

1 *Implementing Code Review in the Scientific Workflow: Insights from*
2 *Ecology and Evolutionary Biology*

3 Edward R. Ivimey-Cook*¹, Joel L. Pick², Kevin Bairos-Novak³, Antica Culina⁴, Elliot Gould⁵, Matt
4 Grainger⁶, Benjamin M. Marshall⁷, David Moreau⁸, Matthieu Paquet⁹, Raphaël Royauté¹⁰,
5 Alfredo Sánchez-Tójar¹¹, Inês Silva¹², and Saras M. Windecker*¹³

6
7 *corresponding author(s)

8 e.ivimeycook@gmail.com & saras.windecker@gmail.com



9
10 Affiliations:

11 1. School of Biodiversity, One Health and Veterinary Medicine, University of Glasgow, Glasgow, UK

12 2. Institute of Ecology and Evolution, University of Edinburgh, Edinburgh, UK

13 3. Australian Research Council Centre of Excellence for Coral Reef Studies, James Cook University, Townsville, Queensland, Australia

14 4. Rudjer Boskovic Institute, Zagreb, Croatia, and the Netherlands Institute of Ecology, NIOO-KNAW, Wageningen, The Netherlands

15 5. School of Ecosystem and Forest Sciences, University of Melbourne, Melbourne, VIC, Australia

16 6. Norwegian Institute for Nature Research, 5685 Torgarden, 7485, Trondheim, Norway

17 7. Biological and Environmental Sciences, Faculty of Natural Sciences, University of Stirling, Stirling, United Kingdom

18 8. School of Psychology and Centre for Brain Research, University of Auckland, Auckland, New Zealand

19 9. Institute of Mathematics of Bordeaux, University of Bordeaux, CNRS, Bordeaux INP, Talence, France

20 10. Université Paris-Saclay, INRAE, AgroParisTech, UMR EcoSys, 91120 Palaiseau, France

21 11. Department of Evolutionary Biology, Bielefeld University, 33615, Bielefeld, Germany

22 12. Department of Evolutionary Biology, Bielefeld University, 33615, Bielefeld, Germany

23 13. School of Ecosystem and Forest Sciences, University of Melbourne, Melbourne, VIC, Australia

24

25 Keywords: reliability, reproducibility, software development, coding errors, research process,
26 open science, transparency.

27

28 *Abstract*

29 Code review increases reliability and improves reproducibility of research. As such, code review
30 is an inevitable step in software development and is common in fields such as computer
31 science. However, despite its importance, code review is noticeably lacking in ecology and
32 evolutionary biology. This is problematic as it facilitates the propagation of coding errors and a
33 reduction in reproducibility and reliability of published results. To address this, we provide a
34 detailed commentary on how to effectively review code, how to set up your project to enable this
35 form of review and detail its possible implementation at several stages throughout the research
36 process. This guide serves as a primer for code review, and adoption of the principles and
37 advice here will go a long way in promoting more open, reliable, and transparent ecology and
38 evolutionary biology.

39

40 *Introduction*

41 Across scientific disciplines, researchers increasingly rely on code written in open-source
42 software, such as R and Python, to clean, manipulate, visualise, and analyse data (Mislán *et al.*,
43 2016; Lai *et al.*, 2019; Peikert & Brandmaier, 2021; Peikert *et al.*, 2021). Such software allows
44 for increased transparency and reproducibility compared to software that operates through
45 point-and-click interfaces (“User Interface” or “UI-based”), such as Minitab and SPSS (Obels *et al.*,
46 2020). One of the key benefits of this code-based software is flexibility, because researchers
47 can tailor analyses to their specific research needs which would otherwise be unavailable.
48 However, the flexibility of code comes at a cost, as it means that it can be more error-prone
49 (Budd *et al.*, 1998). These errors may be conceptual (e.g., implementing the wrong function for

50 a given task), programmatic (e.g., indexing the wrong column of a data frame), or syntactic
51 (e.g., the incorrect spelling of a statement or function). Although UI-based software is also prone
52 to conceptual errors, programmatic and syntactic errors are more common in code-based
53 software. These errors can contribute to a lack of reproducibility or to the propagation of
54 incorrect results (see Obels *et al.*, 2020 for a review of code and data in psychology). Indeed,
55 several high-profile retractions have centred on these types of mistakes (Miller, 2006; Ma &
56 Chang, 2007; Bolnick & Paull, 2009; Huijgen *et al.*, 2012; Williams & Bürkner, 2020). One way
57 to minimise potential errors, besides carefully annotating code and following best coding
58 practices, is to undergo a process of code review. However, unlike in some disciplines (such as
59 in computer science and software development) where code review is routinely implemented
60 (Nelson & Schumann, 2004; Badampudi *et al.*, 2019), it is noticeably absent from the research
61 and publication processes in other academic disciplines that rely on code to make inferences
62 and predictions (Indriasari *et al.*, 2020), including ecology and evolutionary biology.

63
64 To address this, we advocate for a fundamental shift in research culture that brings code review
65 into all stages of the research process, as reviewing of code is necessary to facilitate error
66 correction and to confirm the reproducibility and reliability of reported results. This is particularly
67 important as analyses are becoming ever more complicated, especially in the fields of ecology
68 and evolutionary biology (Touchon & McCoy, 2016). But how can we implement code review?
69 By whom, when, and how can it take place? In this paper, we provide some suggestions about
70 how to conduct a code review and how to produce code that facilitates this form of review.
71 Finally, we discuss the application of code review throughout the entire process of publication,
72 from the early stages of pre-publishing right through to after work is published. Although we
73 focus mainly on issues and techniques related to the R and Python coding languages due to
74 their popularity in the fields of ecology and evolutionary biology (Mislán *et al.*, 2016; Lai *et al.*,
75 2019), the concepts and principles we discuss are widely applicable.

76

77 *What should code review evaluate?*

78 Code review is the process of either formally (as part of the peer review process) or informally
79 (as co-authors or colleagues) checking and evaluating each other's code. It is critical to help
80 avoid conceptual, programmatic, and syntactic errors in code and can take place at any stage of
81 the research cycle; pre-submission, during formal peer review, or post-publication. Although the
82 manner and scope in which code review occurs may vary depending on the position in the
83 research cycle, the core priorities remain the same: to ensure code is as reported in the
84 methods section, is able to successfully run, is reliable, and is able to reproduce stated results.
85 Below we describe these key priorities as the four Rs of code review (Figs. 1 and 2):

86

87 *Is the code as Reported?*

88 Code is a key research output and a critical component of scientific methodology. As such, open
89 code accompanying written methods sections is becoming more common, following similar
90 pushes for Open and FAIR data (Lamprecht *et al.*, 2020). Therefore, it is imperative that code is
91 checked for consistency when presented with the corresponding manuscript. These questions
92 help us avoid conceptual errors in code. Does the code match the description of what is
93 "Reported" within the methods section (Fig. 1, SM Box 1)? Ensuring code matches the methods
94 reported is imperative to evaluate whether the code is doing what is stated in the manuscript
95 and what it is intended to do by the user. For instance, methods may state that an analysis uses
96 a generalised linear model with Poisson error, but the code instead fits a Gaussian error
97 distribution. Reviewing for this mismatch must be part of code review. In addition, and equally
98 important for reproducibility is whether the relevant packages (with appropriate version
99 numbers) are stated somewhere in the manuscript. In general, it is good practice to, at the very
100 least, list the packages (with version numbers) that are integral to the analysis or to visualisation

