

1 **Advice for improving the reproducibility of data extraction in**  
2 **meta-analysis**

3 Edward R. Ivimey-Cook<sup>1\*</sup>, Daniel W. A. Noble<sup>2</sup>, Shinichi Nakagawa<sup>3</sup>, Marc J.  
4 Lajeunesse<sup>4</sup> & Joel L. Pick<sup>5</sup>

5 <sup>1</sup> School of Biodiversity, One Health and Veterinary Medicine, University of Glasgow,  
6 Glasgow, G12 8QQ, UK

7 <sup>2</sup> Division of Ecology and Evolution, Research School of Biology, The Australian National  
8 University, Canberra ACT 2600

9 <sup>3</sup> Evolution and Ecology Research Centre, School of Biological, Earth and Environmental  
10 Sciences, The University of New South Wales, Sydney, NSW 2052

11 <sup>4</sup> Department of Integrative Biology, University of South Florida, Tampa, FL USA

12 <sup>5</sup> Institute of Ecology and Evolution, University of Edinburgh, Charlotte Auerbach Road,  
13 Edinburgh, EH9 3FL, UK

14 \* Corresponding author, email address: e.ivimeycook@googlemail.com

15 **Running title:** Reproducibility in meta-analysis

16 **Keywords:** reproducibility, meta-analysis, data extraction, metaDigitise, juicr, shiny-  
17 Digitise

## 18 Abstract

19 Extracting data from studies is the norm in meta-analyses, enabling researchers to gener-  
20 ate effect sizes when raw data are otherwise not available. While there has been a general  
21 push for increased reproducibility in meta-analysis, the transparency and reproducibility  
22 of the data extraction phase is still lagging behind. Unfortunately, there is little guidance  
23 of how to make this process more transparent and shareable. To address this, we pro-  
24 vide several steps to help increase the reproducibility of data extraction in meta-analysis.  
25 We also provide suggestions of R software that can further help with reproducible data  
26 policies: the *shinyDigitise* and *juicer* packages. Adopting the guiding principles listed  
27 here and using the appropriate software will provide a more transparent form of data  
28 extraction in meta-analyses.

## 29 Introduction

30 In recent years, there has been a push to increase the reproducibility of meta-analyses  
31 (the ability to recreate the same findings if the same project was reconducted; see [Ihle](#)  
32 [et al., 2017](#)), with the expectation that exact search strings, screening steps (e.g. the  
33 PRISMA flowchart: [Moher et al., 2011](#); [O’Dea et al., 2021](#)), and metadata of accepted  
34 papers are included alongside manuscripts. However, unlike the study selection process,  
35 the reproducibility of steps taken during data extraction is typically overlooked, and  
36 no unified reporting guidelines currently exist. Indeed, several papers have highlighted  
37 the prevalence of errors in meta-analysis, particularly surrounding the data extraction  
38 process ([Gøtzsche et al., 2007](#); [Mathes et al., 2017](#); [Wong & Bouchard, 2022](#)). As a  
39 result, if studies provide neither the data needed to reproduce the analysis nor the source  
40 of the effect size within the screened study (e.g. in text, table or figure, reporting of which  
41 is typically low; see [O’Dea et al., 2021](#)), then there can be a lack of repeatability, where  
42 independent screeners are unable to locate and extract the same values (see [Buscemi](#)

43 *et al.*, 2006). Altogether, this suggests that this vital stage of the meta-analysis workflow  
44 lacks both transparency and, importantly, reproducibility.

45 Here, to assess the extent of problems with data extraction reporting, we review  
46 the current state of the literature. Firstly, we review the evidence of reporting of data  
47 extraction software in recent meta-analyses in Ecology and Evolution. Secondly, we  
48 investigate the reporting practices of papers that have cited the R package *metaDigitise* as  
49 a case study. We then introduce a simple five-step guide to help improve the replicability  
50 and reproducibility of data extraction. We note that this will not reduce user-specific  
51 errors made during the data extraction process, but will enable a higher probability  
52 of detecting and correcting any errors made. Finally, we introduce two R-based GUI  
53 packages, *shinyDigitise* and *juicer*, which have both been designed to aid transparency  
54 and reproducibility.

