

A novel replicability survey tool to measure assess and promote improve data availability and research reproducibility in hydrology

James H. Stagge^{1,2*}, David E. Rosenberg¹, Adel M. Abdallah¹,
Hadia Akbar¹, Nour A. Attallah¹, Ryan James¹

January 11, 2019

1. Utah State University, Department of Civil and Environmental Engineering, Logan, UT 84321

2. The Ohio State University, Department of Civil, Environmental and Geodetic Engineering, Columbus, OH 43210

*corresponding author(s): James H. Stagge (stagge.11@osu.edu)

Abstract

There is broad interest to improve the reproducibility of published research. We developed a simple survey tool to assess the 1) availability of availability of digital research artifacts published alongside peer-reviewed journal articles (e.g. data, models, code, directions, and other digital artifacts, and 2) replicability of results published in peer-reviewed journal articles for use) and reproducibility of article results. We used the tool to assess 360 random-sampled articles of the 1,989 articles published in 2017 in six well-regarded by six hydrology and water resources journals. Like several prior studies, we replicated results for just in 2017. Like studies from other fields, we reproduced results for only a small fraction of articles (1.6% estimated between of tested articles) using their available artifacts. We estimated, with 95% confidence, that results might be reproduced for only 0.6% and to 6.8% for of all 1,989 articles with 95% confidence). Unlike prior studies, the tool helped identify identified key bottlenecks to make making work more reproducible. These bottlenecks Bottlenecks include: only some digital artifacts available (44% of articles), no directions (89%), or all artifacts available but results not replicable reproducible (5%). These results suggest the The tool (or extensions) can help authors, journals, funders, and institutions to self-assess the availability and replicability of manuscripts, provide feedback to authors to improve manuscript improve reproducibility, and recognize and reward reproducible articles as examples for others to follow.

1 Introduction

The science community is broadly interested to improve the reproducibility of research [1, 2, 3, 5, 6][1, 2, 3, 4, 5, 6]. While exact definitions of reproducibility vary [7, 8, 1, 9, 10, 11, 12, 13, 14], ~~one overarching theme~~ [7, 8, 9, 10, 11, 12, 13, 14], there are recent attempts to harmonize definitions [15, 16]. One overarching theme of definitions is that reproducibility requires multiple, progressive components such as (i) all data, models, code, directions, and other digital artifacts used in the research are available for others to reuse (hereafter, “availability”; [17, 18, 19]), (ii) the artifacts can be used to exactly ~~replicate-reproduce~~ published results (~~replicability~~reproducibility, sometimes called bit or computational reproducibility; [20, 21]), and (iii) existing and new datasets can be processed using the artifacts to reproduce published conclusions (~~reproducibility~~-replicability). This progression follows the framework laid out out in a new report on reproducibility and replication by the National Science Foundation and the U.S. Department of Education [16]. In this framework, replicability is a higher standard than reproducibility.

~~Replicable and reproducible~~ Reproducible and replicable scientific work is currently uncommon because of misaligned incentives and poor coordination among authors, journals, institutions, and funding agencies that conduct, publish, and support scientific research [22, 23, 9]. For example, making artifacts available requires authors to document additional materials [24, 25] and learn new skills and technologies [26]. Authors may worry that shared materials will never be used [27, 10] or that other scientists will scoop them on follow-up studies [28]. Further, universities typically reward peer-reviewed journal publications, rather than data repositories or documentation, while current scientific culture rewards novelty rather than ~~replicating-reproducing~~ prior efforts [20, 2].

Several efforts are underway to encourage more reproducible science [9]. Authors can share research materials in a growing number of online repositories such as Github, Figshare, Harvard Dataverse, Dryad, ~~and-or~~ HydroShare. Institutional libraries are transitioning to offer online repositories to house digital research artifacts [29, 30, 31, 32]. Within our ~~field-fields~~ of hydrology and water resources, recent tools provide environments to ~~repository-data-store data~~ publicly and allow software to operate on the data as well as create virtual environments that package code, data, and a working operating system to reduce problems of incompatibility (e.g., [33, 34]). Authors can assign digital object identifiers (DOIs) to research packages to create persistent links and use umbrella research licenses to describe the manner in which these digital artifacts and their associated paper may be legally used by others [22]. Additionally, authors can specify the level of reproducibility that readers and reviewers can expect from each publication, for example that a typical reader could easily reproduce the paper’s results. And yet, despite these powerful tools, few authors are making their work available for others to ~~replicate-or-reproduce-reproduce~~ or replicate.

To quantify the current state of reproducible science in hydrology and to understand the factors preventing more ~~replicable-or-reproducible~~ reproducible

or replicable publications, we present here a simple 15-question survey tool (Fig. 1) designed to assess the availability of digital artifacts and replicability of results in peer-reviewed journal articles (see Methods). We use this survey tool to assess 360 stratified random-sampled articles from the 1,989 articles published in 2017 across six reputable hydrology and water resources journals. Sampling The sampling design was stratified by journal and reproducibility keywords of interest to produce a representative population sample, while increasing the likelihood to include articles with reproducible results. Results identify bottlenecks to making digital artifacts available and replicating output. We also use results to generalize replicability-reproducibility for the entire sample of articles, test a hypothesis about use of keywords to identify replicable-reproducible articles, compare the effectiveness of different journal data availability policies, and highlight how authors, journals, funders, and institutions can use the enclosed survey tool to encourage, recognize, and reward the publication of more reproducible articles.

Results

Applying our survey tool to 360 random-sampled hydrology articles published in 2017 shows that a decreasing number of articles are able to satisfy the progressively stricter reproducibility requirements of artifact availability and ultimately replication-of-reproduction of the published results (Fig. 2). For example, 70.3% of the 360 sampled articles stated some materials were available, but we could only access 48.6% of those materials online (Fig. 3). Only 5.6% of sampled articles made data, model/code, and directions publicly available while just 1.1% of sampled articles made artifacts available and were fully replicated-reproduced. We partially replicated-reproduced an additional 0.6% of articles.