101 in the manuscript. These can be obtained by using the “citation()” function in R or using the
102 {setuptools} package in Python. A full list of all packages used (and versions), for instance those
103 involved with cleaning and tidying of data, could be given elsewhere such as in an associated .R
104 or .py file. Importantly, this will allow for any package or module with versions that are found to
105 contain bugs or coding errors to be identified at a later stage. Packages such as {renv} (Ushey,
106 2023; which replaces {packrat}, Ushey *et al.*, 2022), {groundhog} (Simonsohn & Gruson, 2023),
107 or {poetry} (Eustace, 2023) and {pipenv} (Pipenv Maintainer Team, 2023) in Python can help
108 with ensuring a reproducible environment and allow for specific loading of desired package
109 versions. Another option is containerisation through the use of Docker (Boettiger, 2015; N.B.
110 detailed tutorials already exist which highlight the use of this reproducible method in far more
111 detail than we will discuss here).

112

113 Does the code *Run*?

114 Even if code matches the methodology reported in a paper, this does not mean the code is
115 executable (i.e., can “Run”). Programmatic and syntactic errors can make code fail to rerun. For
116 example, code will not be able to be run if it includes calls to libraries (or modules) that are
117 not installed in the current computing environment or if there are spelling mistakes (Fig. 1,
118 SM Box 1). Data sharing, where possible, should accompany code sharing, so that code can be
119 fully rerun with the original data. If data sharing is not possible, simulated data or a data snippet
120 should be provided so that the code can be rerun. In cases where it would take a long period of
121 time to rerun code (for instance with some forms of Bayesian modelling), the code should be
122 accompanied with appropriate model outputs (readily provided by the author, see below “Output
123 reproducibility”).

124

125 Is the code *Reliable*?

126 Errors can still propagate through code that runs and produces an output, because code can
127 produce incorrect results in a reproducible manner (i.e., every time the code is run). For
128 example, if code selects or modifies the wrong column in a dataset, the code will still run, but
129 produce a reproducible yet inaccurate result (i.e., the code is not “Reliable”; Fig. 1, SM Box 1).
130 This type of error could easily be conceptual, arising from a misunderstanding of the dataset, or
131 programmatic, such as from indexing by number and producing a mistaken column order or
132 from user-defined functions. Although some coding techniques, such as explicitly indexing by
133 column name or by performing unit testing of any user-defined function (see Cooper, 2017;
134 relevant packages include {testthat}, (Wickham, 2011) in R or {pytest} in Python (Okken, 2022),
135 can help avoid many of these mistakes, this type of error is common and also extremely difficult
136 to pick up by anyone without deep familiarity with the dataset and code. In particular, these
137 errors are thought to scale with the number of lines and complexity of code (Lipow, 1982).
138 Although intrinsically linked to evaluating whether code can be run (the second “R”), evaluating
139 code reliability means not only ensuring that the code runs to completion without error, but
140 examining intermediate outputs of the code to ensure there are no mistakes. The functions
141 “identical()” in R and “numpy.array_equal()” in Python can be useful at this stage of code review
142 to compare object similarity between newly-generated and previously-saved intermediate
143 outputs.

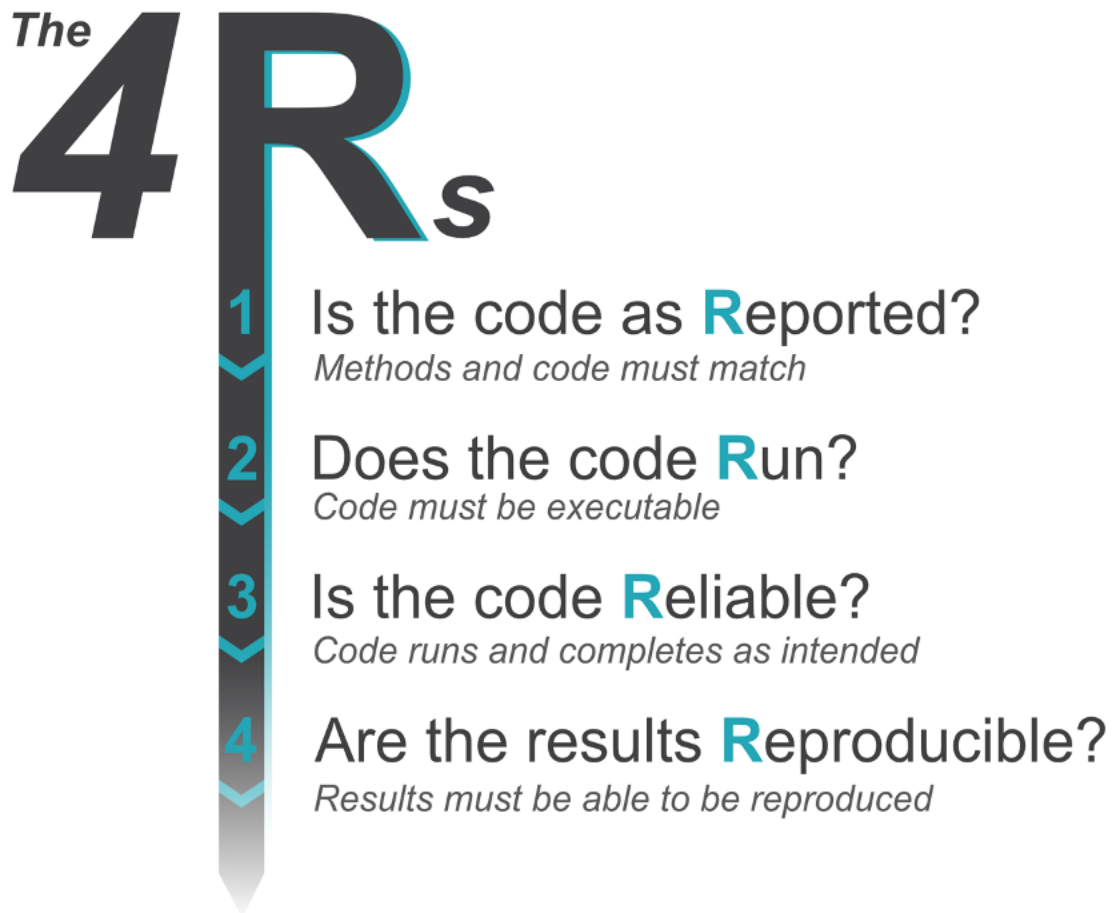
144

145 *Are the results **R**eproducible?*

146 The last “R” of code review builds on the previous code review stages, and is perhaps the most
147 fundamental: can the code produce the output, and thus support the conclusions, given in the
148 paper (Goodman *et al.*, 2016; Fig. 1, SM Box 1)? As several recent papers have highlighted
149 (Archmiller *et al.*, 2020; Obels *et al.*, 2020; Errington *et al.*, 2021; Minocher *et al.*, 2021; Tiwari *et*
150 *al.*, 2021), reproducibility in research results is often very low. Therefore, the final step of code

