



HAL
open science

Understanding the provenance and quality of methods is essential for responsible reuse of FAIR data

Tracey L. Weissgerber, Malgorzata Anna Gazda, Gustav Nilsson, Gerben ter Riet, Kelly D Cobey, Julia Priess-Buchheit, Jorge Noro, Robert Schulz, Joeri K. Tjeldink, Evgeny Bobrov, et al.

► To cite this version:

Tracey L. Weissgerber, Malgorzata Anna Gazda, Gustav Nilsson, Gerben ter Riet, Kelly D Cobey, et al.. Understanding the provenance and quality of methods is essential for responsible reuse of FAIR data. *Nature Medicine*, 2024, 30 (5), pp.1220-1221. 10.1038/s41591-024-02879-x . hal-04531169

HAL Id: hal-04531169

<https://hal.science/hal-04531169v1>

Submitted on 27 Jun 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

Understanding the provenance and quality of methods is essential for responsible reuse of FAIR data

Tracey L. Weissgerber^{1*}, Małgorzata Anna Gazda², Gustav Nilsson^{1,3,4}, Gerben ter Riet^{5,6}, Kelly D. Cobey⁷, Julia Prieß-Buchheit⁸, Jorge Noro⁹, Robert Schulz¹, Joeri K. Tjeldink¹⁰, Evgeny Bobrov¹, Alexandra Bannach-Brown¹, Delwen Franzen¹, Ugo Moschini¹¹, Florian Naudet¹², Ulrich Mansmann¹³, Maia Salholz-Hillel¹, Anita Bandrowski^{13,14}, Malcolm R Macleod¹⁵

¹QUEST Center for Responsible Research, Berlin Institute of Health at Charité – Universitätsmedizin Berlin, Berlin, Germany

²Department of Biological Sciences, University of Montréal, Montréal, Canada

³Department of Clinical Neuroscience, Karolinska Institutet, Stockholm, Sweden

⁴Swedish National Data Service, University of Gothenburg, Gothenburg, Sweden

⁵Center of Expertise Urban Vitality, Faculty of Health, Amsterdam University of Applied Sciences, Amsterdam, The Netherlands

⁶Dept Cardiology, Amsterdam University Medical Center, Amsterdam, The Netherlands

⁷Meta-Research and Open Science Program, University of Ottawa Heart Institute, Ottawa, Canada

⁸Institut für Pädagogik, Kiel University, Kiel, Germany

⁹Institute for Interdisciplinary Research, Center for Business and Economics Research (CeBER), University of Coimbra, Coimbra, Portugal

¹⁰AmsterdamUMC, location VUmc, Department of Ethics, Law and Humanities, Amsterdam, the Netherlands

¹¹Data Analysis Office, Istituto Italiano di Tecnologia, Genoa, Italy

¹²University of Rennes, CHU Rennes, Inserm, Irset (Institut de recherche en santé, environnement et travail)-UMR_S 1085, CIC 1414 (Centre of Clinical Investigation of Rennes), Rennes, France. Institut Universitaire de France (IUF), Paris, France

¹³Department of Medical Information Sciences, Biometry, and Epidemiology, Medical Faculty, Ludwig-Maximilians-Universität München, München, Germany

¹⁴Dept of Neuroscience, University of California – San Diego, San Diego, United States

¹⁴BIH Visiting Professor (funded by Stiftung Charité), Berlin Institute of Health at Charité – Universitätsmedizin Berlin, Berlin, Germany

¹⁵Centre for Clinical Brain Sciences, University of Edinburgh, Edinburgh, Scotland, United Kingdom

*Corresponding author, tracey.weissgerber@bih-charite.de

Data availability and reusability are critical to open research. The FAIR principles provide a minimal set of guiding principles for making data findable, accessible, interoperable and reusable.¹ Open data is not necessarily FAIR, and FAIR data are not necessarily open. Since their publication in 2016¹, the FAIR principles have accelerated the open data movement by inspiring activities and infrastructure development (as shown in refs²⁻⁴, <https://eosc-portal.eu/about/eosc>). The principles are also being adapted for other research outputs, such as software⁵. As funders increasingly demand FAIR practices and researchers work to implement FAIR, additional actions should be taken for responsible data use and reuse.

The FAIR principles indirectly outline the responsibilities of the data depositor by identifying dataset properties that facilitate reuse. However, the data provenance and the quality of the methods and procedures used to generate and validate data are often overlooked. This information is essential for responsible data reuse. FAIR data evaluations typically focus on the question “Can I reuse these data?”. We argue that it’s time to also ask “Should I re-use these data?”; and “How should I re-use these data responsibly?”. These questions allocate responsibilities between the data depositor and the prospective data user. This shift should include several elements.

While FAIR data is necessary for reusability, this does not guarantee scientific rigor, trustworthiness, or research quality. In addition to determining whether data are FAIR, prospective data users should consider whether the data are appropriate to answer their research question. Furthermore, data users must consider the rigor and quality of the study design and procedures used to generate the data, and whether reuse is likely to yield trustworthy results. Sharing FAIR data may encourage others to uncritically reuse data. Reuse of data from poorly designed experiments may yield valuable insights if users address design limitations when analyzing and interpreting the data^{6,7}. However, uncritical reuse of problematic data to generate new, untrustworthy findings may be harmful. We believe that comprehensive FAIR data sharing evaluations (such as EOSC call HORIZON-INFRA-2022-EOSC-01, <https://ec.europa.eu/info/funding-tenders/opportunities/portal/screen/opportunities/topic-details/horizon-infra-2022-eosc-01-01>) should consider the quality of the data provenance.