Artifact Availability

Across all journals sampled publications, the most common primary artifact provided was input data, followed by code/model/software, and then directions to run (Fig. 4). Other artifacts-These three primary artifacts were each needed to reproduce modeled results. Secondary artifacts, such as hardware/software requirements, common file formats, unique and persistent digital identifiers, and metadata, were made available at much lower rates than the primary artifacts. Articles published in Environmental Modeling & Software (EM&S) made artifacts available online at a high rate and these articles also provided data/-model/code, directions, hardware/software requirements, common file formats, and metadata at rates 2 to 3 times the next most frequent journal two times or higher than other journals.

Sampled articles use different methods to make artifacts available and these methods differ markedly across journals (Fig. 4). A majority of sampled EM&S had the largest proportion of sampled articles that articles made at least some artifacts available online (61.9%). By contrast, Hydrology and Earth Systems

87 Sciences (HESS) and Water Resources Research (WRR) had high percentages of
88 articles where materials were only available by first author request (38.5-40.2%).
89 Otherwise, the Journal of Hydrology (JoH), Journal of the American Water
90 Resources Association (JAWRA), and Journal of Water Resources Planning and
91 Management (JWRPM) had large proportions of articles where data ~~was not~~
92 ~~available or was available~~ were available within the article ~~and or as~~ supplemental
93 material. These three journals also had high proportions of sampled papers in
94 which research artifacts were not available.

95 ~~Replication~~ Reproducibility of Results

96 Twenty sampled articles (5.6% of total sampled articles) made the required in-
97 put data, software/model/code, and directions available, allowing an attempt
98 at ~~replicating~~ reproducing published results. We were able to fully ~~replicate~~
99 reproduce results for four articles [35, 36, 37, 38] and partially ~~replicated~~ reproduced
100 results for two additional articles [39, 40]. We were unable to ~~replicate~~ reproduce
101 results for four articles [~~?, 42, 43, 44]~~ [41, 42, 43, 44], which nonetheless appeared
102 to provide the necessary materials. During the process to ~~replicate~~ reproduce
103 results, we found 10 of the 20 articles did not have all the required artifacts,
104 despite being initially considered for ~~replication~~ reproducibility testing. Rea-
105 sons we only partially ~~replicated~~ reproduced results for two articles and did not
106 ~~replicate~~ reproduce results for four articles included unclear directions (4 arti-
107 cles), did not generate results (3 articles), hardware/software error (2 articles),
108 or results differed from the publication (1 article; some articles had multiple
109 reasons). The survey permitted multiple selections for this question. A com-
110 mon issue across cases where we did not generate results was that folder and
111 file locations were hard-coded to work on the author's computer. If these issues
112 were obvious, we tried, with limited success, to fix them. Other cases pointed
113 to general data gateways, like the USGS streamgauge network, with no further
114 details or required expensive or proprietary software. Of the 10 articles that had
115 all artifacts available, five were published in EM&S, two articles were published
116 in HESS and in WRR, and the remaining article was published in JWRP&M.

117 ~~Estimated Replicability~~ Reproducibility for Population

118 Because the stratified sampling method oversampled articles with certain repro-
119 ducibility keywords, we used bootstrap resampling (see Methods) to estimate
120 that 0.6% to 6.8% of all 1,989 articles published in 2017 in the six journals tested
121 here would be ~~replicable~~ reproducible (95% confidence interval). We estimated
122 28.0% (23.1-32.6% confidence interval) of all articles published in these journals
123 during 2017 provided at least some of the artifacts necessary for ~~replicability~~
124 reproducibility (Fig. 5, black). EM&S differed from other journals by having
125 a large proportion of articles with some or all data available (31.8-64.0%) and
126 relatively high estimates of ~~replicability~~ reproducibility (Fig. 5).

127 Using Keywords to Identify ~~Replicable~~ Reproducible Arti- 128 cles

129 We found that five of the six articles with some or complete ~~replicability had~~
130 ~~certain reproducibility-related~~ reproducibility had certain related keywords of
131 interest in their abstracts (full list in Methods). This positive hit rate (4.2%)
132 for articles with reproducibility keywords is significantly greater than articles
133 without (0.4%; 2-sample Chi-Squared test with Yate's continuity correction ($p =$
134 0.014)). These findings confirm the value of reproducibility keywords to identify
135 ~~replicable~~ reproducible articles and reaffirm the difficulty to know at the outset
136 whether the results presented in an article are ~~replicable~~ reproducible.

137 Time Required to Determine Availability and ~~Replicability~~ Reproducibility

138 We surveyed and analyzed the time required to complete the survey to show
139 the incremental effort required to determine the availability of article artifacts
140 and ~~replicability~~ reproducibility of results (Fig. 6). For example, for a single
141 publication it took us as little as 5 to 14 minutes (25-75% range) to deter-
142 mine the availability of input data, model/software/code, and directions ~~and~~ ,
143 Reproducing results for a single paper required upwards of 25 to 86 minutes
144 (25-75% range) ~~to fully replicate results with outliers of 25 and~~ , ~~with an upper~~
145 outlier of 200 minutes ~~(minimum and maximum)~~. There were no statistically
146 detectable differences in the time between journals to determine availability of
147 digital artifacts or ~~replicate~~ to reproduce results.

148 Reproducibility and Journal Policies

149 Among the six hydrology and water resources journals we studied, the HESS and
150 WRR policies effective during the 2017 review period require articles to state
151 how data, models, and code can be accessed. In contrast, the 2017 policies
152 by EM&S, WRRPM, JoH, and JAWRA simply encouraged this practice. EM&S
153 further recommends articles include an explicit "Software and/or data availabil-
154 ity" section within the article and requires authors to make software essential to
155 the paper available to reviewers during the review process (Supplemental Mate-
156 rial ~~Table ??~~). HESS includes an assets tab in each publication, based on the
157 Code and Data Availability sections. EM&S, WRR, and JOH are all signatories
158 of the Transparency and Openness Promotion (TOP) policy framework [45],
159 while HESS participates in the FAIR (Findable, Accessible, Interoperable, and
160 Reusable) data project [46].

161 ~~These Stronger~~ journal data availability policies ~~partially but not fully explain~~
162 ~~differences in~~ and making open data commitments tend to produce higher
163 rates of artifact availability and result ~~replicability by journal~~ reproducibility.
164 However, there is significant variation among these journals, likely due to minor
165 differences in implementation or other factors. For example, EM&S, which only
166 encourages authors to make artifacts available, had the highest rate of articles
167 that made artifacts available (Fig. 3) and this high rate persisted across nearly

every artifact category (Fig. 4). Although EM&S used “should” instead of “must” statements, their policy was by far the most specific for papers with a software component (Supplemental [S1 Material](#)). This may explain their high participation rate. EM&S is also a software-focused journal, which may attract papers and authors that are more conscious of reproducible software. In contrast, HESS and WRR, which require data availability statements, had lower percentages of articles that made artifacts available and more papers that direct readers to the authors or third parties for data, models, or code (Fig. 3). These directional statements tend to appear in the Data Availability section of HESS articles and the Acknowledgements of WRR articles. The final group, JoH, JAWRA, and JWRP&M, that also encouraged authors to make artifacts available, had high proportions of articles without available digital artifacts (Fig. 3). The HESS and WRR policies that require data availability statements appear to encourage authors to select options like contact the author rather than work to provide a research article and supporting materials that are available, ~~repleiability~~, ~~and-reproducible~~ [reproducible, and replicable](#). In July 2018, JWRP&M switched to start requiring authors to state the availability of data, models, and code, similar to HESS and WRR [47].

Discussion

Our findings of low ~~repleiability~~ [reproducibility](#) of research published in six hydrology and water resources journals in 2017 mirrors low rates of ~~repleiability~~ [reproducibility](#) previously reported in psychology (100 experiments; [2]), computer systems research (613 articles; [48]), and ~~aeross~~—articles published in Science (204 articles; [6]). Unlike those studies, our survey tool additionally identified bottlenecks to making all digital artifacts available and ~~repleiating~~ [reproducing](#) results. Here, we discuss how results for our study in hydrology and water resources can inform broader use of the survey tool by authors, journals, funders, and institutions to improve the reproducibility of published scientific research.

Authors. Authors can use the survey tool as a checklist to self-assess the availability of data, models, code, and directions and ~~repleiability~~ [reproducibility](#) of their work before submitting work for publication. The tool can help identify missing components that, if provided, will improve reproducibility. For example, our results showed that, for all journals, the number of sampled articles with code/data/software was consistently 2 to 3 times higher than the number of articles that additionally provided directions (Fig. 4). If authors used the tool and subsequently provided directions to use their materials, the tool could potentially double the number of articles which could reasonably be tested for ~~repleiability~~ [reproducibility](#). Another bottleneck was a large fraction of authors who chose not to make artifacts available or only made artifacts available by author or third party request. Authors can look to the 10 articles we found that made all digital artifacts available to see easy-to-access platforms to provide access. These platforms included GitHub (6 articles), HydroShare (1

article), journal ~~supplementary~~-material (1 article), a custom website (1 article), or Figshare (1 article). Authors who bundled their code and data together in a single GitHub repository further allowed us to download the entire project, with a higher likelihood that code pointed to valid file directories.

Journals. Journals can use the survey tool to assess the availability of data/model/code and directions, ~~replicability~~-~~reproducibility~~ of results of new submissions, and provide feedback to authors. Alternatively, journals can require that authors use the survey tool to self-check their work prior to submission. Feedback can be crucial as our study showed that a very low fraction of articles provided all the required artifacts. However, when artifacts were available, we fully or partially ~~replicated~~-~~reproduced~~ results for 6 out of 10 articles. We also found that time to assess the availability of artifacts was typically less than 15 minutes, while time to ~~replicate~~-~~reproduce~~ results was much longer. The combination of these findings suggests that promoting inclusion of digital artifacts through a relatively quick availability survey could pay significant dividends for ~~replicability~~~~reproducibility~~. We leave open whether responsibility for assessment should fall on a journal editor's assistant, associate editor for reproducibility, reviewers, or others. With a tool to measure reproducibility of published articles, journals could also track reproducibility over time. Tracking and publishing this information would benefit the journal as a promotional tool to show journal commitment to reproducible science. Tracking would also allow journals to acknowledge articles and authors that reach certain reproducibility standards, as implemented by Psychological Science [49]. For example, show a bronze, silver, or gold medal icon on article webpages to recognize and reward progressively greater reproducibility corresponding to availability, ~~replicability~~, ~~and~~-~~reproducibility~~~~reproducibility~~, and ~~replicability~~. These badges would simultaneously communicate the expected level of reproducibility of published work. In our study, we are excited to award silver badges to the four articles whose results we fully ~~replicated~~-~~reproduced~~ [35, 36, 37, 38] and bronze badges to six articles that made all artifacts available [39, 41, 42, 43, 44, 40]. This recognition also makes it easier for authors to find and emulate excellent reproducibility practices. We propose these recognition programs as voluntary to encourage authors to make their artifacts available and results reproducible, but not required in cases of proprietary data or code. Cross-journal indices could further aggregate reproducibility metrics and encourage journals and authors to improve the reproducibility of their research portfolios. To oversee these journal efforts, we envision a new role for an Associate Editor of Reproducibility to develop journal data availability and reproducibility policy, manage reproducibility evaluations, and advocate for best reproducibility practices.

Funders and Institutions. Similar to journals, funders and institutions can use the survey tool to assess artifact availability, verify ~~replicability~~-~~reproducibility~~ of results, and recognize or reward authors whose work achieves bronze, silver, and gold levels of reproducibility. Alternatively, funders and institutions could use reproducibility assessments made by journals. Funders can encourage authors to use the survey tool to self-check work prior to submitting progress or final reports or use the tool to check the reproducibility of work authors submit.

257 Use of the tool could help verify that author submissions fulfill requirements of
258 funder data management policies and help direct authors to improve the reproducibility of their work. Institutions could also use the survey tool to determine
259 and post the expected level of reproducibility of author work deposited into
260 institutional repositories.