151 review is ensuring that final outputs when code is rerun match those reported in the analysis
152 and results sections (including any relevant figures and narrative text contained within these
153 sections). With that said, at times obtaining the *exact* same result is not possible. Some level of
154 tolerance must therefore be applied especially when dealing with stochastic methods in which
155 parameter estimates will change between subsequent runs or with techniques that are
156 computationally demanding and slow. This can occur for example if the “set.seed()” function in
157 R or “random.seed” function in Python has not been used prior to stochastic sampling. Providing
158 model outputs can go some way in helping with this (see above), however it does not allow for
159 the code to be explicitly run to see if you can obtain similar results as stated in the paper
160 (regardless of potential time taken). In this case, newly generated results should be compared
161 to the originally reported results. The assessment of how well these two match can then be
162 done using different methods. For example, Archmiller *et al.*, (2020) suggest comparing the
163 conclusion (the direction and significance level) and the numbers (intervals matching within one
164 significant figure) of the original and reproduced results. A useful example of this is given in
165 Archmiller *et al.*'s., (2020) supplementary material, in which they state that a mean of 4.12 and
166 interval of 3.45 to 4.91 reproduces the conclusion and numbers of a study with a mean of 4.00
167 and interval of 3.30 to 5.00. Similar conclusions would be drawn if these means (and CIs) were
168 higher (e.g., 6.5, 6.0 to 7.0), but the numbers would not be considered quantitatively
169 reproduced. However, the conclusions and numbers would not be reproduced if the model
170 instead produced a mean of 4.10 with an interval of -1.00 to 8.40 (as the confidence interval
171 here overlaps with 0). However, the use of one significant figure for comparison of quantitative
172 reproducibility is highly dependent on scale. One alternative, which is unaffected by differences
173 in scale, is provided by Hardwicke *et al.*, (2021) (see also Kambouris *et al.*, 2023) who suggest
174 using % error ($PE = (new - original)/original * 100$). Error is then classified as non-existent when
175 this value is 0, minor if between 0-10 and major (and not reproducible) if 10+.

176 Importantly, as well as the overlap with 0 (or null hypothesis) not changing, the reported and
177 reproduced estimate and intervals should not significantly differ from each other. With stochastic
178 MCMC methods, the variation between chains is expected to be much lower than the variation
179 within chain (i.e., the credible intervals), and so similarly the reproduced value should be well
180 within the stated uncertainty of the reported estimate. It is worth noting and mentioning in your
181 review how closely the numbers and conclusion matched with the reported results.



182

183 Figure 1. The four “Rs” of code review. Figure design by B.M.M.

184

185 *Setting up your code for effective code review*

186 Code review should evaluate if code matches reported methods, whether code runs and is
187 reliable, and lastly, if results can be reproduced. But in order for these questions to be
188 addressed, code must be written and shared in a way that it is possible for someone else to
189 rerun an analysis; both to allow for code to be reviewed and to be reused in the future when
190 properly maintained and contained (Boettiger, 2015). For this to happen, all necessary scripts
191 must be shared along with appropriate metadata indicating how the scripts interact with one
192 another, along with describing all other necessary software and appropriate versions. Often,
193 researchers lack formal training in coding, and learn to code in an *ad-hoc* fashion that excludes
194 training on general styling, appropriate use of workflows, and project organisation. As a result,
195 researchers may often not be aware of the steps necessary to set up code for a project in a
196 manner that reflects best coding practices. Therefore, below we list key principles (Fig. 2) that
197 will help make code reviewable at any stage of the research cycle.

198

199 Project organisation

200 Every project needs some form of directory organisation and folder structure. This is likely to be
201 largely driven by the function and form that your research takes, but an efficient and transparent
202 folder structure that keeps raw data separate from code and intermediate outputs should be
203 created. This helps to ensure that raw data is not accidentally modified or overwritten if any data
204 cleaning or wrangling techniques are applied. A simple folder and file structure such as this will
205 go a long way to help researchers from all coding skill levels understand the order and flow of
206 the data analysis, particularly when the user creates sequentially labelled subfolders and scripts
207 where someone following the code knows which order things must be run (e.g., files beginning
208 with “01...”) in addition to dividing and naming folders to fit their purpose (e.g., data, scripts,
209 function). Several incredibly useful examples already exist (Cooper, 2017; Alston & Rick, 2021;
210 Chure, 2023 see also <https://coderefinery.github.io/reproducible-research/> and

211 https://lakens.github.io/statistical_inferences/14-computationalreproducibility.html). Project code
212 should be stored and available on any data or code repository. Another option for organising a
213 project is to use pipeline or workflow tools (for instance see
214 <https://github.com/pditommaso/awesome-pipeline>), such as the {targets} (Landau, 2021) and
215 {workflowr} R packages (Blischak *et al.*, 2019) or the {luigi} package (The Luigi Authors, 2023)
216 in Python (see
217 <https://www.martinalarcon.org/2018-12-31-a-reproducible-science/>). These tools allow users to
218 automate the process of data analysis, taking a raw dataset through the steps necessary to
219 produce data analysis and visualisation. The advantage to the user is that the code is
220 compartmentalised into logical steps (e.g., import raw data, data cleaning, data wrangling, data
221 analysis, data visualisation) and any changes to the code only affects the downstream steps.
222 For example, if we change the type of analysis we do, we do not need to re-import the data or
223 clean it again. This saves time in computation (especially important for complex, long-running
224 pipelines) but is also advantageous for reproducibility, and sharing and reuse of code.
225 Reviewers can effectively rerun the steps needed to produce a data analysis or figure without
226 having to rerun time consuming pre-processing steps.

227

228 *Project and input metadata*

229 Projects will instantly have better organisation and increased reproducibility when users know
230 *how* they should work through the various folders and subfolders. A README text file and
231 additional metadata gives users the signposts required to facilitate rerunning of code. This can
232 contain information on the packages used (e.g., the package name and version number), along
233 with a detailed description of the various data files, project aim, contact information of the
234 authors, and any relevant licences in place for code or for data (see
235 <https://choosealicense.com/licenses/> for more information). Furthermore, key information about
236 source data is critical for reproducing analysis code. If sharing data is inappropriate to your

237 study (for example when dealing with sensitive confidential data) or if data is so large it cannot
238 be easily shared, then a user can provide a sample of simulated data or a primer so that the
239 code can be checked and read (Quintana, 2020; Hennessy *et al.*, 2022). However, if data is
240 readily available, then providing detailed information about what the data is (preferably in an
241 associated README) and where the data is (e.g., stored on a free data repository such as The
242 Open Science Framework (OSF), Zenodo, or for ecology data, the Knowledge Network for
243 Biocomplexity) should be provided. Metadata should include information such as where the data
244 comes from, who the owners are, as well as what each column header entails, and any relevant
245 acronyms or shorthand notation (ideally following FAIR principles, so data is Findable,
246 Accessible, Interoperable, and Reusable; see Lamprecht *et al.*, 2020). This is particularly useful
247 when controlled vocabulary is used throughout, and R packages such as {codemeta} (Boettiger,
248 2017) and {dataReporter} (Petersen & Ekstrøm, 2019) or Python packages such as
249 {CodeMetaPy} (Gompel, 2023) and {cookiecutter} (Feldroy, 2022) can help with this. Lastly, it is
250 also crucial to explain what data cleaning or curation occurred before the analysis code. For
251 instance, outlining what previous data manipulation or pre-processing steps have been taken to
252 obtain the data in its current state or when an intermediate data file was used.

