

How to do meta-analysis of open datasets

The amount of open data in ecology and evolution is increasing rapidly, yet this resource remains underused. Here, we introduce a new framework and case study for conducting meta-analyses of open datasets, and discuss its benefits and current limitations.

Antica Culina, Thomas W. Crowther, Jip J. C. Ramakers, Phillip Gienapp and Marcel E. Visser

In recent decades, the meta-analysis approach has emerged as the most valuable avenue for scientific progress, along with empirical studies and theoretical models^{1,2}. Traditional meta-analysis combines results from a number of studies (ideally all conducted on the same research question, to statistically summarize findings, evaluate discrepancies and detect generalizable effects². The ability to detect overarching patterns makes meta-analyses extremely relevant to evolutionary ecology, which is characterized by highly complex systems, heterogeneous environments and variable methodologies^{3,4}.

Systematic advances in the meta-analysis approach over the past decade have been intended to improve the transparency, replicability, reliability and impact of data synthesis efforts^{2,5–7}. However, despite these advances, the major outstanding limitation of any synthesis remains the challenge of accessing a comprehensive range of available data on the topic⁷. Conventionally, meta-analyses are conducted using effect sizes (that is, measure of the strength and direction of effects) extracted from the values reported in published studies. These meta-analyses are often limited to studies that focus specifically on the topic of interest (we term these ‘target studies’). However, a wealth of useful data is often available in various ‘non-target studies’ that have attained relevant information to address different research questions. Additional data from non-target studies can enhance the statistical power of meta-analyses (a fact that has been widely accepted and embraced in medical research⁸), as well as considerably reduce current issues with biased effect sizes. These data can be used either on their own, or in a combination with data from target studies. Until now, the complex and variable research landscape in ecology and evolution has restricted such data ingestion from non-target studies. However, the increase in data made openly accessible, as now required by many journals, is transforming our capacity to access, evaluate and use raw data from both target and non-target studies. Hence, our potential to survey the data landscape

and gain a comprehensive understanding of the available information has never been greater⁹. Yet, unlike other scientific fields, this resource remains relatively unexploited in the field of ecology and evolution^{10,11}.

Data retrieval for meta-analysis

Here, we describe how to transparently retrieve and select data, when the information retrieval starts from published (open) datasets, rather than from published studies. Our standard is based on existing guidelines for the information retrieval in ecological/evolutionary meta-analysis^{5,6,12,13}, but adapted specifically for open data. The retrieval and selection process should be highly transparent — we provide a checklist of the information that needs to be recorded (Table 1). This information should ideally be supported by the Preferred Reporting Items for Systematic Reviews and Meta-Analyses¹⁴ diagram (Supplementary Fig. 1).

In the first step of the approach (Step 1), researchers need to identify the type of

data needed to answer the meta-analysis question (or test hypothesis), set appropriate exclusion/inclusion criteria and choose the search terms (used in a search for the relevant data). This is followed by the data search. In evolutionary ecology, datasets are usually scattered across various repositories (for example, Dryad, Figshare, Zenodo) or published in the supplementary materials associated with a paper. Thus, an effective search should be conducted using data-harvesting platforms that crawl through many different research data repositories that host research data (like Web of Science crawls through journals in a search for articles); some also explore supplementary materials of published papers for additional information. A complete overview of how to navigate the data landscape by using data search platforms can be found in ref. ¹⁵. We suggest using DataCite, BASE search engine and DataONE (see Box 1). The original search terms usually need to be adjusted according to the output of the initial

Table 1 | Checklist of the main steps in conducting meta-analysis that starts from datasets

Step	What to record (report)
Step 1: what type of data are needed and where/how to obtain them?	Research question/questions The exact exclusion/inclusion criteria Platform(s) used in search Search terms and syntax (for every platform; whether and how search terms were adjusted)
Step 2: screening the results according to the meta-data provided (keywords, dataset title, description of the dataset and/or subject area)	What meta-data screening was based on Number of excluded results Reasons for exclusion (optional)
Step 3: open and screen remaining datasets	Number of excluded results Reasons for exclusion (optional)
Step 4: detailed examination of the datasets. Contacting the authors of the dataset about missing/unclear information	Number of excluded results Reasons for exclusion Whether the authors were contacted and with what outcome
Step 5: calculate the effect sizes	Statistical procedures to calculate effect sizes
Step 6: contact the authors to check if they agree with the approach	Contact letter, author responses, dates of contact Datasets excluded based on authors' feedback and reasons why
Step 7: conduct the statistical part of meta-analysis	The dataset used in meta-analysis Exact models/formulas

Following these steps will ensure a transparently conducted meta-analysis that complies with the current scientific standards.

Thus, it is important to record and report on which meta-data the screening was based. This step is equivalent to the initial screening of the title, abstract and keywords in the 'traditional' meta-analysis that starts from published studies. The main difference is that the standards to describe datasets are less well established than the standards to describe articles (title, abstract, keywords, subject areas). Thus, this screening might be more time-consuming, and lead to the retention of more irrelevant datasets. Next (Step 3), each potentially relevant dataset should be opened and screened to identify whether it corresponds to meta-analysis requirements.

The remaining datasets are relevant according to the dataset type, but some will be excluded (Step 4) as they do not match the specific inclusion criteria or are not fit for use because information crucial to run the desired analysis (to obtain the effect sizes) is missing (equivalent to under-reported effects in the approach that starts from published studies). At this stage, researchers might decide to contact the dataset owner(s).

The final list of datasets is then used to calculate the effect sizes (Step 5). Ideally, all effect sizes are calculated in the same, standardized way. This process can take several sub-steps. In line with good scientific practice, and to address the issues of data misinterpretation¹⁶, owners of the datasets should be contacted, at the latest, when analysing their data, and asked whether they agree with the way in which the data were processed (Step 6). Some data owners specifically ask (in the meta-data files) to be contacted directly if there are plans to use the data. Some datasets might be excluded after this step. Statistical analyses can then be conducted using these effect sizes (Step 7) following the existing guidelines (choose an appropriate model, explore the sources of heterogeneity, account for non-independencies and, if considered necessary, test for publication bias^{12,13}). Statistical analysis can also be conducted using both, effect sizes calculated from raw data and those calculated using values reported in published articles (when possible). In this case, information retrieval protocol should be recorded separately for the data and article selection process^{6,12,13}, and we would further advise controlling for the source of effect size (data or article) when conducting the statistical analysis.

To demonstrate the information retrieval framework, in Box 1 we outline the search for pedigree datasets for the meta-analysis that aimed at evaluating the strength of the evidence for the environmental coupling of heritability and selection¹⁷.

Benefits of open data meta-analysis

Our case study (Box 1) demonstrates an obvious benefit of the information retrieval that starts from published data (rather than published studies): the considerable increase in the data available to conduct meta-analysis (and thus in the number of research questions that can be addressed¹⁸). These data can be used on their own to calculate effect sizes for the meta-analysis, or used alongside effect sizes extracted from published studies. In our example, a traditional meta-analysis was impossible (only two published studies on the research question, see Box 1). Use of open data from studies that themselves addressed another question enabled us to collect enough evidence for meta-analysis. Given that the number of published datasets is greatly increasing across evolutionary and ecological fields^{9,15}, the scope of evolutionary ecology meta-analysis can be extended, and not limited only to target studies in the published literature.

An additional benefit of open data is the reduced publication bias that stems from the selective reporting of 'significant' or 'interesting' results⁷. The under-reporting of weak, negative or unwanted effects (or ambiguous results) is common across scientific disciplines: two reviews showed that basic information (sample size and variance) was missing from generally half of otherwise relevant primary studies collected for meta-analysis in conservation ecology¹⁹ and evolutionary ecology²⁰. Even more worrying is that these under-reported results seem to be a biased sample of all results²⁰. However, datasets, and effect sizes calculated using published datasets, are less likely to suffer from this kind of issue. Datasets that support published studies can be also used to verify or supplement the reported results of the study, increasing the number of effect sizes that can be calculated (missing or contradictory reported results).

Finally, meta-analyses conducted using the values reported in studies have to combine effect sizes calculated in a different way (as primary studies analyse their data and report the results differently). Effect sizes can be calculated in a consistent manner if the original data are used (such as in our case study; Box 1), thus leading to directly compatible effect sizes¹⁸.

Limitations of open data meta-analysis

Despite the apparent benefits, our meta-analysis conducted using non-target research data suffers several limitations. These should not discourage data-driven meta-analysis, but rather be acknowledged and, if possible, adequately resolved.

First, as our case study demonstrates, the description of datasets is often insufficient to enable a sensitive and targeted search. This means that data searches may retrieve a substantial number of irrelevant datasets, while also missing some relevant ones. However, this has always been a limitation of meta-analyses, and we believe this will only improve as the scientific community continues to embrace the advised data standards (for example, ref. ²¹), supported by improvements in the data curation by research institutions and scientific repositories. The second and related issue is the quality of the retrieved datasets, a number of which might need to be excluded due to the lack of sufficient information. In our case, this reduced the number of species for the analysis and led to loss of a number of taxonomic groups (panel b of the Box 1 figure). The third issue is the misinterpretation of data used in meta-analysis¹⁶, especially when using non-target studies that addressed different questions from the proposed study. Contacting data owners is probably the best approach to address this issue (for example, we excluded 4 out of 18 datasets based on owner comments) and should thus be standard in open-data meta-analysis. The outlined issues might make meta-analysis based on data more time-consuming compared with traditional meta-analysis, but based on our experience this will vary from case to case.

Conclusion

The meta-analysis approach has become increasingly important across ecological and evolutionary research fields, having a strong impact on future research, interventions and policies. Here, we introduce a new standard on how to conduct a data-driven meta-analysis that, in contrast to the conventional meta-analysis, uses research data rather than published studies. This new standard is now possible given that the amount of open research data has been steadily increasing across evolutionary and ecological fields. We show that new questions can be addressed with the use of this ever-growing data landscape, broadening the scope of meta-analysis in evolutionary ecology. In addition, by embracing open data, evolutionary ecology has the potential to benefit from a spectrum of higher standards and reporting practices brought in the new era of open science. □

Antica Culina^{1*}, Thomas W. Crowther^{1,2}, Jip J. C. Ramakers¹, Phillip Gienapp¹ and Marcel E. Visser¹

¹Netherlands Institute of Ecology (NIOO-KNAW), Wageningen, Netherlands. ²Institute of Integrative Biology, ETH Zurich, Zurich, Switzerland.

*e-mail: A.Culina@nioo.knaw.nl

Published online: 18 June 2018

<https://doi.org/10.1038/s41559-018-0579-2>

References

- Cadotte, M. W., Mehrkens, L. R. & Menge, D. N. L. *Evol. Ecol.* **26**, 1153–1167 (2012).
- Gurevitch, J., Koricheva, J., Nakagawa, S. & Stewart, G. *Nature* **555**, 175–182 (2018).
- Jennions, M. D., Kahn, A. T., Kelly, C. D. & Kokko, H. *Evol. Ecol.* **26**, 1119–1151 (2012).
- Stewart, G. B. & Schmid, C. H. *Res. Synth. Methods* **6**, 109–110 (2015).
- Lortie, C. J., Stewart, G., Rothstein, H. & Lau, J. *Res. Synth. Methods* **6**, 246–264 (2015).
- Bayliss, H. R. & Beyer, F. R. *Res. Synth. Methods* **6**, 136–148 (2015).
- Parker, T. H. et al. *Trends Ecol. Evol.* **31**, 711–719 (2016).
- Simmonds, M. C. et al. *Clin. Trials* **2**, 209–217 (2005).
- Roche, D. G., Kruuk, L. E. B., Lanfear, R. & Binning, S. A. *PLoS Biol.* **13**, e1002295 (2015).
- Wallis, J. C., Rolando, E. & Borgman, C. L. *PLoS ONE* **8**, e67332 (2013).
- Evans, S. R. *PLoS Biol.* **14**, 1–9 (2016).
- Koricheva, J., Gurevitch, J. & Mengersen, K. (eds) *Handbook of Meta-Analysis in Ecology and Evolution* (Princeton Univ. Press, Princeton, Oxford, 2013).
- Nakagawa, S., Noble, D. W. A., Senior, A. M. & Lagisz, M. *BMC Biol.* **15**, 18 (2017).
- Moher, D., Liberati, A., Tetzlaff, J., Altman, D. G. & The PRISMA Group *PLoS Med.* **151**, 264–269 (2009).
- Culina, A. et al. *Nat. Ecol. Evol.* **2**, 420–426 (2018).
- Mills, J. A. et al. *Trends Ecol. Evol.* **30**, 581–589 (2015).
- Ramakers, J. J. C., Culina, A., Visser, M. E. & Gienapp, P. *Nat. Ecol. Evol.* <https://doi.org/10.1038/s41559-018-0577-4> (2018).
- Mengersen, K., Gurevitch, J. & Schmid, M. D. in *Handbook of Meta-Analysis in Ecology and Evolution* (eds Koricheva, J., Gurevitch, J. & Mengersen, K.) 300–313 (Princeton Univ. Press, Princeton, Oxford, 2013).
- Cote, I. M. & Reynolds, J. D. *Evol. Ecol.* **26**, 1237–1252 (2012).
- Cassey, P., Ewen, J. G., Blackburn, T. M. & Møller, A. P. *Proc. R. Soc. Lond. B* **271**, 451–454 (2004).
- Wilkinson, M. D. et al. *Sci. Data* **3**, 160018 (2016).
- Wood, C. & Brodie, E. *Ecol. Lett.* **19**, 1189–1200 (2016).
- Husby, A., Visser, M. E. & Kruuk, L. E. B. *PLoS Biol.* **9**, e1000585 (2011).
- Wilson, A. J. et al. *PLoS Biol.* **4**, e216 (2006).

Author contributions

A.C. collected the data and wrote the majority of the manuscript; T.W.C., J.J.C.R., P.G. and M.E.V. all contributed to the discussion and the manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information is available for this paper at <https://doi.org/10.1038/s41559-018-0579-2>.