262 Together, these actions by authors, journals, funders, and institutions can
263 help nudge authors further along the reproducibility continuum to make their
264 digital artifacts more available and to ~~replicate~~ reproduce published results.
265 While these proposed policy nudges represent small shifts targeted at particular
266 actors in the science community, this approach can produce large effects collectively [50], particularly when all parties agree that the shift will provide a net
267 benefit, as for reproducible science. Each individual nudge is made possible by
268 using a survey (or similar) tool to measure and quantify the availability of digital
269 artifacts, ~~replicability~~ reproducibility of published results, and ~~reproducibility~~
270 replicability of findings. We welcome discussion to improve the survey tool
271 aimed at improving the reproducibility of our science.

273 Methods

274 Online Survey Tool

275 The authors translated definitions of availability, ~~replicability, and reproducibility~~
276 reproducibility, and replicability into a 15-question Qualtrics Research Core
277 (Qualtrics) online survey (Fig. 1). The Qualtrics survey format has been converted into a publicly available ~~version~~ Google form, provided here as an example (<https://goo.gl/forms/95S4y9BdPmVqMtm02>). The survey progressed
279 from soliciting metadata about the article (Questions 1-4), to testing availability of artifacts (Q5-9), and ultimately testing ~~replicability~~ reproducibility of
282 results (Q10-14). Green or yellow shaded answers (Fig. 1) were required to progress to the next question so that survey questions followed the availability and ~~replicability~~ reproducibility progression. Selecting a red-shaded answer
285 stopped progression and directed the reviewer to a final question that asked how many minutes the reviewer spent to reach their stopping point (Q15). This
287 time to complete was self-reported by reviewers rather than using the built-in Qualtrics timer so reviewers could consider the entire time spent reading and
288 assessing the published article, rather than the time completing the survey.

290 The authors developed the tool over four months in Fall 2017 and pre-tested it in early 2018 on a sub-sample of five articles that spanned the availability and ~~replicability~~ reproducibility progression. From our experience pre-testing and
293 to improve use of the tool, we reworded some questions, altered the survey logic, discussed and addressed inter-reviewer variability. Later, after we had reviewed
295 23% of sampled articles, we added a final question (Q15) to ask how much time it took to complete the survey. We did not re-analyze the time spent for the initial 23% of papers, as reviewers were already familiar with those papers. Instead, we calculated time spent using papers from the remaining sample.

299 Selection of Articles

300 360 peer-reviewed articles were random stratified sampled from the 1,989 arti-
 301 cles found in Scopus that were published in 2017 ~~in by~~ six well-regarded hydrol-
 302 ogy and water resources journals(~~measured by impact factor and reputation~~).
 303 Journals were selected based on impact and to cover a range of hydrology and
 304 water resources topics. The six journals were Hydrology ~~of and~~ Earth Systems
 305 Sciences (HESS), Environmental Modeling & Software (EM&S), Journal of the
 306 American Water Resources Association (JAWRA), Journal of Hydrology (JoH),
 307 Journal of Water Resources Planning and Management (JWRPM), and Water
 308 Resources Research (WRR). Stratified random sampling was approximately pro-
 309 portional to the number of articles published in each journal in 2017, with extra
 310 weight placed on articles with a set of reproducibility-related keywords (Table
 311 ??).

312 We further adjusted the stratification so each journal had at least 30 ar-
 313 ticles (JAWRA and JWRPM were oversampled). Similarly, we oversampled
 314 articles with the keywords: analytical software, application programs, C++,
 315 cloud computing, computational reproducibility, computer modeling, computer
 316 programming, computer software, computer software usability, computer-based
 317 models, development and testing, engineering software, fortran, freely available
 318 data, freely available software, github, hardware and software, java, open code,
 319 open source, replicative validation, scientific software, code, python, cran, and
 320 http. Of the 120 articles published in the six journals in 2017 that had at
 321 least one keyword, we sampled 119 articles, principally to retain at least 15
 322 non-keyword articles for each journal with an approximately 2:1 non-keyword
 323 to keyword ratio overall.

324 ~~Number of articles published in 2017 and number of articles sampled by~~
 325 ~~journal.~~

326 Each author was randomly assigned 60 articles stratified by journal to assess
 327 the availability of article artifacts (Q1-9). After identifying all publications that
 328 had the available artifacts, we re-assigned reviewers to assess whether the pub-
 329 lished results could be ~~replicated reproduced~~ (Q1-15). ~~Reassignments matched~~
 330 ~~article software to the~~ We carried through responses of "Not sure" or "Not
 331 familiar with resources" to Q9 and reassigned these articles to match article
 332 software with a reviewer most familiar with those tools. The Qualtrics online
 333 format allowed us to both simultaneously and asynchronously assess journal ar-
 334 ticles and store survey responses in an accompanying Qualtrics database. After
 335 all availability and ~~replicability reproducibility~~ assessments were complete, we
 336 exported results from the Qualtrics database to a text file which was processed
 337 in R to generate figures, tables, and results. Time spent to complete the survey
 338 (Q15) was analyzed for three key stopping points: no artifacts available (Q5),
 339 availability of artifacts (Q9), and ~~replicability reproducibility~~ of results (Q13).

340 Population Estimates

341 ~~Bootstrap resampling~~ Resampling was used to estimate the overall percent-

age of articles from all $n=1,989$ articles published in 2017 in the six journals while adjusting for keyword and journal sampling. Sampled articles were sorted into six mutually exclusive categories that were stopping points in the survey: Data-less or review, Author or Third Party Only, No availability, Some ~~Availability~~availability, Available but not reproducible, and Some or ~~All Replicable~~all reproducible. "Some availability" included articles with one or two data/model/code, and direction elements of the three required elements (Q7). "Available but not ~~replicable~~reproducible" articles had all three required elements available on the initial review, but either could not be ~~replicated~~reproduced or were found to be missing key elements when ~~replication was attempted~~reviewers attempted to reproduce the results.

