1182–1182.
- 415 Culina, A., Crowther, T. W., Ramakers, J. J. C., Gienapp, P., and Visser, M. E. (2018). How to do
416 meta-analysis of open datasets. *Nature Ecology & Evolution*, 2(7):1053–1056.
- 417 Desjardins-Proulx, P., White, E. P., Adamson, J. J., Ram, K., Poisot, T., and Gravel, D. (2013).
418 The case for open preprints in biology. *PLOS Biology*, 11(5):e1001563.

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

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

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

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

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

431 Haddaway, N. R. and Verhoeven, J. T. A. (2015). Poor methodological detail precludes
432 experimental repeatability and hampers synthesis in ecology. *Ecology and Evolution*,
433 5(19):4451–4454.

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

436 Ioannidis, J. P. A. (2005). Why most published research findings are false. *PLOS Medicine*,
437 2(8):e124.

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

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

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

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

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

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

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

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

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

457 Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D., Breckler, S. J., Buck, S.,
458 Chambers, C. D., Chin, G., Christensen, G., Contestabile, M., Dafoe, A., Eich, E., Freese, J.,
459 Glennerster, R., Goroff, D., Green, D. P., Hesse, B., Humphreys, M., Ishiyama, J., Karlan, D.,
460 Kraut, A., Lupia, A., Mabry, P., Madon, T., Malhotra, N., Mayo-Wilson, E., McNutt, M.,

461 Miguel, E., Paluck, E. L., Simonsohn, U., Soderberg, C., Spellman, B. A., Turitto, J.,
462 VandenBos, G., Vazire, S., Wagenmakers, E. J., Wilson, R., and Yarkoni, T. (2015). Promoting
463 an open research culture. *Science*, 348(6242):1422–1425.

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

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

468 Peng, R. D. (2011). Reproducible research in computational science. *Science*,
469 334(6060):1226–1227.

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

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

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

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

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

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

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

- 482 Sullivan, I., DeHaven, A., and Mellor, D. (2019). Open and reproducible research on Open
483 Science Framework. *Current Protocols Essential Laboratory Techniques*, 18(1):e32.
- 484 Whitlock, M., McPeck, M., Rausher, M., Rieseberg, L., and Moore, A. (2010). Data archiving.
485 *American Naturalist*, 175(2):145–146.
- 486 Whitlock, M. C. (2011). Data archiving in ecology and evolution: best practices. *Trends in*
487 *Ecology & Evolution*, 26(2):61–65.
- 488 Wickham, H. (2014). Tidy data. *Journal of Statistical Software*, 59(i10).
- 489 Xie, Y. (2015). *Dynamic documents with R and knitr*. CRC Press.

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

^a free to use, but paid premium options with more features are available

491 **Figure Captions**

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

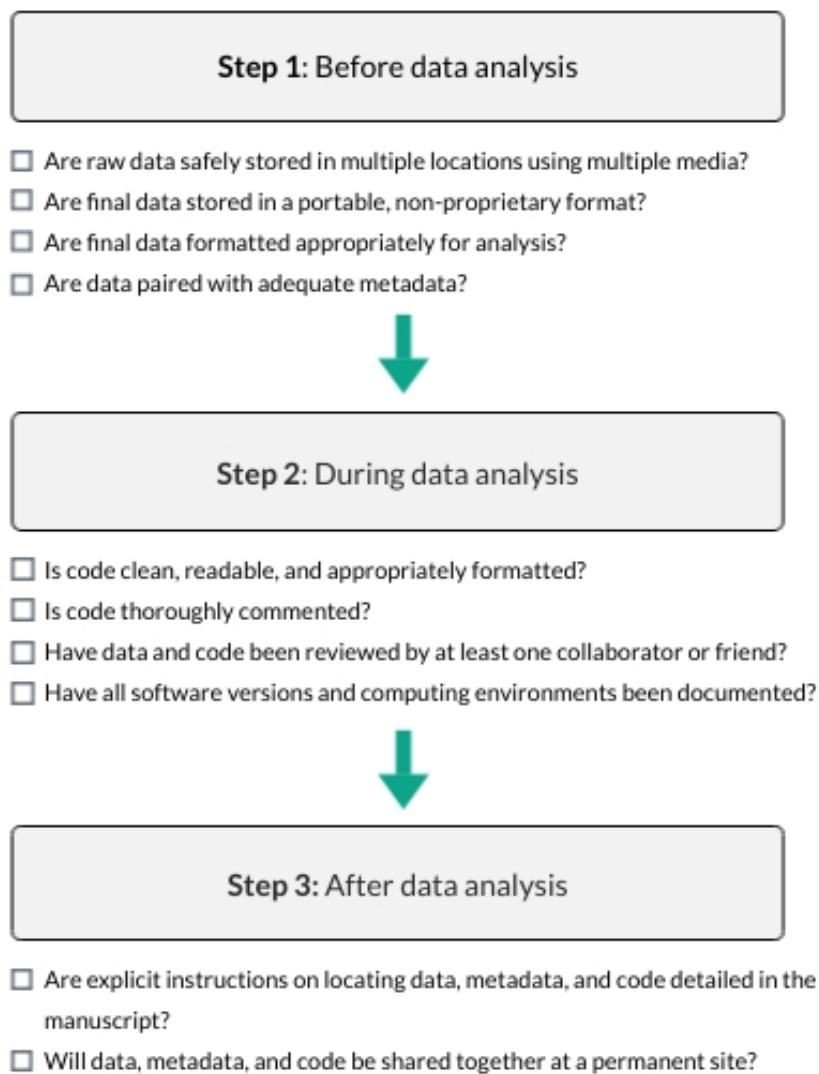


Figure 1: