

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

Edward R. Ivimey-Cook¹  | Daniel W. A. Noble²  | Shinichi Nakagawa³  |
Marc J. Lajeunesse⁴  | Joel L. Pick⁵ 

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

²Division of Ecology and Evolution, Research School of Biology, The Australian National University, Canberra, Australian Capital Territory, Australia

³Evolution and Ecology Research Centre, School of Biological, Earth and Environmental Sciences, The University of New South Wales, Sydney, New South Wales, Australia

⁴Department of Integrative Biology, University of South Florida, Tampa, Florida, USA

⁵Institute of Ecology and Evolution, University of Edinburgh, Edinburgh, UK

Correspondence

Edward R. Ivimey-Cook, School of Biodiversity, One Health and Veterinary Medicine, University of Glasgow, Glasgow, G12 8QQ, UK.
Email: e.ivimeycook@googlegmail.com

Abstract

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

KEYWORDS

data extraction, juicr, meta-analysis, metaDigitise, reproducibility, shinyDigitise

Highlights

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

1 | INTRODUCTION

In recent years, there has been a push to increase the reproducibility of meta-analyses (the ability to recreate the

same findings if the same project was reconducted¹), with the expectation that exact search strings, screening steps (e.g., the PRISMA flowchart^{2,3}), and metadata of accepted papers are included alongside manuscripts. However,

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2023 The Authors. *Research Synthesis Methods* published by John Wiley & Sons Ltd.

unlike the study selection process, the reproducibility of steps taken during data extraction is typically overlooked, and no unified reporting guidelines currently exist. Indeed, several papers have highlighted the prevalence of errors in meta-analysis, particularly surrounding the data extraction process.^{4–6} As a result, if studies provide neither the data needed to reproduce the analysis nor the source of the effect size within the screened study (e.g., in text, table or figure, reporting of which is typically low; see³), then there can be a lack of repeatability, where independent screeners are unable to locate and extract the same values (see⁷). Altogether, this suggests that this vital stage of the meta-analysis workflow lacks both transparency and, importantly, reproducibility.

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

2 | STATE OF THE FIELD

To start, we quantified the percentage of meta-analyses that reported any software packages used to extract data from figures. To do this, Y.Y., M.L., and J.R. re-examined the 102 meta-analyses reviewed in the 2021 PRISMA-EcoEvo guidelines paper.³ From these 102 studies, only 39 cited the data extraction software that was used to extract data from figures (representing 38% of the total number). We note that while this survey and results focus on meta-analysis within the fields of Ecology and Evolution, no such survey has yet been conducted in other disciplines despite the common nature of figure-based data extraction.

Next, to assess transparency of the data extraction process itself, E.I-C reviewed all studies listed as citing the R package *metaDigitise*⁸ in August 2022 (for full methodology, see SM1). The *metaDigitise* package (on CRAN in 2018, associated paper published in 2019) was in part designed to help improve transparency and reproducibility of data extraction.⁸ It provides a simple way of storing figures and associated extraction data

which can easily be uploaded as part of the data archiving process. Papers citing *metaDigitise*, therefore, provide good insight into the transparency of data-extraction and reporting in recently published meta-analyses. In total, 55 published meta-analyses were obtained that covered several subject areas, ecology and evolution, medicine, environmental science, and psychology.

The results of this survey are shown in Figure S1. 78% of the 55 meta-analyses using *metaDigitise* ($n = 43$) had available data in an interoperable format, despite the open access policy of many journals and increased awareness of the importance of open-data. From these, only 24 (44% of the total) readily provided information about the origin of the effect sizes which is in line with the 39% reported from a recent survey in ecology and evolution meta-analyses.³ Of these studies between 2% and 96% (median = 28%) of all effect sizes were generated from figures. Finally, only four studies (7%) provided the figures from which data were extracted and only two provided the calibration data needed to recreate the extraction (5%) in addition to the figure and metadata required to reproduce the analysis (Figure S1). The low reporting rates are even more extreme when one considers only 38% of meta-analyses reviewed by O'Dea et al.³ reported the software used to extract these effect sizes from figures.

3 | ADVICE FOR DATA EXTRACTION

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

1. *Provide data.* As discussed at length elsewhere,⁹ providing data is a minimum requirement for reproducibility. We found that 78% of meta-analyses provide data, similar to the 77% in a recent review of ecology and evolution meta-analyses (2010–2019³). Although this shows an improvement over the last decade (from 31% shared data in Ecology meta-analyses between the years 1996–2013¹⁰), and is substantially greater than in other fields (e.g., 3% of studies provided interoperable data in clinical psychological meta-analyses from 2000 to 2020¹¹), data in meta-analysis typically come from open sources (i.e., published literature) and so there are few obvious reasons why data should not be made public. Meta-analysts should, therefore, be expected to lead by example and provide their own data.
2. *Clearly state where each effect size was extracted.* In addition to providing other relevant metadata, it

should be clearly stated where effect sizes were extracted from (e.g., text, table, figure or supplementary material), including a reference to the exact location, for example, “fig. 2a,” “tab. 3,” “main text p275.” Curtis et al.¹² suggested a shorthand for reporting this information in tabular form (e.g., F2a, T3), and we extend this formatting to T = table, M = main text, F = figure, A = appendix, S = supplementary material, R = raw data, followed by the figure and/or page number where the data was extracted. In addition to providing copies of the extracted figures, uploading a screenshot or section of PDF which clearly highlights the location of the extracted effect size would be useful, particularly when considering data in text or in table (although note the caveats listed below). Lastly, under some circumstances, data might be provided from unpublished studies through personal contact with authors. In this case, it is still important to provide a location of where or how the effect size was obtained (i.e. personal communication or unpublished data), in order to allow for others to similarly obtain the data.