The resampling approach ~~was based on~~generated random 5,000 ~~simulations of the population~~. For each simulation, populations. Each population had 1,989 articles. In each population, we inserted the 360 sampled articles were inserted directly based on the survey resultsarticles we manually assessed, assuming that ~~these were known exactly~~we exactly knew the reproducibility of these articles. Estimates for the remaining 1,629 unsampled articles were simulated based on survey results for the sampled articles in their stratified category, i.e. journal and keyword/non-keyword. ~~Uncertainty for the unsampled articles was calculated using the~~For each random sample population, the proportion of unsampled articles in each reproducibility category was randomly simulated using the multinomial uncertainty approach of Sison and Glaz [51], assuming the six potential levels of availability/replicability represent multinomial proportions of categorical data [52][51, 52]. This produced 5,000 sample populations equal in size and distribution (journal and keyword) to the true population of articles published in 2017, while incorporating uncertainty due to unsampled papers.

Acknowledgements

This material is based upon work supported by Utah Mineral Lease Funds, the National Science Foundation, funded through OIA – 1208732, and the U.S. Fullbright Program. Additional funding was by National Science Foundation grant #1633756. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of any of the funding organizations.

The authors thank Amber S Jones for providing feedback on an early draft, Stephen Maldonado and Marcos Miranda for external review of the code repository, and Ayman Alaffi for participation in early discussions to develop the survey tool.

Author Contributions

Stagge, Rosenberg, James, and Abdallah conceived the idea for a survey tool to measure the reproducibility of journal articles. Stagge and Rosenberg led

the design and refinement of the tool and Abdallah implemented the design. All authors participated equally in the use of the survey tool to evaluate the reproducibility of articles. Stagge led the results analysis, visualization, and writing of the first draft. Stagge and Abdallah made all article digital artifacts available online. Rosenberg rewrote the initial draft. All authors reviewed and approved the final draft prior to submission.

Code Availability

The survey tool, Qualtrics results, and all code used for analysis presented in this article are available online through ~~GitHub~~ (~~https://github.com/jstagge/reproduc_hyd~~) ~~or through a permanent repository [?]~~ the permanent repository [53]. Please cite this repository for any use of the related data or code. Additionally, results can be ~~replicated~~ reproduced using RStudio deployed in the cloud using MyBinder through the GitHub website.

Data Availability

All relevant data presented in this article are available online through the permanent repository [53]. Please cite this repository for any use of the related data or code. An open Google Forms version of the survey tool is available for readers to use, modify, and extend ~~—~~(<https://goo.gl/forms/95S4y9BdPmVqMtm02>).

Competing financial interests

The author(s) declare no competing financial interests.

Paper Metadata

Q1. Assessor's name
Q2. Journal name
Q3. Article DOI
Q4. Full paper citation

Availability

Q5. How accessible to users?

Q6. Where available?

Q7. What is present?

Required

Optional

Q8. Comments on availability [open response].

Q9. Do you estimate you and readers could use the available artifacts to generate results?

Q10. Continue to reproduce results?

Reproducibility

Q11. Do the outputs verify published results (in text, figures, and tables)?

Q12. If yes, explain what made the work reproducible and other comments [open response].

Q13. If no, why did reproducing the work fail?

Q14. Other comments on why reproducing the work failed [open response].

Time to Complete

Q15: How many minutes did the survey take?

Figure 1: Survey questions to assess journal article data availability and ~~replicability~~reproducibility. Green and grey answers continue to the next question, while red answers skip to question 15.

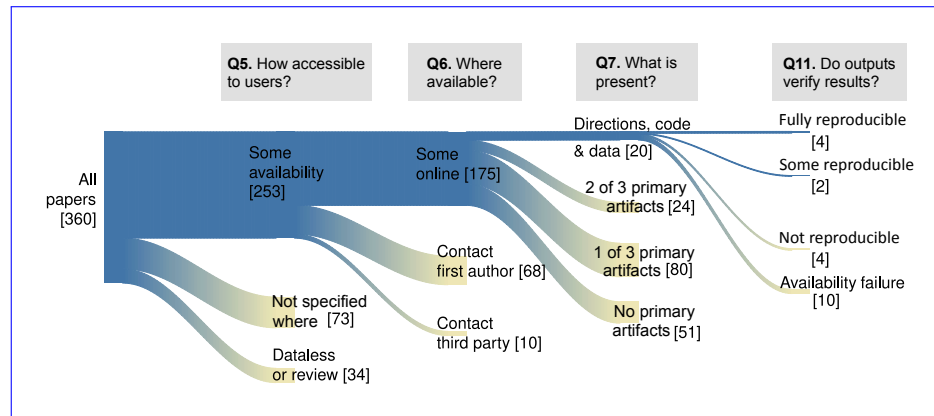


Figure 2: Number of papers progressing through the survey questions to determine availability and ~~reproducibility~~ reproducibility requirements.