253

254 Code Readability

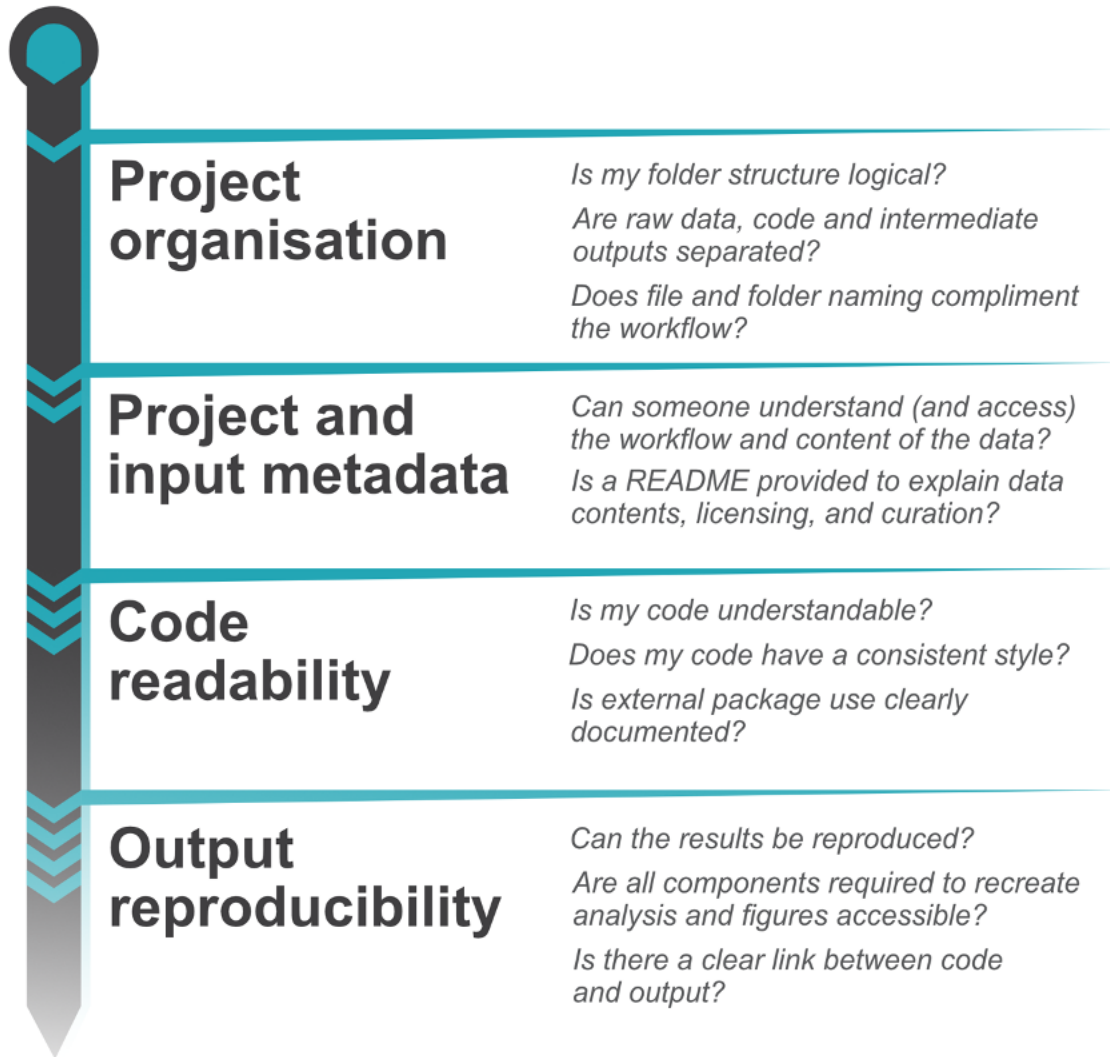
255 Good readability of code is extremely important in enabling effective code review. Several quick
256 solutions exist to provide increased clarity: 1) explicitly calling packages (via a package's
257 namespace, e.g., `package::function()` in R or `package.module.function` in Python); 2) using
258 relative file paths (for instance using the {here} package (Müller, 2020) and preferably with an
259 associated R project file, if using R and RStudio or in a virtual environment if using Python); 3)
260 removing redundant packages; and 4) writing analysis code with clear subheadings, intuitive
261 coding comments, and easy-to-understand object names. Best practice coding tips can be
262 implemented by R packages such as {styler} (Müller & Walthert, 2020) or {pycodestyle} in Python

263 (Rocholl, 2022) and can format code in a number of standardised styles (e.g., Google, tidyverse
264 in R, or PeP8 in Python) with a single line of code or a click of a button. Furthermore, the use of
265 R/Python Markdown, Quarto or the open-source integrated development environment, Jupyter
266 Notebook (or its extension, JupyterLab) enables users to present code in chunks which, along with
267 suitably descriptive comments, allows for far easier readability. In addition, several recent guides
268 and primers have been written that focus on increasing coding cleanliness (Sweigart, 2020;
269 Hunter-Zinck *et al.*, 2021; Filazzola & Lortie, 2022), so we urge the reader to consult these
270 guidelines for tips and advice on improving code readability.

271

272 Output reproducibility

273 One of the key principles and requirements of code is the ability to correctly reproduce
274 published graphs, statistics, and results. In order to do so, a user's code needs to provide a
275 clear link between each section of the code and the various reported graphs and outputs to
276 enable comparison of code to paper and to results. This should then facilitate checking that the
277 results produced by the code match the stated results in the publication. In some cases,
278 reproducing analysis from models can take considerable time to complete, for instance when re-
279 running complicated Bayesian models or other techniques involving long computational time. In
280 this case the "exact" reproducibility of results is not always possible if code must simulate a
281 stochastic process (e.g., Monte Carlo sampling methods). In this case using `set.seed()` or
282 saving simulation outputs still allows for reproducible results (e.g., with the "saveRDS" function
283 in R or the "pickle.dump" function in Python) and can enable code reviewers to check the
284 reproducibility of the reported results.



286

287 Figure 2. A basic workflow for reviewable code that can be adopted from the onset of a project.

288 See Supplementary Material for a printable checklist of the points listed here. Figure and

289 checklist design by B.M.M.

290

291 *Pre-Publication: Setting up a code review group*

292 Informal training coupled with insufficient time and incentives (Touchon & McCoy, 2016), means
293 that coding and subsequent analysis are often the responsibility of a single member of a team
294 throughout a project's entire lifetime. This is in stark contrast to the research-team wide
295 collaboration typical when developing methodology and experimental design. The individual
296 nature of writing research code is part of what makes pre-publication code review so unlikely,
297 but even more critical. Although code review has a place in the formal peer review process and
298 post-publication, one of the most important places for code review to take place is before
299 publication.

