*Tales from the AGU Data Help Desk:*
# Data/Software Availability Statements, Citations, and More

## 21 April 2022

**Chris Erdmann**
**Assistant Director, Data Leadership**
**American Geophysical Union**
**0000-0003-2554-180X**
**@libcce | cerdmann@agu.org**

# AGU Data & Software Sharing Guidance

What is covered:

- What data needs to be available?
- Repository Selection
- Availability Statement
- Data & Software Citation
- Citation Formatter
- Models & Simulations
- Journal Specific Guidance
- International Geo Sample Numbers
- **Data Help Desk**



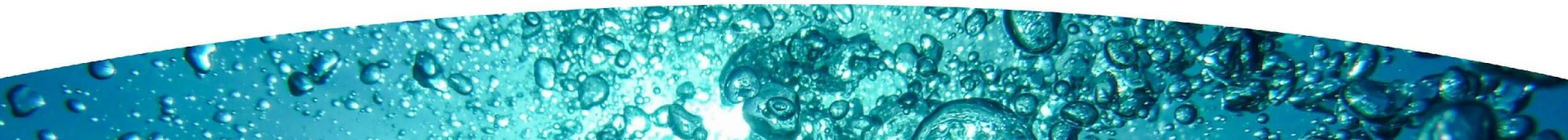https://www.agu.org/Publish-with-AGU/Publish/Author-Resources/Data-and-Software-for-Authors

**Main Goal:** Avoid parachuting researchers into your data/software and instead guide them as best as you can





Learning to get from A to B, one windy afternoon in the New Forest by Annie Spratt
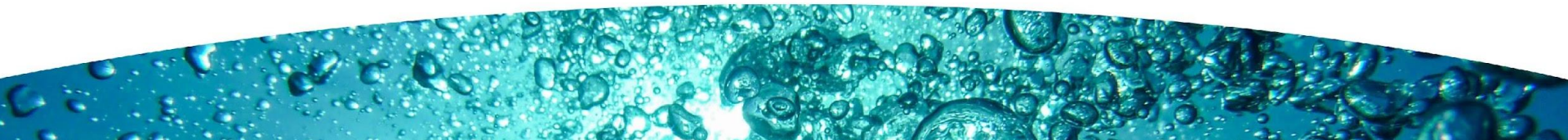Paragliding above the Chartreuse massif by Nicolas Tissot

# What is included in an Availability Statement?

1. A brief description of the type(s) of data or software
2. Repository Name(s) where they are deposited
3. Version (of software)
4. **DOI, Persistent Identifier** Link to Data or Software
5. Link to publicly accessible development platform (in the case of Software, e.g. GitHub)
6. Access Conditions (e.g. if Registration is Required)
7. Licensing/Permissions (e.g. Creative Commons Attribution)
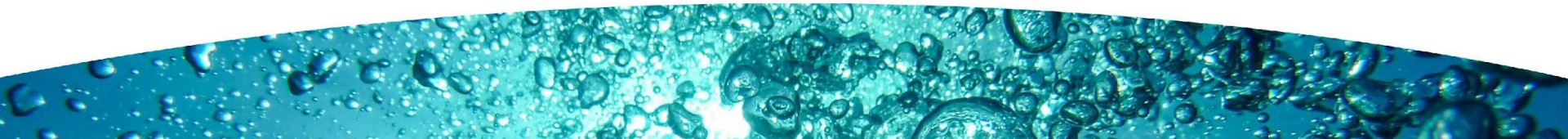8. **In-text citation in References (optional)**

# Availability Statement Templates

- The [type of data] data used for [brief context, description] in the study are available at [repository, source name] via [DOI, persistent identifier link] with [license, access conditions] [optional in-text citation in References]

- [Version number] of the [software name] used for [brief context, description of what the software was used for] is preserved at [DOI, persistent identifier link], available via [license type, access conditions] and developed openly at [software development platform link].* [optional in-text citation in References]

# What is included in a data/software citation?

1. Author(s) or project name(s)
2. Date / Software published
3. Title / Software name
4. Data or software release/version (optional)
5. **Bracketed description type (e.g., [Dataset], [Software], [Collection], [ComputationalNotebook])**
6. Repository name / Publication venue
7. **DOI, persistent identifier, URL**
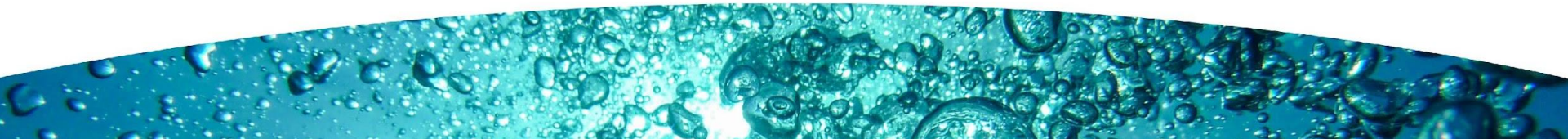8. Retrieved date (required when using URL)