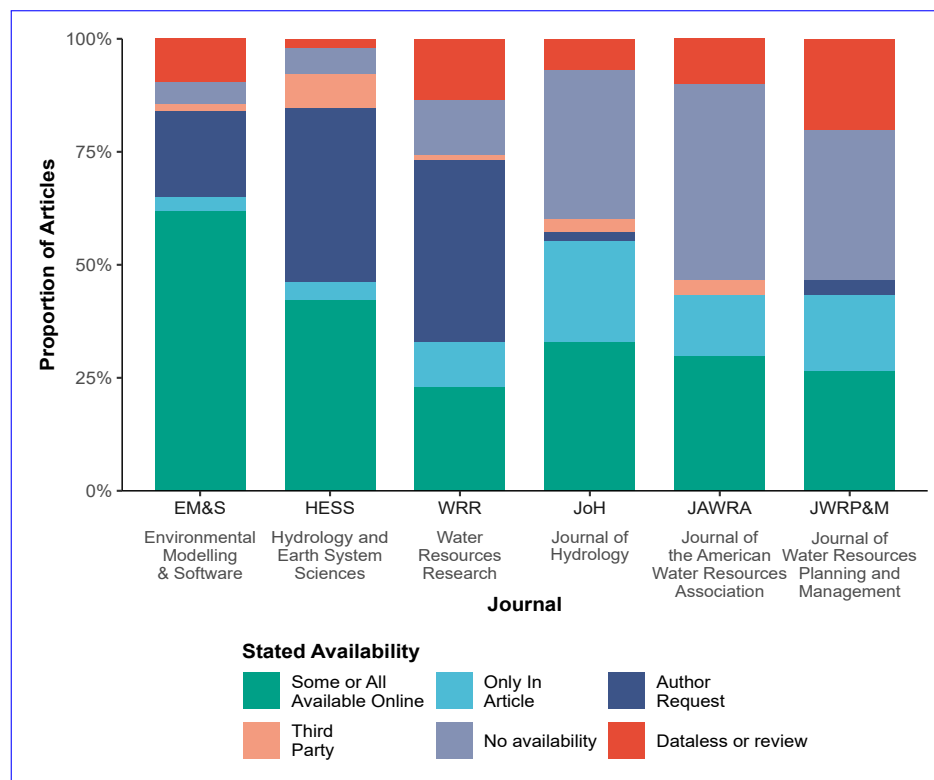


Figure 3: Data, model, code availability by journal (summary of Q4 and Q5).

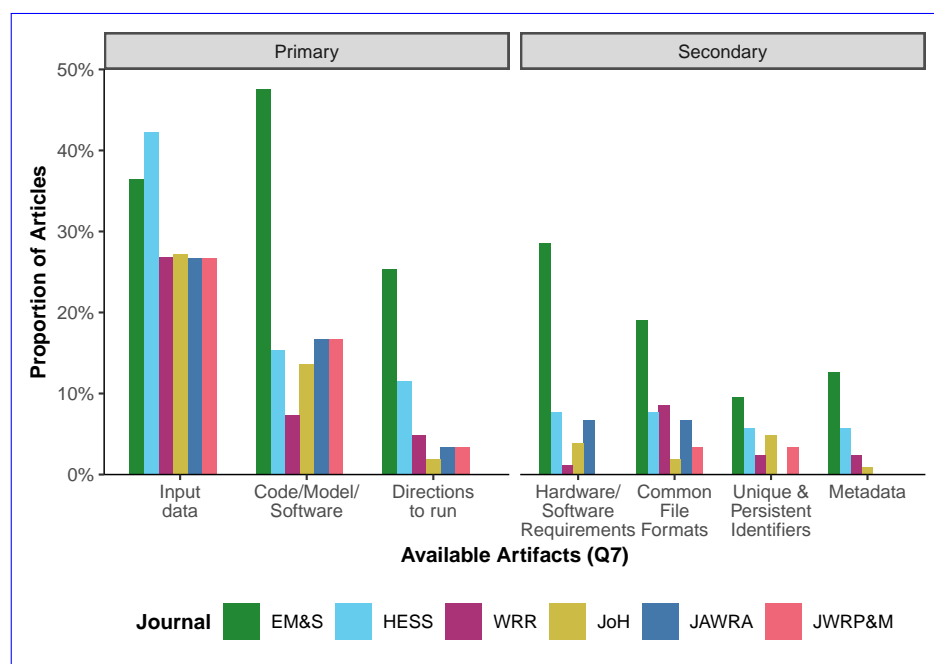


Figure 4: Availability artifacts organized by journal. All percentages are based on the total number of sampled papers for each journal. [Refer to Figure 3 or the text for full journal titles.](#)

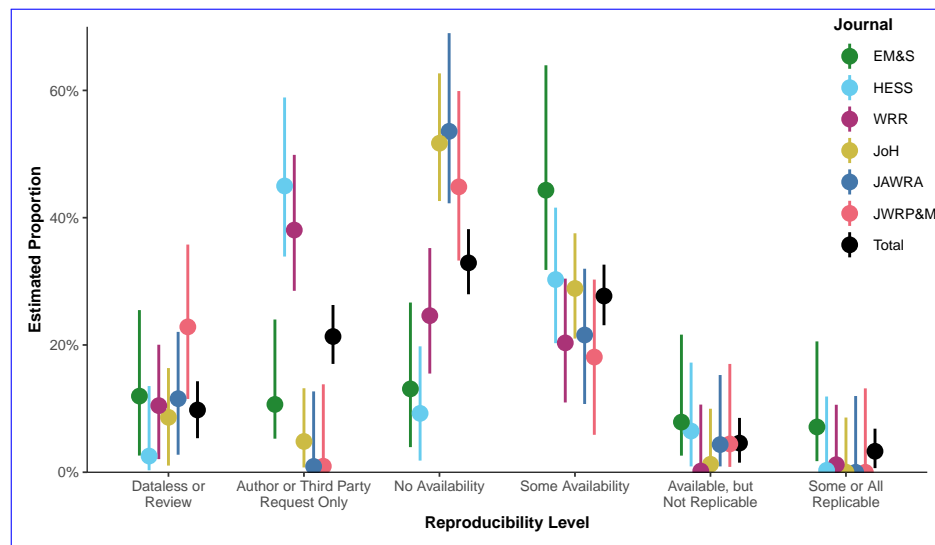


Figure 5: Population estimate of reproducibility for all papers published in 2017. Results are sorted by journal, with “Total” representing all 6 journals. Median estimate is represented by a point, vertical bars show the 95% confidence interval. [Refer to Figure 3 or the text for full journal titles.](#)