## 55 State of the field

56 To start, we quantified the percentage of meta-analyses that reported any software pack-  
57 ages used to extract data from figures. To do this, Y.Y., M.L., and J.R. re-examined the  
58 102 meta-analyses reviewed in the 2021 PRISMA-EcoEvo guidelines paper (O’Dea *et al.*,  
59 2021). From these 102 studies, only 39 cited the data extraction software that was used  
60 to extract data from figures (representing 38% of the total number). We note that whilst  
61 this survey and results focus on meta-analysis within the fields of Ecology and Evolution,  
62 no such survey has yet been conducted in other disciplines despite the common nature of  
63 figure-based data extraction.

64 Next, to assess transparency of the data extraction process itself, E.I-C reviewed  
65 all studies listed as citing the R package *metaDigitise* (Pick *et al.*, 2019) in August  
66 2022 (for full methodology, see SM1). The *metaDigitise* package (on CRAN in 2018,  
67 associated paper published in 2019) was in part designed to help improve transparency

68 and reproducibility of data extraction (Pick *et al.*, 2019). It provides a simple way of  
69 storing figures and associated extraction data which can easily be uploaded as part of  
70 the data archiving process. Papers citing *metaDigitise* therefore provide a good insight  
71 into the transparency of data-extraction reporting in recently published meta-analyses.  
72 In total, 55 published meta-analyses were obtained that covered several subject areas,  
73 ecology and evolution, medicine, environmental science, and psychology.

74 The results of this survey are shown in Figure S1. 78% of the 55 meta-analyses  
75 using metaDigitise (n = 43) had available data in an interoperable format, despite the  
76 open access policy of many journals and increased awareness of the importance of open-  
77 data. From these, only 24 (44% of the total) readily provided information about the  
78 origin of the effect sizes which is in line with the 39% reported from a recent survey  
79 in ecology and evolution meta-analyses (O’Dea *et al.*, 2021). Of these studies between  
80 2-96% (median = 28%) of all effect sizes were generated from figures. Finally, only four  
81 studies (7%) provided the figures from which data was extracted and only two provided  
82 the calibration data needed to recreate the extraction (5%) in addition to the figure and  
83 metadata required to reproduce the analysis (Figure S1). The low reporting rates are  
84 even more extreme when one considers only 38% of meta-analyses reviewed by O’Dea  
85 *et al.* (2021) reported the software used to extract these effect sizes from figures.

## 86 Advice for data extraction

87 Based on this survey it is clear that we need to improve the transparency and repro-  
88 ducibility of data extraction in meta-analyses. To achieve this, we introduce a simple  
89 five-step guide.

90 1. **Provide data.** As discussed at length elsewhere (Miyakawa, 2020), providing data  
91 is a minimum requirement for reproducibility. We found that 78% of meta-analyses  
92 provide data, similar to the 77% in a recent review of ecology and evolution meta-

93 analyses (2010 - 2019; O’Dea *et al.*, 2021). Although this shows an improvement  
94 over the last decade (from 31% shared data in Ecology meta-analyses between  
95 the years 1996–2013 Koricheva & Gurevitch, 2014), and is substantially greater  
96 than in other fields (e.g. 3% of studies provided interoperable data in clinical  
97 psychological meta-analyses from 2000–2020; López-Nicolás *et al.*, 2022), data in  
98 meta-analysis typically come from open sources (i.e. published literature) and so  
99 there are few obvious reasons why data should not be made public. Meta-analysts  
100 should therefore be expected to lead by example and provide their own data.

101 **2. Clearly state where each effect size was extracted.** In addition to providing  
102 other relevant metadata, it should be clearly stated where effect sizes were extracted  
103 from (e.g. text, table, figure or supplementary material), including a reference to the  
104 exact location, e.g. "Figure 2a", "table 3", "main text p275". Curtis *et al.* (2013)  
105 suggested a shorthand for reporting this information in tabular form (e.g., F2a, T3),  
106 and we extend this formatting to T=table, M=main text, F=figure, A=appendix,  
107 S=supplementary material, R=raw data, followed by the figure and/or page num-  
108 ber where the data was extracted. In addition to providing copies of the extracted  
109 figures, uploading a screenshot or section of PDF which clearly highlights the loca-  
110 tion of the extracted effect size would be useful, particularly when considering data  
111 in text or in table (although note the caveats listed below). Lastly, under some  
112 circumstances, data might be provided from unpublished studies through personal  
113 contact with authors. In this case, it is still important to provide a location of where  
114 or how the effect size was obtained (i.e. personal communication or unpublished  
115 data), in order to allow for others to similarly obtain the data.

