Generalist Repository Ecosystem Initiative

# GREI Data citation best practices for repositories

## Introduction

One of the objectives of the Generalist Repository Ecosystem Initiative (GREI) is to implement data metrics that enable reporting on the reach and impact of NIH-funded research data.

Data citations are a key component of the measures of data usage, as they bring benefits to the data creators, the data users, and the scholarly communication ecosystem more broadly:

- Data citations are a signal of a dataset being used in research (beyond mere exploration), providing valuable information to evaluate data usage.
- Data citations provide credit for the data producer, the citation recognizes the individual(s) or organization(s) that collected and shared the data used in the citing work.
- For academic researchers, accruing citations to datasets can also be valuable as part of research evaluation frameworks (e.g. for hiring or promotion), as they provide evidence of the reach of their open datasets.
- Surveys of researchers regularly show that getting citations to their research papers as well as citations to the datasets themselves are among the biggest motivators for them to publish their data (see for example The State of Open Data report).
- Data citations increase the rigor and reproducibility of research, enabling data users to document the source of the data they employed as part of research activities.
- Data citations can enable the development of tools to aid search and discovery of research works; increasing visibility on what outputs cite a dataset can help researchers find other scholarly objects relevant to their work.

Providing visibility for data citations is therefore a way of increasing the information on data usage available to the community, it signals the added value of data repositories, and can create an incentive for researchers to share more of their data and to cite open data.

The importance of data citations has been recognized by earlier community work, for example, in the form of the [FORCE11 Data Citation Principles,](#) which cover the purpose, function and attributes of citations to data. The principles note as part of its [guidance](#) that while citations may vary in style, citations to data "*should be included in the full reference list along with citations to other types of works.*"

All of the GREI repositories (Dataverse, Dryad, Figshare, Open Science Framework, Mendeley Data, Vivli, and Zenodo) already collect data citations or have it on their roadmap to add this feature. This group of repositories has also been considering approaches to capture those citations in a manner that provides transparent and consistent information to the community on citations for data.

We outline here, from our perspective as generalist repositories, our recommendations for handling data citations in repositories. We hope that this will be useful for other repositories that are seeking to implement or update their practices for data citations.

## What do we mean by data citation?

Citations can be created for any scholarly output, but for the scope of this document, we are focusing specifically on citations to datasets. For the purpose of our recommendations, we refer to a data citation as:

> ***A link between a dataset in a repository and another scholarly output where both have persistent identifiers.***

Examples of data citations include those below (please note this is not an exhaustive list):

- A journal article that cites the dataset underlying the results
- A preprint that cites a dataset re-analysed using a different methodology
- A software tool that cites the training dataset used during development
- A conference presentation or poster that cites the dataset produced during the study
- A clinical guideline cites a dataset that informed the recommendations
- A derived dataset (e.g. harmonized data, an analysis-ready dataset) that cites the original dataset on which it builds

## Collecting data citations

Repositories can collect citations to datasets through self report by authors as they deposit and/or update data records, or harvest data citations from external sources. Data sources employed by GREI repositories to harvest data citations are listed below:

- DataCite / Crossref

- Dimensions
- Europe PMC
- NASA ADS (Astrophysics Data System)

We recommend that repositories submit data citations they collect to DataCite (see DataCite's documentation for contributing data citations to DataCite) so that these connections can be aggregated and discoverable by others in the community.

## Storing data citations

Repositories collect data citations via the metadata for the datasets they host. We recommend that repositories use the metadata fields recommended by GREI to collect data citations, specifically, those listed below to establish the relationship between the dataset and the citing object:

12. relatedIdentifier
12.a relatedIdentifierType
12.b relationType
12.f resourceTypeGeneral

The above fields align to the properties in the DataCite metadata schema for establishing citations. Note that under DataCite's schema, the following relationType entities as part of the metadata for the dataset are recorded as a data citation:

IsCitedBy
IsReferencedBy
IsSupplementTo

### 'Cite As' template for data citation

In order to encourage researchers and other parties to cite datasets they use, data repositories should provide a citation template on the landing page of the dataset. We include below a citation template example from the Dryad repository:

The NLM Style Guide recommends that citations of datasets include the following elements: Author, Title, Type of Medium, Publisher, Date of Publication Date, Date of Update/Revision, Date of Citation[1], Availability (i.e. url or DOI). An example citation is included below[2]:

*Di Stefano B, Collombet S, Graf T. Time-resolved gene expression profiling during reprogramming of C/EBPα-pulsed B cells into iPS cells [dataset]. 2014 May 22 [cited 2014 Oct 6]. In: figshare [Internet]. Available from: http://dx.doi.org/10.6084/m9.figshare.939408*
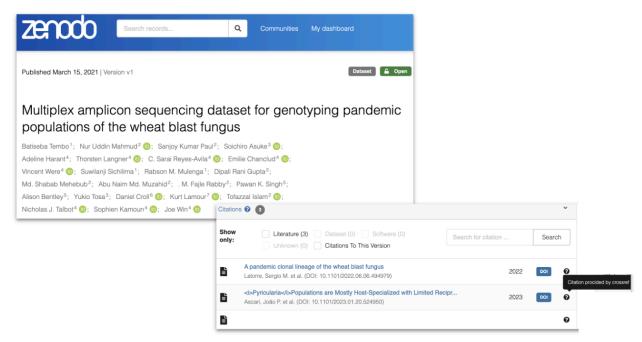
## Exposing data citations

Repositories should expose the citations for individual datasets on the landing page for the dataset record. To increase transparency and trust in the information displayed for data citations, we recommend that repositories indicate the provenance for the data citation, i.e. the source for the asserted citation. The citation information listed at the repository should include:
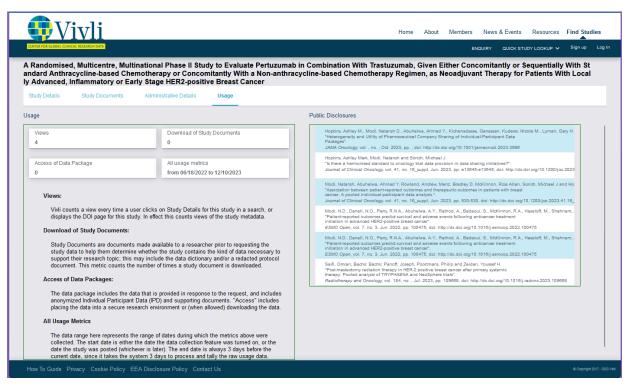
- The number of citations to the dataset
- The list of citing objects and their associated identifiers, or a link to where such a list of citing objects can be accessed
- A breakdown of the citations according to their provenance, in text format - where the repository harvests data citations from external sources (see examples in section Collecting data citations), those should be listed

---

1 Note the date of citation will only be created at the time at which the citation is introduced into the reference list, and thus, it will not be part of the template provided by the repository.

2 Citation adapted from https://www.ncbi.nlm.nih.gov/books/NBK7273/#A57904.

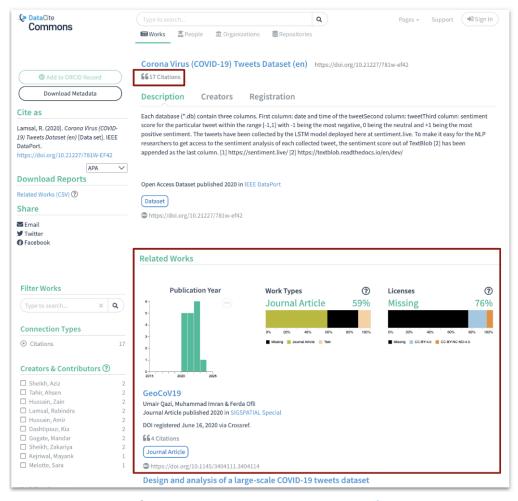Data citations for dataset on Zenodo: https://doi.org/10.5281/zenodo.4605959.



Data citations for dataset on Vivli, displayed in the 'Public Disclosures' box for the dataset at https://doi.org/10.25934/00005758. The term "Public Disclosures" was chosen at a time when the term "Citations" had a narrower connotation than the definition above.

The repository should also document, as part of the information pages for contributors and users, information on the mechanisms by which the repository collects and stores data citation, including the sources it employs.

## Data citations on DataCite Commons

Data citations captured via the metadata deposit with DataCite, as well as some citations available from Crossref metadata, are also displayed on the DataCite portal DataCite Commons. This web interface enables searches for scholarly resources with persistent identifiers as well as connections to metadata provided by DataCite, Crossref, ORCID, and others. DataCite Commons provides the number of citations for the dataset as well as a list of citing objects.



DataCite Commons record for https://doi.org/10.21227/781w-ef42, displaying the number of citations and citing objects.

## Data Citation Corpus

There have been challenges in consistently collecting information about data citations across repositories and the literature, due to the fact that workflows (e.g. by publishers) have not been optimized to collect citations to datasets, and that different parties are using their own mechanisms to capture citations and store that information in different locations. This has limited the ability of stakeholders such as institutions and funders to incorporate data as part of their evaluation frameworks.

To address these challenges and scale the data citations that are consistently made available to the community, DataCite is working on the development of the [Data Citation Corpus](#), which will provide a centralized resource that compiles data citations from a variety of sources, and make data citation information readily and openly available to the community.