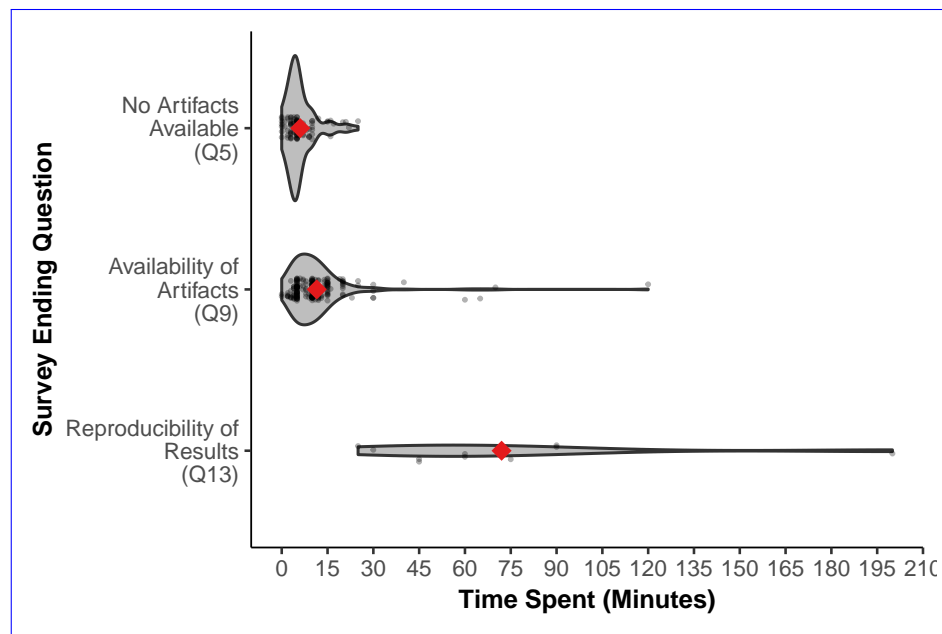


Figure 6: Self-reported time to complete survey organized by the survey's ending question. Each reviewed paper is shown by a ~~black~~-dot, while the mean is represented by a red ~~dot~~diamond. Distribution density is shown by width.

References

- [1] Sandve, G. K., Nekrutenko, A., Taylor, J. & Hovig, E. Ten Simple Rules for Reproducible Computational Research. *PLOS Computational Biology* **9**, e1003285 (2013).
- [2] Aarts, A. *et al.* Estimating the reproducibility of psychological science. *Science* **349**, 1–8 (2015).
- [3] Baker, M. 1,500 scientists lift the lid on reproducibility. *Nature* **533**, 452–454 (2016).
- [4] Gil, Y. *et al.* [Toward the geoscience paper of the future](#). *Earth and Space Science* **3**, 388–415 (2016).
- [5] Brembs, B. Prestigious Science Journals Struggle to Reach Even Average Reliability. *Frontiers in Human Neuroscience* **12** (2018).
- [6] Stodden, V., Seiler, J. & Ma, Z. An empirical analysis of journal policy effectiveness for computational reproducibility. *Proceedings of the National Academy of Sciences* **115**, 2584–2589 (2018).
- [7] Kovacevic, J. How to Encourage and Publish Reproducible Research. ~~in~~ *vol.*, 2007 IEEE International Conference on Acoustics, Speech and Signal Processing **4**, 1273–1276 (2007).
- [8] Peng, R. D. Reproducible research and Biostatistics. *Biostatistics* **10**, 405–408 (2009).
- [9] Stodden, V., Borwein, J. & Bailey, D. H. Setting the default to reproducible in computational science research. *SIAM News* **46**, 4–6 (2013).
- [10] Easterbrook, S. M. Open code for open science? *Nature Geoscience* **7**, 779–781 (2014).
- [11] Leek, J. T. & Peng, R. D. Opinion: Reproducible research can still be wrong: Adopting a prevention approach. *Proceedings of the National Academy of Sciences* **112**, 1645–1646 (2015).
- [12] Pulverer, B. Reproducibility blues. *The EMBO Journal* **34**, 2721–2724 (2015).
- [13] Goodman, S. N., Fanelli, D. & Ioannidis, J. P. A. What does research reproducibility mean? *Science Translational Medicine* **8**, 341ps12 (2016).
- [14] ~~Melsen, L. A.~~ Torfs, P. J. J. F., Uijlenhoet, R. & Teuling, A. J. Comment on “Most computational hydrology is not reproducible, so is it really science?” by Christopher Hutton *et al.* *Water Resources Research* **53**, 2568–2569 (2017).

- [15] Plesser, H. E. [Reproducibility vs. Replicability: A Brief History of a Confused Terminology](#). *Frontiers in Neuroinformatics* **11** (2018).
- [16] Institute of Education Sciences (IES), U.S. Department of Education & National Science Foundation (NSF). [Companion Guidelines on Replication & Reproducibility in Education Research: A Supplement to the Common Guidelines for Education Research and Development](#) <https://www.nsf.gov/pubs/2019/nsf19022/nsf19022.pdf> (2018).
- [17] Akmon, D., Zimmerman, A., Daniels, M. & Hedstrom, M. The application of archival concepts to a data-intensive environment: working with scientists to understand data management and preservation needs. *Archival Science* **11**, 329–348 (2011).
- [18] Hutton, C. *et al.* Most computational hydrology is not reproducible, so is it really science? *Water Resources Research* **52**, 7548–7555 (2016).
- [19] Añel, J. A. Comment on “Most computational hydrology is not reproducible, so is it really science?” by Christopher Hutton *et al.* *Water Resources Research* **53**, 2572–2574 (2017).
- [20] Casadevall, A. & Fang, F. C. Reproducible Science. *Infection and Immunity* **78**, 4972–4975 (2010).
- [21] Drummond, C. Reproducible research: a minority opinion. *Journal of Experimental & Theoretical Artificial Intelligence* **30**, 1–11 (2018).
- [22] Stodden, V. The Legal Framework for Reproducible Scientific Research: Licensing and Copyright. *Computing in Science & Engineering* **11**, 35–40 (2009).
- [23] Fary, M. & Owen, K. Developing an Institutional Research Data Management Plan Service *EDUCAUSE, ACTI DMWG–Advanced Core Technologies Initiative Data Management Working Group*. (2013).
- [24] Shen, Y. Research Data Sharing and Reuse Practices of Academic Faculty Researchers: A Study of the Virginia Tech Data Landscape. *International Journal of Digital Curation* **10**, 157–175 (2016).
- [25] Shiffrin, R. M., Börner, K. & Stigler, S. M. Scientific progress despite irreproducibility: A seeming paradox. *Proceedings of the National Academy of Sciences* **115**, 2632–2639 (2018).
- [26] Diekema, A., Wesolek, A. & Walters, C. The NSF/NIH Effect: Surveying the Effect of Data Management Requirements on Faculty, Sponsored Programs, and Institutional Repositories. *Data* (2014).
- [27] Wallis, J. C., Rolando, E. & Borgman, C. L. If We Share Data, Will Anyone Use Them? Data Sharing and Reuse in the Long Tail of Science and Technology. *PLoS ONE* **8**, e67332 (2013).