116 **3. Provide transformation information.** Providing effect sizes alone does not  
117 give information on how they have been generated. For example, transformations  
118 have to be used to generate means and standard deviations from the quantiles in  
119 a boxplot (e.g. Wan *et al.*, 2014). Other transformations include converting stan-

120 dard errors (SE) to standard deviations (SD), or calibrations of extracted data by  
121 back-transforming logarithms. Generating effect sizes from figures always requires  
122 additional steps in order to make them usable in meta-analysis. These details are  
123 more challenging to report succinctly, as they may require equations, but a textual  
124 description alongside raw data and code is better than nothing. Indeed, [O'Dea \*et al.\*](#)  
125 [\(2021\)](#) showed that only 39% of papers provided the raw data used to generate  
126 effect sizes, compared with the 77% that provided processed effect sizes.

127 **4. Provide figures alongside a record of the data extraction process.** A con-  
128 siderable amount of data for meta-analysis comes from figures (e.g., in the above  
129 survey, 28% of effect sizes, on average, originated from figures). Therefore, every  
130 figure that has undergone data extraction should be provided in a digital data repos-  
131 itory (e.g. Open Science Framework, Zenodo, or Dryad) alongside the generated  
132 effect size. Data extraction files including calibration data are also needed for any  
133 researcher to be able to recreate and check the extraction process. Importantly, it  
134 is also worth considering (and noting in the metadata) whether the source paper  
135 was open or non-open access. Whilst a breach of copyright may not be an issue  
136 with figures from open access papers, this could be a potential problem with non-  
137 open access papers. In this case, we suggest three actions: 1) note in the metadata  
138 which figures might be restricted due to copyright infringement; 2) seek permission  
139 from the journal and/or author of the paper; 3) store all of the figures on a private  
140 repository (such as those listed above) which can be made available upon request.  
141 It is also a requirement, regardless of whether the paper is open or non-open access,  
142 to appropriately cite the primary literature where the figure has been obtained.

143 **5. For software developers, enable the saving and reloading of the data**  
144 **extraction process.** Whilst there exists a multitude of data extraction tools,  
145 few allow users to easily save and reload the data extraction process. Therefore, to  
146 increase reproducibility, the development of new tools or software for data extraction

147 should ensure this functionality. The file format of extractions should be also tool  
148 agnostic with a format accessible to all (interoperable; e.g., a .csv file).

## 149 **Tools for increasing reproducibility in figure-based ex-** 150 **traction**

151 Here we highlight two R-based packages that are being developed that allow for repro-  
152 ducible figure-based data extraction. Firstly, *shinyDigitise*, a GUI for the *metaDigitise*  
153 (Pick *et al.*, 2019) package, and secondly, *juicr* (Lajeunesse, 2021). We focus on these  
154 packages because R is one of the most widely used statistical environments for analysing  
155 meta-analytic data. We note that whilst these packages should be suitable for extraction  
156 of many of the commonly used figures across disciplines (scatterplots, mean-error plots,  
157 boxplots, and histograms), they may not be as well equipped to extract data from highly  
158 specialised domain-specific figures.

159 *shinyDigitise* (developed by E.I-C & J.L.P) is a streamlined and intuitive GUI inter-  
160 face which is built upon the functions of the *metaDigitise* package (Pick *et al.*, 2019). This  
161 includes the ability to extract data from a wide variety of plot-types (scatterplots, mean-  
162 error plots, boxplots, and histograms), and automatically saves calibration data so users  
163 have a historical record of the data extraction process. *shinyDigitise* should reduce the  
164 barrier of entry by requiring very little experience of writing code or the R coding software.  
165 To install this package, see the GitHub: <https://github.com/EIvimeyCook/shinyDigitise>.