300 To achieve this, there must be a culture of peer code review among research teams. One of the
301 most effective methods by which researchers can establish a culture of peer code review in
302 research labs or among colleagues is by setting up a code review group. Here we draw on our
303 experience building a code review club (which we set up in collaboration with the Society for
304 Open, Reliable, and Transparent Ecology and Evolutionary Biology, SORTEE) to present tips
305 for establishing this type of community. In particular, we focus on advice for removing the
306 barriers people have towards sharing their code and receiving feedback; be these due to a lack
307 of time and incentive, a lack of technical knowledge and unclear workflows, or due to social
308 pressures and the fear of being judged by peers (Gomes *et al.*, 2022).

309

310 *Encourage collaboration from the start of a project*

311 Code review can begin as early as the first initiation of a project and play a role beyond
312 publication; it is useful to keep continuous code collaboration at all stages of a manuscript.
313 Collaboration can be facilitated through various code-sharing platforms such as GitHub where
314 users can submit and comment on pull requests (Braga *et al.*, 2023). At SORTEE we
315 established a peer review group and used GitHub issues to summarise discussion of an
316 individual's code during an interactive zoom session (see

317 <https://github.com/SORTEE/peer-code-review/issues/8> for an example including a summary).

318 However, it is important to find a method of facilitating code review that works for your group.

319

320 *Set clear goals for the review*

321 Setting out what you want to achieve with each code review session is particularly important

322 when it comes to organising peer review meetings. Is the focus on general learning and

323 improving readability or is it to error-check and scrutinise the reproducibility of your code?

324 Having a clear structure and goal for each peer review session is important in order to focus

325 comments and advice to address the precise reason for review. Similarly, unless the aim of a

326 code review is to evaluate different analytical options, it would be better to leave methodological

327 questions aside to ensure code review is streamlined.

328

329 *Normalize coding errors and establish a judgement-free environment*

330 Code review volunteers often feel very anxious about showing code that may have errors. It is

331 therefore vital to normalise the existence of errors and highlight that perfection is never possible.

332 It is also useful to stress that there is no such thing as bad code (Barnes, 2010) and there are

333 usually multiple ways to approach the same problem (Silberzahn *et al.*, 2018; Botvinik-Nezer *et*

334 *al.*, 2020). One of the most important statements for peer code review is that there is no single

335 way to code. It is important for code review not to get bogged down by modifying or

336 homogenising style; as long as code is readable, then coding diversity should be encouraged. It

337 is important to create a relaxed environment where people can learn and correct mistakes

338 without judgement or fear of failure and everyone in the peer review group should have a

339 chance to contribute and speak.

340

341 *Carefully consider group size*

342 Usually, a smaller group is a friendly starting point for peer code review because it allows
343 people to feel more comfortable speaking up and participating. Small peer review groups
344 (potentially even one-to-one) can better facilitate peer-to-peer learning and a more focused
345 review of code. However, there are also times when larger groups are more effective, such as
346 having wider discussions on general themes and tips. It is worth considering the aims in
347 establishing the group to help guide the ideal size. For instance, if your goal is to facilitate more
348 general discussions, then a big group size is more likely to enable this. However, if your goal is
349 to enable more focused review of code, then perhaps it is better to reduce the size of the peer
350 review group for this purpose.

351

352 Consider the incentives

353 Code review, outside of paper submission and the formal peer review process, can have a large
354 impact on an individual's project, from error-checking, to validation of appropriate statistical
355 analyses. This then poses the question: what incentives should reviewers of code get? If
356 deemed appropriate, the reviewer could be acknowledged using the MeRIT (Method Reporting
357 with Initials for Transparency) system (Nakagawa *et al.*, 2023), "e.g., J.L.P. ran a linear mixed
358 model with a Gaussian error distribution. Code was checked by E.I.C.". In some circumstances,
359 it may even be appropriate for the reviewer to obtain co-authorship of the paper, if the review
360 fundamentally altered the project and subsequent paper. For instance, a situation may arise
361 where a code reviewer(s) finds a major coding error which, when fixed after highlighting and
362 reproducing the issue to the author(s), alters the subsequent results and conclusions of the
363 manuscript. Ultimately, incentives should be relative to the impact of the reviewer on the
364 project.

365

366 *During Publication: Formal code review*

367 One of the most crucial aspects of code review can take place during the formal peer review
368 process. This is where reviewers are able to carefully follow and understand the logic of
369 analyses, much like the flow of writing from the introduction to the discussion of a paper
370 (Powers & Hampton, 2019). In some journals, such as The Royal Society (“Data sharing and
371 mining | Royal Society,” 2023), Behavioural Ecology and Sociobiology (Bakker & Traniello,
372 2020), and The American Naturalist (Bolnick, 2022) *both* code and data are requested for
373 review at from the submission stage (although, based on our experience, this is not enforced).
374 In some cases, such as in Journal of Open Source Software, the entire process of formal peer
375 review, including that of code and manuscript is hosted on GitHub and implemented via GitHub
376 issues (see <https://github.com/openjournals/joss-reviews/issues> for several useful examples).
377 This, as Fernández-Juricic (2021) points out, has several benefits. For authors, providing code
378 during peer review could lead to an increase in the quality of the manuscript, and for reviewers,
379 available code allows for a far deeper insight into the manuscript as there is a clearer link
380 between experimental methodology and statistical analysis (the First R; code as “Reported”).
381 These benefits are substantial and could ultimately contribute to the adoption of code review
382 during the publication process by journals.

383 However, beyond the availability of code during submission, there are numerous other hurdles
384 before effective and in-depth code review can be reasonably formalised as part of the peer
385 review process. One of the most pressing issues is finding suitable individuals to review code
386 given there is already a lack of willing reviewers in the current system. It is reasonable to expect
387 reviewers to check that code is as reported, but anything more in-depth could take up the time
388 of already overworked academics, who may not necessarily have the exact expertise needed to
389 check other people’s code. Reviewers could be asked to state if they are competent to review
390 the code is as reported, and the journal could ensure that at least one reviewer has checked the
391 code. Another vital issue is how to provide code and data during peer review within a double-

392 blind peer review system (which has been shown to significantly reduce peer review bias; see
393 Fox *et al.*, 2023). There are several solutions to this. Anonymised data and code could be
394 submitted directly to the journal during review (e.g., as a zip file). Alternatively, the anonymised
395 data and code could be uploaded to repositories such as figshare, the Open Science
396 Framework (OSF), GitHub or Dryad (although we note the latter may not be freely available).
397 These repositories allow the authors to provide an anonymised link for peer review (see also
398 [https://methodsblog.com/2023/08/23/double-anonymous-peer-review-frequently-asked-](https://methodsblog.com/2023/08/23/double-anonymous-peer-review-frequently-asked-questions/)
399 [questions/](https://methodsblog.com/2023/08/23/double-anonymous-peer-review-frequently-asked-questions/) for useful links to repositories that enable double-blinding). While there are still
400 issues that need to be fully considered before any kind of extended code review becomes a
401 standard part of the peer review process, the mandatory provision of code for peer review
402 alongside the explicit expectation that at least one reviewer checks the code matches the
403 reported methodology (i.e., *Is the code as reported?*) would make a strong start in shifting
404 publication culture and increasing the reliability of published research.