We propose that creating guiding principles that outline the responsibilities of the data user will facilitate and enhance responsible data reuse. These responsibilities might include performing a systematic search to identify datasets relevant to the new research question, assessing and describing the scientific rigor of the methods and procedures used to generate these data, and determining whether identified datasets are appropriate to answer the research question. If the underlying study has a high risk of bias, users should develop an analysis plan to address this bias or avoid using the dataset. Researchers who aim to combine datasets should determine whether any datasets have properties that preclude combination. Pre-registration of secondary analyses is important, as data reuse allows many exploratory hypotheses to be tested, at low cost, and only those supporting a particular view or reaching a particular evidence threshold to be published. Additionally, users might share products resulting from studies that reused data, including protocols, modified data, code, software or tools. Further discussion is needed to clearly define guiding principles.

Responsible data reuse requires knowing how the data were generated. However, progress on open and reusable methods and procedures lags far behind open data. Detailed methods may facilitate a broader spectrum of data reuse, including for purposes not anticipated by the data depositors. FAIR highlights the importance of metadata in helping potential data users to understand the dataset,¹ and several groups have set domain-specific metadata standards (<https://fairsharing.org/>). The term 'metadata', however, is poorly understood by many researchers. Furthermore, the importance of metadata is often disregarded, as sharing of high-quality metadata and data are typically not incentivized or rewarded. Data sharing requirements are often unfunded mandates, introduced without adequate training. Many depositors provide little metadata, or simply cite a paper describing the study. This is often inadequate, as publications regularly lack essential methodological details.⁸

The research community can take several steps to solve these problems. When using the technical term, metadata, researchers should state that metadata include detailed methods. Data management plans should include sharing of high-quality information on methods. Methods contextualizing datasets should include detailed information about the study aim, design, methods used, any additional measures taken to reduce the risk of bias, and study limitations. Data depositors should also share guidance for responsible dataset reuse. Research assessment systems must reward and incentivize sharing of methods, data and code as separate research outputs. The academic assessment system primarily values papers, and so researchers who share methods, data and code are doing more work without additional reward. This must change.

Data repositories can contribute by providing fields for data depositors to link detailed methods shared in methods repositories. This could include pre-registrations, study design protocols, reusable step-by-step protocols and data validation or analysis plans. Many researchers use generalist repositories, which allow unstructured depositing of data, methods and other materials. Further research is needed to determine whether the structured fields in methods repositories improve reporting. Generalist repositories should have machine-readable systems for determining what materials (such as methods, data, and code) an entry contains,

Methods are crucial to scientific advancement; they are not simply a tool to contextualize datasets. If properly shared, methods could be more widely reused than data. Open and reusable methods should be shared as separate, essential research products. A vibrant open methods community is needed to champion this, as exists for open data and open code.

FAIR data should increase the opportunities for secondary analysis. Previously, these analyses have been conducted by, or in close collaboration with, the researchers who collected the data, or involved large, well-documented publicly available datasets, such as population studies or government registries. FAIR data sharing can further expand the number and types of available datasets, while reducing the need for collaboration between the data depositor and the data user or re-user, but this will require changes to data depositing and reuse strategies.

We encourage those with relevant expertise who are interested in contributing to principles for responsible data reuse to contact us.

Competing interests

TLW leads PRO-MaP (Promoting Reusable and Open Methods and Protocols), which aims to improve the quality and reporting of methods and protocols in the life sciences. AB is the founder and CEO of SciCrunch Inc a company devoted to improving the scientific literature. UM (Ulrich Mansmann) is coordinating the doctoral network MSCA-DN SHARE-CTD (HORIZON-MSCA-2022-DN-01 101120360), funded by the EU. UM (Ulrich Mansmann) is member of the national working group Data Sharing of and supported by the German Medical Informatics Initiative (BMBF 01ZZ23048, 01ZZ2316). FN received funding from the French National Research Agency (ANR-17-CE36-0010), the French ministry of health and the French ministry of research. He is a work package leader in the OSIRIS project (Open Science to Increase Reproducibility in Science). The OSIRIS project has received funding from the European Union's Horizon Europe research and innovation programme under the grant agreement No. 101094725. He is a work package leader for the doctoral network MSCA-DN SHARE-CTD (HORIZON-MSCA-2022-DN-01 101120360), funded by the EU. KDC is the co-chair of DORA (Declaration of Research Assessment).

References

1. Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and stewardship. *Sci. Data* **3**, 160018 (2016).
2. Gaiarin, S. P. *FAIR Assessment and Certification in the EOSC region*.
<https://zenodo.org/record/4486280> (2020) doi:10.5281/ZENODO.4486280.
3. Devaraju, A. & Huber, R. An automated solution for measuring the progress toward FAIR research data. *Patterns* **2**, 100370 (2021).
4. Patel, B. & Soundarajan, S. Making biomedical research software findable, accessible, interoperable, reusable (FAIR) with FAIRshare. (2022)
doi:10.7490/F1000RESEARCH.1119055.1.
5. Chue Hong, N. P. *et al.* FAIR Principles for Research Software (FAIR4RS Principles). (2021) doi:10.15497/RDA00065.
6. Nielson, J. L. *et al.* Topological data analysis for discovery in preclinical spinal cord injury and traumatic brain injury. *Nat. Commun.* **6**, 8581 (2015).
7. Torres-Espín, A. *et al.* Topological network analysis of patient similarity for precision management of acute blood pressure in spinal cord injury. *eLife* **10**, e68015 (2021).
8. Errington, T. M., Denis, A., Perfito, N., Iorns, E. & Nosek, B. A. Challenges for assessing replicability in preclinical cancer biology. *eLife* **10**, e67995 (2021).