166 Alongside *shinyDigitise*, *juicr* (developed by M.J.L.) offers savable and shareable  
167 records of retrieved data from images. *juicr* offers a point-and-click solution to extract-  
168 ing data from images; however for some tasks, decision-making of what to extract can  
169 be delegated to automated (full algorithmic) or semi-automated (algorithmic with user  
170 assistance) tools. The *juicr* package extends the automated extraction tools first devel-  
171 oped in the *metagear* package for research synthesis (Lajeunesse, 2016); to install this

172 package, see the GitHub: <https://github.com/mjlajeunesse/juicr>.

173 Importantly, these software packages provide the user with an effect size *in addition*  
174 to a record of the extraction process for each figure. After depositing into an appropriate  
175 data repository, these can be subsequently viewed and error checked by the user or by  
176 anyone with access to both the figure and record files. Whilst this is an important step  
177 for reproducibility, and directly adheres to step four above, very few people have adopted  
178 the use of this archiving functionality. Figure S1 highlights the low percentage of studies  
179 that share source figures, their extracted data, and information as to when and what  
180 extraction software tool was used, in addition to providing records of the data extraction  
181 process. Clearly there is an urgent need to increase transparency of data extraction, and  
182 the steps outlined above should go some way to addressing this.

## 183 **Acknowledgements**

184 We are grateful to Malgorzata Lagisz, Yefeng Yang and Joanna Rutkowska who surveyed  
185 the 102 meta-analyses as a part of a larger project. We also thank two anonymous  
186 reviewers for helpful comments on the manuscript. Lastly, we thank Stuart Taylor from  
187 the Royal Society for guidance on issues with copyright and non-open access papers. Note,  
188 we refer to authors in text using the MeRIT system (Method Reporting with Initials for  
189 Transparency) as per (Nakagawa *et al.*, 2023).

## 190 **Data Availability Statement**

191 The data that support the findings of this study are openly available on GitHub:  
192 <https://github.com/EIvimeyCook/DataExtraction> or Zenodo: 10.5281/zenodo.8187175



## 193 Conflict of Interest

194 The authors declare no conflict of interest.

## 195 Highlights

- 196 • In meta-analysis, large quantities of data need to be extracted from published liter-  
197 ature. However, the transparency and reproducibility of the data extraction process  
198 is often limited, both in terms of its description in the methods section and also  
199 when data is later uploaded to an open data repository.
- 200 • In order to increase the reproducibility of data extraction in meta-analysis, we  
201 introduce a simple five-step guide which includes suggestions for future research.  
202 Furthermore, we highlight two packages in R that readily facilitate reproducible  
203 workflows and allow for shareable records of the data extraction process.
- 204 • Adopting the principles and suggestions provided here will help to make the entire  
205 meta-analysis process more transparent, open, and reproducible.

## 206 References

- 207 Buscemi, N., Hartling, L., Vandermeer, B., Tjosvold, L. & Klassen, T.P.  
208 (2006) Single data extraction generated more errors than double data extrac-  
209 tion in systematic reviews. *Journal of Clinical Epidemiology*, **59**, 697–703.  
210 <https://dx.doi.org/https://doi.org/10.1016/j.jclinepi.2005.11.010>.
- 211 Curtis, P.S., Mengersen, K., Lajeunesse, M.J., Rothstein, H.R. & Stewart, G.B. (2013)  
212 Extraction and critical appraisal of data. J. Koricheva, J. Gurevitch & K. Mengersen,  
213 eds., *Handbook of meta-analysis in ecology and evolution*, pp. 52–60. Princeton Univer-  
214 sity Press, Princeton, New Jersey, USA.