405

406 *Post-Publication: Reviewing code after publication*

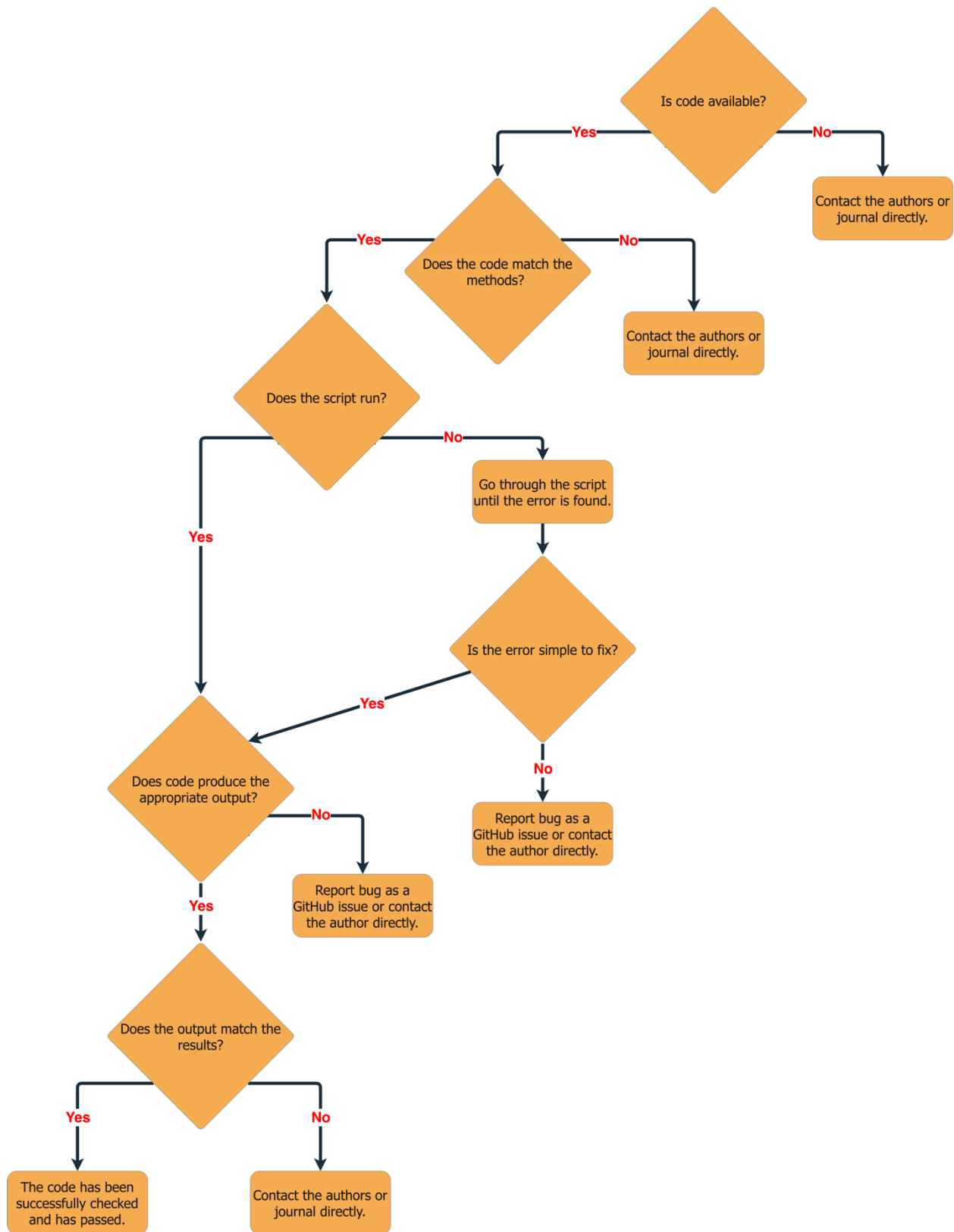
407 Reviewing code post-publication is another facet of code review that has been much less
408 discussed. Although it does not prevent *publication* of incorrect results, it does enable checking
409 if code is indeed adhering to the R's listed above (Fig. 1). The initial question is 'has all code
410 used to produce the results been made available'? An increasing number of journals are now
411 *requesting* code be shared alongside scientific articles (Culina *et al.*, 2020), such as in
412 supplemental materials or by linking to an online repository. This then allows for any open and
413 shared code to be checked and verified alongside methods section statements (Stodden, 2011;
414 Light *et al.*, 2014). Unfortunately, unlike data, code is a lot less likely to be made available
415 regardless of these mandatory journal policies. As Figure 2 from Culina *et al.*, (2020) shows,

416 although the number of journals that possesses a mandatory code rule is increasing (from 15%
417 in 2015 to 79% in 2020) the number of articles that actually provide open code is still around
418 27% (although this number varies considerably among journals). This suggests that not many
419 authors are adhering to this policy, which is an impediment to computational reproducibility
420 (Culina *et al.*, 2020). However, there is hope to be found here. As Culina *et al.* have shown,
421 journals requiring code to be shared are increasing in number yearly and, as a field, we already
422 have improved substantially (Mislán *et al.*, 2016; Culina *et al.*, 2020; Jenkins *et al.*, 2023). In
423 some cases, journals have implemented far stricter (and rightly so) data and code requirements
424 along with assigning corresponding data editors (Bolnick, 2022). However, the first necessary
425 step is for all journals to make it a requirement for *both* code and data to be present from the
426 very start of the submission stage (Powers & Hampton, 2019; Fernández-Juricic, 2021). But
427 what happens if the code is not available? In this case, the main option is to reach out to the
428 corresponding author (or perhaps the journal itself) and ask if the code could be made available;
429 similar to data being made available “upon reasonable request”.

430 The next part is relevant to the previous section above (“What should code review evaluate?). If
431 you find that the code associated with a manuscript does not adhere to any of the “R”s listed
432 above, then the first step is to contact the corresponding author (or if the paper uses the MeRIT
433 system (Nakagawa *et al.*, 2023), the person who actually conducted the analysis). This could be
434 in the form of a GitHub issue if there is a repository for the code or an email (see Fig. 3). If there
435 is indeed an error in code, and it is not due to differences in software version (e.g., differences
436 in R and package versions) or due to inherent stochasticity (e.g., simulations or MCMC
437 sampling), then the authors should be given a chance to contact the journal themselves to
438 highlight and correct their mistakes. For instance, as per American Naturalist’s stance (Bolnick,
439 2022) authors who contact the journal to correct code or data errors will not be penalised and
440 corrections are encouraged (when warranted).

441 However, in cases where updated results would alter the narrative of a published paper, corrections
442 may be more difficult to address without newer methods of documenting changes. Publication
443 versioning or “living” documents may present a solid first step in such a scenario (Kane & Amin,
444 2023). By encouraging post-publication code review, we can both decrease the proliferation of
445 coding errors and also increase the reliability of published science.

446



447

448 Figure 3. An example peer code review flowchart that can occur post-publication. Figure design

449 by J.L.P and E.I.C.