3. *Provide transformation information.* Providing effect sizes alone does not give information on how they have been generated. For example, transformations have to be used to generate means and standard deviations from the quantiles in a boxplot (e.g.,¹³). Other transformations include converting standard errors (SE) to standard deviations (SD), or calibrations of extracted data by back-transforming logarithms. Generating effect sizes from figures always requires additional steps in order to make them usable in meta-analysis. These details are more challenging to report succinctly, as they may require equations, but a textual description alongside raw data and code is better than nothing. Indeed, O’Dea et al.³ showed that only 39% of papers provided the raw data used to generate effect sizes, compared with the 77% that provided processed effect sizes.
4. *Provide figures alongside a record of the data extraction process.* A considerable amount of data for meta-analysis comes from figures (e.g., in the above survey, 28% of effect sizes, on average, originated from figures). Therefore, every figure that has undergone data extraction should be provided in a digital data repository (e.g., Open Science Framework, Zenodo, or Dryad) alongside the generated effect size. Data extraction files including calibration data are also needed for any researcher to be able to recreate and check the extraction process. Importantly, it is also worth considering (and noting in the metadata) whether the source paper was open or non-open access. While a breach of copyright may not be an

issue with figures from open access papers, this could be a potential problem with non-open access papers. In this case, we suggest three actions: (1) note in the metadata which figures might be restricted due to copyright infringement; (2) seek permission from the journal and/or author of the paper; (3) store all of the figures on a private repository (such as those listed above) which can be made available upon request. It is also a requirement, regardless of whether the paper is open or non-open access, to appropriately cite the primary literature where the figure has been obtained.

5. *For software developers, enable the saving and reloading of the data extraction process.* While there exists a multitude of data extraction tools, few allow users to easily save and reload the data extraction process. Therefore, to increase reproducibility, the development of new tools or software for data extraction should ensure this functionality. The file format of extractions should be also tool agnostic with a format accessible to all (interoperable; e.g., a .csv file).

4 | TOOLS FOR INCREASING REPRODUCIBILITY IN FIGURE-BASED EXTRACTION

Here, we highlight two R-based packages that are being developed that allow for reproducible figure-based data extraction. Firstly, *shinyDigitise*, a GUI for the *metaDigitise*⁸ package, and secondly, *juicr*.¹⁴ We focus on these packages because R is one of the most widely used statistical environments for analysing meta-analytic data. We note that while these packages should be suitable for extraction of many of the commonly used figures across disciplines (scatterplots, mean-error plots, boxplots, and histograms), they may not be as well equipped to extract data from highly specialised domain-specific figures.

shinyDigitise (developed by E.I-C & J.L.P) is a streamlined and intuitive GUI interface which is built upon the functions of the *metaDigitise* package.⁸ This includes the ability to extract data from a wide variety of plot-types (scatterplots, mean-error plots, boxplots, and histograms), and automatically saves calibration data so users have a historical record of the data extraction process. *shinyDigitise* should reduce the barrier of entry by requiring very little experience of writing code or the R coding software. To install this package, see the GitHub: <https://github.com/ElvimeyCook/shinyDigitise>.

Alongside *shinyDigitise*, *juicr* (developed by M.J.L.) offers savable and shareable records of retrieved data from images. *juicr* offers a point-and-click solution to extracting data from images; however, for some tasks,

decision-making of what to extract can be delegated to automated (full algorithmic) or semi-automated (algorithmic with user assistance) tools. The *juicr* package extends the automated extraction tools first developed in the *metagear* package for research synthesis¹⁵; to install this package, see the GitHub: <https://github.com/mjlajeunesse/juicr>.

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

AUTHOR CONTRIBUTIONS

Edward R. Ivimey-Cook: Conceptualization; writing – original draft; writing – review and editing; methodology; data curation; formal analysis. **Daniel W. A. Noble:** Writing – review and editing; formal analysis; supervision; writing – original draft; software; methodology. **Shinichi Nakagawa:** Writing – review and editing; formal analysis; supervision; writing – original draft; software; methodology. **Marc J. Lajeunesse:** Software; writing – review and editing. **Joel L. Pick:** Conceptualization; writing – original draft; writing – review and editing; methodology; supervision; software.

ACKNOWLEDGMENTS

We are grateful to Malgorzata Lagisz, Yefeng Yang and Joanna Rutkowska who surveyed the 102 meta-analyses as a part of a larger project. We also thank two anonymous reviewers for helpful comments on the manuscript. Lastly, we thank Stuart Taylor from the Royal Society for guidance on issues with copyright and non-open access papers. Note, we refer to authors in text using the MeRIT system (Method Reporting with Initials for Transparency) as per.¹⁶

CONFLICT OF INTEREST STATEMENT

The authors declare no conflict of interest.