The data sources for the data citation corpus include:

- Persistent identifier authorities: Sources that collect citations as part of their DOI registration workflow, such as DataCite and Crossref.
- Third-party aggregators: Sources that aggregate or discover citations through various techniques, such as full-text mining and curation. For example, the Chan Zuckerberg Initiative (CZI) has contributed data citations identified through the CZI Knowledge Graph, which mines the text of publications via a machine-learning algorithm.

The Data Citations Corpus will include data citations in DataCite; this incorporates citations deposited by DataCite-member repositories, including the GREI repositories. We recommend that all data repositories contribute their data citations to DataCite so that those citations can be integrated into the Data Citations Corpus.

DataCite will engage closely with GREI repositories during the development process of the corpus and invites other repositories and interested parties to get in touch if they are interested in contributing data citations to the corpus or in providing feedback as potential users of this resource.

## Appendix 1: Additional reading/Resources

- Curtin, L., Feri, L., Gautier, J., Gonzales, S., Gueguen, G., Scherer, D., Scherle, R., Stathis, K., Van Gulick, A., & Wood, J. (2023). GREI Metadata and Search Subcommittee Recommendations_V01_2023-06-29. Zenodo. https://doi.org/10.5281/zenodo.8101957
- Data Citation Synthesis Group: Joint Declaration of Data Citation Principles. Martone M. (ed.) San Diego CA: FORCE11; 2014 https://doi.org/10.25490/a97f-egyk
- Fenner, M., Crosas, M., Grethe, J.S. *et al.* A data citation roadmap for scholarly data repositories. *Sci Data* 6, 28 (2019). https://doi.org/10.1038/s41597-019-0031-8
- Data Citation Corpus: https://makedatacount.org/data-citation/
- JATS4R Data Citations recommendations, Version 2.0: https://doi.org/10.3789/niso-rp-36-2020
- Science, Digital; Hahnel, Mark; Smith, Graham; Schoenenberger, Henning; Scaplehorn, Niki; Day, Laura (2023). The State of Open Data 2023. Digital Science. Report. https://doi.org/10.6084/m9.figshare.24428194.v1

## Appendix 2: Description of citation metadata for human readers and contributors

To support understanding of how repositories collect and store citations via metadata, and how this is aggregated in DataCite, we include below descriptions for specific instances in metadata that designate a data citation or where a dataset cites another scholarly object.

| Metadata for citation<br>*The examples below reflect the properties within the dataset DOI metadata that designate a citation* | Description |
| --- | --- |
| Dataset DOI: 10.5066/F70Z71H2<br><br><resourceType resourceTypeGeneral="Dataset">Dataset</resourceType><br><br><relatedIdentifier relatedIdentifierType="DOI" relationType="IsCitedBy">https://doi.org/10.1021/acs.est.7b00812</relatedIdentifier> | The dataset is cited by the article with DOI 10.1021/acs.est.7b00812 |
| Dataset DOI: 10.5281/zenodo.5112001<br><br><resourceType resourceTypeGeneral="Dataset">Dataset</resourceType><br><br><relatedIdentifier relatedIdentifierType="DOI" relationType="IsCitedBy">https://doi.org/10.1101/2021.01.04.425314</relatedIdentifier> | The dataset is cited by the preprint with DOI 10.1101/2021.01.04.425314 |
| Dataset DOI: 10.5281/zenodo.8091721 | The dataset cites the preprint with |

| | |
|---|---|
| `<resourceType resourceTypeGeneral="Dataset">Dataset</resourceType>`<br><br>`<relatedIdentifier relatedIdentifierType="DOI" relationType="Cites">https://doi.org/10.1101/2022.08.09.503355</relatedIdentifier>` | DOI 10.1101/2022.08.09.503355 |
| Dataset DOI: 10.5281/zenodo.10149985<br><br>`<resourceType resourceTypeGeneral="Dataset">Dataset</resourceType>`<br><br>`<relatedIdentifier relatedIdentifierType="DOI" relationType="IsReferencedBy">https://doi.org/10.1101/2023.10.03.560679</relatedIdentifier>` | The dataset is cited by the preprint with DOI 10.1101/2023.10.03.560679 |
| Dataset DOI: 10.5281/zenodo.8399018<br><br>`<resourceType resourceTypeGeneral="Dataset">Dataset</resourceType>`<br><br>`<relatedIdentifier relatedIdentifierType="PMID" relationType="IsSupplementedBy">10.5194/egusphere-2023-2295</relatedIdentifier>` | The dataset is supplemented by the article with DOI 10.5194/egusphere-2023-2295 |
| Dataset DOI 10.6084/m9.figshare.24565030<br><br>`<resourceType resourceTypeGeneral="Dataset">Dataset</resourceType>`<br><br>`<relatedIdentifier relatedIdentifierType="DOI" relationType="IsSupplementTo">10.1159/000534801</relatedIdentifier>` | The dataset supplements the article with DOI 10.1159/000534801 |