450 *Concluding remarks*

451 In this brief overview, we have provided a basic set of guidelines for peer code review,
452 recommendations for producing reviewable code, and considerations for how it should be
453 adopted at every level of research throughout the publication process. The principles and advice
454 listed here should form a baseline for code review that should be improved upon. We hope that
455 this encourages coders at all levels to try and promote more reproducible, transparent, and
456 open coding practices. In addition, we hope that this provides a primer to start a code reviewing
457 club of your own.

458

459 *Acknowledgements*

460 This work began during workshops at the 2021 and 2022 annual conferences of the Society for
461 Open, Reliable, and Transparent Ecology and Evolutionary biology (SORTEE) run by S.M.W.
462 and E.I-C. (in 2022). We also would like to acknowledge Fonti Kar for her help with organising
463 and delivering the 2021 workshop. This work was partially funded by the Center of Advanced
464 Systems Understanding (CASUS), which is financed by Germany's Federal Ministry of
465 Education and Research (BMBF) and by the Saxon Ministry for Science, Culture and Tourism
466 (SMWK) with tax funds on the basis of the budget approved by the Saxon State Parliament. We
467 also thank Corina Logan, Serena Caplins, the PREreview group (Gracielle Higinio, Varina
468 Crisfield, Mobina Gholamhosseini, Katherine Hébert, and Tanya Strydom), and one anonymous
469 reviewer for thoughtful comments on the manuscript. We thank Melina de Souza Leite for
470 allowing us to use their GitHub issue as an example of SORTEE peer code review.

471

472 *References*

473 Alston, J.M. & Rick, J.A. (2021) A Beginner's Guide to Conducting Reproducible Research.
474 *Bulletin of the Ecological Society of America*, **102**, 1–14.

475 Archmiller, A.A., Johnson, A.D., Nolan, J., Edwards, M., Elliott, L.H., Ferguson, J.M., *et al.*
476 (2020) Computational Reproducibility in The Wildlife Society's Flagship Journals. *The Journal*
477 *of Wildlife Management*, **84**, 1012–1017.

478 Badampudi, D., Britto, R. & Unterkalmsteiner, M. (2019) Modern code reviews - Preliminary
479 results of a systematic mapping study. In *Proceedings of the Evaluation and Assessment on*
480 *Software Engineering*, EASE '19. Association for Computing Machinery, New York, NY, USA,
481 pp. 340–345.

482 Bakker, T.C.M. & Traniello, J.F.A. (2020) Ensuring data access, transparency, and preservation:
483 mandatory data deposition for Behavioral Ecology and Sociobiology. *Behavioral Ecology and*
484 *Sociobiology*, **74**, 132.

485 Barnes, N. (2010) Publish your computer code: it is good enough. *Nature*, **467**, 753.

486 Blischak, J.D., Carbonetto, P. & Stephens, M. (2019) Creating and sharing reproducible research
487 code the workflow way. *F1000Research*, **8**, 1749.

488 Boettiger, C. (2015) An introduction to Docker for reproducible research. *ACM SIGOPS*
489 *Operating Systems Review*, **49**, 71–79.

490 Boettiger, C. (2017) Generating CodeMeta Metadata for R Packages. *The Journal of Open*
491 *Source Software*, **2**, 454.

492 Bolnick, D. (2022) EIC Update: American Naturalist policy on data and code archiving.

493 Bolnick, D.I. & Paull, J.S. (2009) Morphological and dietary differences between individuals are
494 weakly but positively correlated within a population of threespine stickleback. *Evolutionary*
495 *Ecology Research*.

496 Botvinik-Nezer, R., Holzmeister, F., Camerer, C.F., Dreber, A., Huber, J., Johannesson, M., *et*
497 *al.* (2020) Variability in the analysis of a single neuroimaging dataset by many teams. *Nature*,
498 **582**, 84–88.

499 Braga, P.H.P., Hébert, K., Hudgins, E.J., Scott, E.R., Edwards, B.P.M., Sánchez Reyes, L.L., *et*
500 *al.* (2023) Not just for programmers: How GitHub can accelerate collaborative and reproducible
501 research in ecology and evolution. *Methods in Ecology and Evolution*, 1–17.

502 Chure, G. (2023) Git + GitHub As A Platform For Reproducible Research.

503 Cooper, N. (2017) A Guide to Reproducible Code in Ecology and Evolution.

504 Culina, A., Berg, I. van den, Evans, S. & Sánchez-Tójar, A. (2020) Low availability of code in
505 ecology: A call for urgent action. *PLOS Biology*, **18**, e3000763.

506 Data sharing and mining | Royal Society [WWW Document]. (2023) . URL
507 <https://royalsociety.org/journals/ethics-policies/data-sharing-mining/> [accessed on 2023].

508 Errington, T.M., Denis, A., Perfito, N., Iorns, E. & Nosek, B.A. (2021) Challenges for assessing
509 replicability in preclinical cancer biology. *eLife*, **10**, e67995.

510 Eustace, S. (2023) poetry: Python dependency management and packaging made easy.

511 Feldroy, A. (2022) cookiecutter: A command-line utility that creates projects from project
512 templates, e.g. creating a Python package project from a Python package project template.

513 Fernández-Juricic, E. (2021) Why sharing data and code during peer review can enhance
514 behavioral ecology research. *Behavioral Ecology and Sociobiology*, **75**, 103.

515 Filazzola, A. & Lortie, C. (2022) A call for clean code to effectively communicate science.
516 *Methods in Ecology and Evolution*.

517 Fox, C.W., Meyer, J. & Aimé, E. (2023) Double-blind peer review affects reviewer ratings and
518 editor decisions at an ecology journal. *Functional Ecology*, **37**, 1144–1157.