DATA AVAILABILITY STATEMENT

The data that support the findings of this study are openly available on GitHub: <https://github.com/EIvimeyCook/>

Data Extraction or Zenodo: <https://doi.org/10.5281/zenodo.8187175>.

ORCID

Edward R. Ivimey-Cook  <https://orcid.org/0000-0003-4910-0443>

Daniel W. A. Noble  <https://orcid.org/0000-0001-9460-8743>

Shinichi Nakagawa  <https://orcid.org/0000-0002-7765-5182>

Marc J. Lajeunesse  <https://orcid.org/0000-0002-9678-2080>

Joel L. Pick  <https://orcid.org/0000-0002-6295-3742>

REFERENCES

1. Ihle M, Winney IS, Krystalli A, Croucher M. Striving for transparent and credible research: practical guidelines for behavioural ecologists. *Behav Ecol*. 2017;28:348-354.
2. Moher D, Altman DG, Liberati A, Tetzlaff J. PRISMA statement. *Epidemiology*. 2011;22:128.
3. O'Dea RE, Lagisz M, Jennions MD, et al. Preferred reporting items for systematic reviews and meta-analyses in ecology and evolutionary biology: a PRISMA extension. *Biol Rev*. 2021;96:1695-1722. doi:10.1111/brv.12721
4. Gøtzsche PC, Hróbjartsson A, Marić K, Tendal B. Data extraction errors in meta-analyses that use standardized mean differences. *JAMA*. 2007;298:430-437. doi:10.1001/jama.298.4.430
5. Mathes T, Klaffen P, Pieper D. Frequency of data extraction errors and methods to increase data extraction quality: a methodological review. *BMC Med Res Methodol*. 2017;17:152. doi:10.1186/s12874-017-0431-4
6. Wong JS, Bouchard J. Do meta-analyses of intervention/prevention programs in the field of criminology meet the tests of transparency and reproducibility? *Trauma Violence Abuse*. 2022;24:15248380211073839. doi:10.1177/15248380211073839
7. Buscemi N, Hartling L, Vandermeer B, Tjosvold L, Klassen TP. Single data extraction generated more errors than double data extraction in systematic reviews. *J Clin Epidemiol*. 2006;59:697-703. doi:10.1016/j.jclinepi.2005.11.010
8. Pick JL, Nakagawa S, Noble DWA. Reproducible, flexible and high-throughput data extraction from primary literature: the metaDigitise r package. *Methods Ecol Evol*. 2019;10:426-431. doi:10.1111/2041-210X.13118
9. Miyakawa T. No raw data, no science: another possible source of the reproducibility crisis. *Mol Brain*. 2020;13:24. doi:10.1186/s13041-020-0552-2
10. Koricheva J, Gurevitch J. Uses and misuses of meta-analysis in plant ecology. *J Ecol*. 2014;102:828-844. doi:10.1111/1365-2745.12224
11. López-Nicolás R, López-López JA, Rubio-Aparicio M, Sánchez-Meca J. A meta-review of transparency and reproducibility-related reporting practices in published meta-analyses on clinical psychological interventions (2000–2020). *Behav Res Methods*. 2022;54:334-349. doi:10.3758/s13428-021-01644-z

12. Curtis PS, Mengersen K, Lajeunesse MJ, Rothstein HR, Stewart GB. Extraction and critical appraisal of data. In: Koricheva J, Gurevitch J, Mengersen K, eds. *Handbook of Meta-Analysis in Ecology and Evolution*. Princeton University Press; 2013:52-60.
13. Wan X, Wang W, Liu J, Tong T. Estimating the sample mean and standard deviation from the sample size, median, range and/or interquartile range. *BMC Med Res Methodol*. 2014;14:135. doi:[10.1186/1471-2288-14-135](https://doi.org/10.1186/1471-2288-14-135)
14. Lajeunesse MJ. Squeezing data from scientific images using the juicr package for R. R Package Version 0.1. 2021.
15. Lajeunesse MJ. Facilitating systematic reviews, data extraction and meta-analysis with the metagear package for R. *Methods Ecol Evol*. 2016;7:323-330.
16. Nakagawa S, Ivimey-Cook ER, Grainger MJ, et al. Method reporting with initials for transparency (merit) promotes more granularity and accountability for author

contributions. *Nat Commun*. 2023;14:1788. doi:[10.1038/s41467-023-37039-1](https://doi.org/10.1038/s41467-023-37039-1)

SUPPORTING INFORMATION

Additional supporting information can be found online in the Supporting Information section at the end of this article.

How to cite this article: Ivimey-Cook ER, Noble DWA, Nakagawa S, Lajeunesse MJ, Pick JL. Advice for improving the reproducibility of data extraction in meta-analysis. *Res Syn Meth*. 2023; 14(6):911-915. doi:[10.1002/jrsm.1663](https://doi.org/10.1002/jrsm.1663)