- 215 Gøtzsche, P.C., Hróbjartsson, A., Marić, K. & Tendal, B. (2007) Data Extraction  
216 Errors in Meta-analyses That Use Standardized Mean Differences. *JAMA*, **298**.  
217 <https://dx.doi.org/10.1001/jama.298.4.430>.
- 218 Ihle, M., Winney, I.S., Krystalli, A. & Croucher, M. (2017) Striving for transparent and  
219 credible research: practical guidelines for behavioral ecologists. *Behavioral Ecology*,  
220 **28**, 348–354.
- 221 Koricheva, J. & Gurevitch, J. (2014) Uses and misuses of meta-analysis in plant ecology.  
222 *Journal of Ecology*, **102**, 828–844. ISBN: 1365-2745, [https://dx.doi.org/10.1111/1365-](https://dx.doi.org/10.1111/1365-2745.12224)  
223 [2745.12224](https://dx.doi.org/10.1111/1365-2745.12224).
- 224 Lajeunesse, M.J. (2016) Facilitating systematic reviews, data extraction and meta-  
225 analysis with the metagear package for R. *Methods in Ecology and Evolution*, **7**,  
226 323–330.
- 227 Lajeunesse, M.J. (2021) Squeezing data from scientific images using the juicr package for  
228 R. R package version 0.1.
- 229 López-Nicolás, R., López-López, J.A., Rubio-Aparicio, M. & Sánchez-Meca, J. (2022)  
230 A meta-review of transparency and reproducibility-related reporting practices in pub-  
231 lished meta-analyses on clinical psychological interventions (2000–2020). *Behavior Re-*  
232 *search Methods*, **54**, 334–349. <https://dx.doi.org/10.3758/s13428-021-01644-z>.
- 233 Mathes, T., Klaffen, P. & Pieper, D. (2017) Frequency of data extraction errors and  
234 methods to increase data extraction quality: a methodological review. *BMC Medical*  
235 *Research Methodology*, **17**, 152. <https://dx.doi.org/10.1186/s12874-017-0431-4>.
- 236 Miyakawa, T. (2020) No raw data, no science: another possible source of the reproducibil-  
237 ity crisis. *Molecular Brain*, **13**, 24. <https://dx.doi.org/10.1186/s13041-020-0552-2>.
- 238 Moher, D., Altman, D.G., Liberati, A. & Tetzlaff, J. (2011) PRISMA statement. *Epi-*  
239 *demiology*, **22**, 128.

240 Nakagawa, S., Ivimey-Cook, E.R., Grainger, M.J., O’Dea, R.E., Burke, S., Drobniak,  
241 S.M., Gould, E., Macartney, E.L., Martnig, A.R., Paquet, M., Morrison, K., Pick, J.L.,  
242 Pottier, P., Ricolfi, L., Wilkinson, D.P., Willcox, A., Williams, C., Wilson, L.A.B.,  
243 Windecker, S.M., Yang, Y. & Lagisz., M. (2023) Method reporting with initials for  
244 transparency (merit) promotes more granularity and accountability for author contri-  
245 butions. *Nature Communications*, **14**, 1788. [https://dx.doi.org/10.1038/s41467-023-](https://dx.doi.org/10.1038/s41467-023-37039-1)  
246 [37039-1](https://dx.doi.org/10.1038/s41467-023-37039-1).

247 O’Dea, R.E., Lagisz, M., Jennions, M.D., Koricheva, J., Noble, D.W., Parker, T.H.,  
248 Gurevitch, J., Page, M.J., Stewart, G., Moher, D. & Nakagawa, S. (2021) Pre-  
249 ferred reporting items for systematic reviews and meta-analyses in ecology and  
250 evolutionary biology: a PRISMA extension. *Biological Reviews*, **96**, 1695–1722.  
251 <https://dx.doi.org/10.1111/brv.12721>.

252 Pick, J.L., Nakagawa, S. & Noble, D.W.A. (2019) Reproducible, flexible and high-  
253 throughput data extraction from primary literature: The metaDigitise r package.  
254 *Methods in Ecology and Evolution*, **10**, 426–431. [https://dx.doi.org/10.1111/2041-](https://dx.doi.org/10.1111/2041-210X.13118)  
255 [210X.13118](https://dx.doi.org/10.1111/2041-210X.13118).

256 Wan, X., Wang, W., Liu, J. & Tong, T. (2014) Estimating the sample mean and standard  
257 deviation from the sample size, median, range and/or interquartile range. *BMC medical*  
258 *research methodology*, **14**, 135. ISBN: 1471-2288 (Electronic)\r1471-2288 (Linking),  
259 <https://dx.doi.org/10.1186/1471-2288-14-135>.

260 Wong, J.S. & Bouchard, J. (2022) Do Meta-Analyses of Intervention/Prevention  
261 Programs in the Field of Criminology Meet the Tests of Transparency and  
262 Reproducibility? *Trauma, Violence, & Abuse*, p. 15248380211073839.  
263 <https://dx.doi.org/10.1177/15248380211073839>.