519 Gomes, D.G.E., Pottier, P., Crystal-Ornelas, R., Hudgins, E.J., Foroughirad, V., Sánchez-Reyes,
520 L.L., *et al.* (2022) Why don't we share data and code? Perceived barriers and benefits to public
521 archiving practices. *Proceedings of the Royal Society B: Biological Sciences*, **289**, 20221113.
522 Gompel, M. van. (2023) CodeMetaPy: Generate and manage CodeMeta software metadata.
523 Hardwicke, T.E., Bohn, M., MacDonald, K., Hembacher, E., Nuijten, M.B., Peloquin, B.N., *et*
524 *al.* (2021) Analytic reproducibility in articles receiving open data badges at the journal
525 Psychological Science: an observational study. *Royal Society Open Science*, **8**, 201494.
526 Hennessy, E.A., Acabchuk, R.L., Arnold, P.A., Dunn, A.G., Foo, Y.Z., Johnson, B.T., *et al.*
527 (2022) Ensuring Prevention Science Research is Synthesis-Ready for Immediate and Lasting
528 Scientific Impact. *Prevention Science*, **23**, 809–820.
529 Huijgen, R., Boekholdt, S.M., Arsenault, B.J., Bao, W., Davaine, J.-M., Tabet, F., *et al.* (2012)
530 RETRACTED: Plasma PCSK9 Levels and Clinical Outcomes in the TNT (Treating to New
531 Targets) Trial: A Nested Case-Control Study. *Journal of the American College of Cardiology*,
532 **59**, 1778–1784.
533 Hunter-Zinck, H., Siqueira, A.F. de, Vásquez, V.N., Barnes, R. & Martinez, C.C. (2021) Ten
534 simple rules on writing clean and reliable open-source scientific software. *PLOS Computational*
535 *Biology*, **17**, e1009481.
536 Indriasari, T.D., Luxton-Reilly, A. & Denny, P. (2020) A Review of Peer Code Review in
537 Higher Education. *ACM Transactions on Computing Education*, **20**, 1–25.
538 Jenkins, G.B., Beckerman, A.P., Bellard, C., Benítez-López, A., Ellison, A.M., Foote, C.G., *et*
539 *al.* (2023) Reproducibility in ecology and evolution: Minimum standards for data and code.
540 *Ecology and Evolution*, **13**, e9961.
541 Kambouris, S., Wilkinson, D.P., Smith, E.T. & Fidler, F. (2023) Computationally reproducing
542 results from meta-analyses in Ecology and Evolutionary Biology using shared code and data.
543 Kane, A. & Amin, B. (2023) Amending the literature through version control. *Biology Letters*,
544 **19**, 20220463.
545 Lai, J., Lortie, C.J., Muenchen, R.A., Yang, J. & Ma, K. (2019) Evaluating the popularity of R in
546 ecology. *Ecosphere*, **10**, e02567.
547 Lamprecht, A.-L., Garcia, L., Kuzak, M., Martinez, C., Arcila, R., Martin Del Pico, E., *et al.*
548 (2020) Towards FAIR principles for research software. *Data Science*, **3**, 37–59.
549 Landau, W.M. (2021) The targets R package: a dynamic Make-like function-oriented pipeline
550 toolkit for reproducibility and high-performance computing. *Journal of Open Source Software*, **6**,
551 2959.
552 Light, R.P., Polley, D.E. & Börner, K. (2014) Open data and open code for big science of science
553 studies. *Scientometrics*, **101**, 1535–1551.
554 Lipow, M. (1982) Number of Faults per Line of Code. *IEEE Transactions on Software*
555 *Engineering*, **SE-8**, 437–439.
556 Ma, C. & Chang, G. (2007) Retraction for Ma and Chang, Structure of the multidrug resistance
557 efflux transporter EmrE from Escherichia coli. *Proceedings of the National Academy of*
558 *Sciences*, **104**, 3668–3668.
559 Miller, G. (2006) A Scientist's Nightmare: Software Problem Leads to Five Retractions. *Science*,
560 **314**, 1856–1857.
561 Minocher, R., Atmaca, S., Bavero, C., McElreath, R. & Beheim, B. (2021) Estimating the
562 reproducibility of social learning research published between 1955 and 2018. *Royal Society*
563 *Open Science*, **8**, 210450.

564 Mislan, K.A.S., Heer, J.M. & White, E.P. (2016) Elevating The Status of Code in Ecology.
565 *Trends in Ecology & Evolution*, **31**, 4–7.

566 Müller, K. (2020) *here: A Simpler Way to Find Your Files*.

567 Müller, K. & Walthert, L. (2020) Styler: Non-invasive pretty printing of R code. *R package*
568 *version 1.3. 2*.

569 Nakagawa, S., Ivimey-Cook, E.R., Grainger, M.J., O’Dea, R.E., Burke, S., Drobniak, S.M., *et al.*
570 (2023) Method Reporting with Initials for Transparency (MeRIT) promotes more granularity and
571 accountability for author contributions. *Nature Communications*, **14**, 1788.

572 Nelson, S. & Schumann, J. (2004) What makes a code review trustworthy? In *37th Annual*
573 *Hawaii International Conference on System Sciences, 2004. Proceedings of the*. Presented at the
574 37th Annual Hawaii International Conference on System Sciences, 2004. Proceedings of the, p.
575 10 pp.-.

576 Obels, P., Lakens, D., Coles, N.A., Gottfried, J. & Green, S.A. (2020) Analysis of Open Data
577 and Computational Reproducibility in Registered Reports in Psychology. *Advances in Methods*
578 *and Practices in Psychological Science*, **3**, 229–237.

579 Okken, B. (2022) *Python Testing with pytest*. Pragmatic Bookshelf.

580 Peikert, A. & Brandmaier, A.M. (2021) A Reproducible Data Analysis Workflow With R
581 Markdown, Git, Make, and Docker. *Quantitative and Computational Methods in Behavioral*
582 *Sciences*, 1–27.

583 Peikert, A., Lissa, C.J. van & Brandmaier, A.M. (2021) Reproducible Research in R: A Tutorial
584 on How to Do the Same Thing More Than Once. *Psych*, **3**, 836–867.

585 Petersen, A.H. & Ekstrøm, C.T. (2019) **dataMaid** : Your Assistant for Documenting Supervised
586 Data Quality Screening in R. *Journal of Statistical Software*, **90**.

587 Pipenv maintainer team. (2023) pipenv: Python Development Workflow for Humans.

588 Powers, S.M. & Hampton, S.E. (2019) Open science, reproducibility, and transparency in
589 ecology. *Ecological Applications*, **29**, e01822.

590 Quintana, D.S. (2020) A synthetic dataset primer for the biobehavioural sciences to promote
591 reproducibility and hypothesis generation. *eLife*, **9**, e53275.

592 Rocholl, J.C. (2022) pycodestyle: Python style guide checker.

593 Silberzahn, R., Uhlmann, E.L., Martin, D.P., Anselmi, P., Aust, F., Awtrey, E., *et al.* (2018)
594 Many analysts, one data set: Making transparent how variations in analytic choices affect results.
595 *Advances in Methods and Practices in Psychological Science*, **1**, 337–356.

596 Simonsohn, U. & Gruson, H. (2023) groundhog: Version-Control for CRAN, GitHub, and
597 GitLab Packages.

598 Stodden, V. (2011) Trust Your Science? Open Your Data and Code, 2.

599 Sweigart, A. (2020) *Beyond the Basic Stuff with Python: Best Practices for Writing Clean Code*.
600 No Starch Press.

601 The Luigi Authors. (2023) luigi: Workflow mgmt + task scheduling + dependency resolution.

602 Tiwari, K., Kananathan, S., Roberts, M.G., Meyer, J.P., Sharif Shohan, M.U., Xavier, A., *et al.*
603 (2021) Reproducibility in systems biology modelling. *Molecular Systems Biology*, **17**, e9982.

604 Touchon, J.C. & McCoy, M.W. (2016) The mismatch between current statistical practice and
605 doctoral training in ecology. *Ecosphere*, **7**, e01394.

606 Ushey, K., McPherson, J., Cheng, J., Atkins, A., Allaire, J.J. & Allen, T. (2022) packrat: A
607 Dependency Management System for Projects and their R Package Dependencies.

608 Wickham, H. (2011) testthat: Get Started with Testing. *The R Journal*, **3**, 5–10.

609 Williams, D. & Bürkner, P.-C. (2020) Coding Errors Lead to Unsupported Conclusions: A
610 critique of Hofmann et al. (2015). *Meta-Psychology*, 4.
611