- [28] Kaufman, D. & PAGES 2k special-issue editorial team. Technical Note: Open-paleo-data implementation pilot – The PAGES 2k special issue. *Clim. Past Discuss.* **2017**, 1–10 (2017).
- [29] Gabridge, T. The Last Mile: Liaison Roles in Curating Science and Engineering Research Data. *Research Library Issues: A Bimonthly Report from ARL, CNI, and SPARC* **265**, 15–21 (2009).
- [30] Bracke, M. S. Emerging Data Curation Roles for Librarians: A Case Study of Agricultural Data. *Journal of Agricultural & Food Information* **12**, 65–74 (2011).
- [31] Pinfield, S., Cox, A. M. & Smith, J. Research Data Management and Libraries: Relationships, Activities, Drivers and Influences. *PLoS ONE* **9**, e114734 (2014).
- [32] Weller, T. & Monroe-Gulick, A. Differences in the Data Practices, Challenges, and Future Needs of Graduate Students and Faculty Members. *Journal of eScience Librarianship* **4**, 2 (2015).
- [33] Horsburgh, J. S. *et al.* HydroShare: Sharing Diverse Environmental Data Types and Models as Social Objects with Application to the Hydrology Domain. *JAWRA Journal of the American Water Resources Association* **52**, 873–889 (2016).
- [34] Essawy, B. T. *et al.* Integrating scientific cyberinfrastructures to improve reproducibility in computational hydrology: Example for HydroShare and GeoTrust. *Environmental Modelling & Software* **105**, 217–229 (2018).
- [35] Gillman, M. A., Lamoureux, S. F. & Lafrenière, M. J. Calibration of a modified temperature-light intensity logger for quantifying water electrical conductivity. *Water Resources Research* **53**, 8120–8126 (2017).
- [36] Horsburgh, J., Leonardo, M., Abdallah, A. & Rosenberg, D. Measuring water use, conservation, and differences by gender using an inexpensive, high frequency metering system. *Environmental Modelling and Software* **96**, 83–94 (2017).
- [37] Neuwirth, C. System dynamics simulations for data-intensive applications. *Environmental Modelling and Software* **96**, 140–145 (2017).
- [38] Quinn, J. *et al.* Detecting spatial patterns of rivermouth processes using a geostatistical framework for near-real-time analysis. *Environmental Modelling and Software* **97**, 72–85 (2017).
- [39] Buscombe, D. Shallow water benthic imaging and substrate characterization using recreational-grade sidescan-sonar. *Environmental Modelling and Software* **89**, 1–18 (2017).

- [40] Yu, C.-W., Liu, F. & Hodges, B. Consistent initial conditions for the Saint-Venant equations in river network modeling. *Hydrology and Earth System Sciences* **21**, 4959–4972 (2017).
- [41] Di Matteo, M., Dandy, G. & Maier, H. Multiobjective optimization of distributed stormwater harvesting systems. *Journal of Water Resources Planning and Management* **143**, 1–10 (2017).
- [42] Engdahl, N., Benson, D. & Bolster, D. Lagrangian simulation of mixing and reactions in complex geochemical systems. *Water Resources Research* **53**, 3513–3522 (2017).
- [43] Güntner, A. *et al.* Landscape-scale water balance monitoring with an iGrav superconducting gravimeter in a field enclosure. *Hydrology and Earth System Sciences* **21**, 3167–3182 (2017).
- [44] Sattar, A., Jasak, H. & Skuric, V. Three dimensional modeling of free surface flow and sediment transport with bed deformation using automatic mesh motion. *Environmental Modelling and Software* **97**, 303–317 (2017).
- [45] Nosek, B.A. *et al.* Promoting an open research culture. *Science* **348**, 1422–1425 (2015).
- [46] Wilkinson, M.D. *et al.* A design framework and exemplar metrics for FAIRness. *Scientific Data* **5**, 180118 (2018).
- [47] Rosenberg, D. E. & Watkins, D. W. New Policy to Specify Availability of Data, Models, and Code. *Journal of Water Resources Planning and Management* **144**, 01618001 (2018).
- [48] Collberg, C. *et al.* Measuring reproducibility in computer systems research. *University of Arizona, Tech. Rep 37* (2014).
- [49] Kidwell, M.C. *et al.* Badges to Acknowledge Open Practices: A Simple, Low-Cost, Effective Method for Increasing Transparency. *PLoS Biology* **14**, e1002456 (2016).
- [50] Thaler, R. H. & Sunstein, C. R. *Nudge: Improving decisions about health, wealth, and happiness* ~~–Nudge: Improving decisions about health, wealth, and happiness–~~ (Yale University Press, New Haven, CT, US, 2008).
- [51] Sison, C. P. & Glaz, J. Simultaneous Confidence Intervals and Sample Size Determination for Multinomial Proportions. *Journal of the American Statistical Association* **90**, 366–369 (1995).
- [52] May, W. L. & Johnson, W. D. Constructing two-sided simultaneous confidence intervals for multinomial proportions for small counts in a large number of cells. *Journal of Statistical Software* **5**, 1–24 (2000).
- ~~–&–~~

- 550 [53] Stagge, J., Abdallah, A. & Rosenberg, D. `—jstagge/reproduc_hyd`: Source
551 code accompanying A survey tool to assess and improve data availabil-
552 ity and research reproducibility. <https://zenodo.org/record/1467693>
553 (2018).