# Data Citation Examples

- Fiechter, J., & Cheresh, J. (2020). Physical and biogeochemical drivers of alongshore pH and oxygen variability in the California Current System (Version 7) [Dataset]. Dryad. https://doi.org/10.7291/D1D96Q

- Edmunds, P. J., Didden, C., & Frank, K. (2021). Mean percentage cover of corals and Porites astreoides at each site by year at St. John, VI from 1992 to 2019 (Version 1) [Dataset]. Biological and Chemical Oceanography Data Management Office (BCO-DMO). https://doi.org/10.26008/1912/BCO-DMO.843284.1

- Alwarda, R., & Smith, I. (2021). Elevation data for Reflectors within the CO2 Deposit in Planum Australe, Mars [Dataset]. Zenodo. https://doi.org/10.5281/ZENODO.4639669

- Gries, C., Downs, R. R., O'Brien, M., Parr, C., Duerr, R., Koskela, R., et al. (2019). Return on Investment Metrics for Data Repositories in Earth and Environmental Sciences [Dataset]. Environmental Data Initiative. https://doi.org/10.6073/PASTA/D49BEC63F51603512EFA7E0FD2717203

# Software Citation Examples

- Lab for Exosphere and Near Space Environment Studies. (2019, March 20). lenses-lab/LYAO_RT-2018JA026426: Original Release (Version 1.0.0) [Software]. Zenodo. http://doi.org/10.5281/zenodo.2598836

- Bell, S. W. (2020). samwbell/saturn_counts: April 26, 2020 Release (Version 1.1.0) [Software]. Zenodo. https://doi.org/10.5281/ZENODO.3766959

- Shaoqian Hu. (2019, December 25). Direct surface wave radial anisotropy tomography package (Version 1.0) [Software]. Zenodo. http://doi.org/10.5281/zenodo.3592528

# DOI Citation Formatter

**Paste your DOI:**

https://doi.org/10.5061/dryad.v18jj97

For example 10.1145/2783446.2783605

**Select Formatting Style:**

apa

Begin typing (e.g. Chicago or IEEE.) or use the drop down menu.

**Select Language and Country:**

en-US

Begin typing (e.g. en-GB for English, Great Britain) or use the drop down menu.

Format

Mistry, R., & Ackerman, J. D. (2019). Data from: Flow, flux and feeding in freshwater mussels (Version 1) [Data set]. Dryad. https://doi.org/10.5061/DRYAD.V18JJ97

Copy to clipboard

https://citation.crosscite.org

# Availability Statement/Citation Paper Example

https://doi.org/10.1029/2021EA001675

# What repository?

https://data.agu.org/resources/useful-domain-repositories

# Help Desk Challenges

- Government Sites, Similar - Technical, Permissions
- Firewalls, Authentication - Openness, Availability, Anonymity
- Supplemental Information - Tradition, Peer Review
- FTP, Directories, Storage - Institutional, Compliant Solution
- Curation, Deposit Workflows  - Service, Publication Workflows
- Web Sharing, Dev Platforms - Citation Information
- Databases / Dynamic Services - Direct Access, Linking
- Available Upon Request - Culture
- Citation Nothingness (Paper not the Data) - Culture
- Website Home (Parachuting) - Laziness
- English Language - Language Diversity, Translation
- Many Data Links/Citations - Tables, Supplements (See Data Citation Community of Practice)
- ...

# Preserving Large Data!

## Preserving very large data is a challenge. Spoilers, there are no easy answers!

OCTOBER 01, 2021

When it comes to large datasets, we are often asked by authors and editors how they should preserve the data. These questions come via datahelp@agu.org and our data and software guidance discussions. Spoilers, there are no easy answers, yet! Here we offer our experience, share the current limitations, and the approaches we recommend with what is possible right now.

AGU requires that primary and processed data used for your research should be preserved and made available. This can range from observational data to the data used to generate your figures. The raw data may be needed, but usually, the processed or refined data that support and lead to the described results and allow other readers to assess your conclusions and build off your work should be preserved.

**For data that is large, over 1 Terabyte (TB)**, authors run into the challenge of finding a suitable repository. Many repositories have file size limitations but also costs associated with deposits over certain limits. This generalist repository comparison chart provides an overview of the limitations. Discipline-specific and institutional repositories are often a place to turn to for assistance with preserving large data but they also have limitations and potential costs. This emphasizes the importance of avoiding surprises at the time of publication by:

# QC/Rules Before Manuscript Submission



Example: https://curvenote.com/

# Notebooks Now!



Example: https://jupyter.org/

# Cookiecutter Data Science

# Indexing/Filtering

https://twitter.com/iainh_z/status/1509091657131638792

# Takeaways

- We need everyone's help advancing data/software sharing policies, requirements, guidance (e.g., societies, publishers)
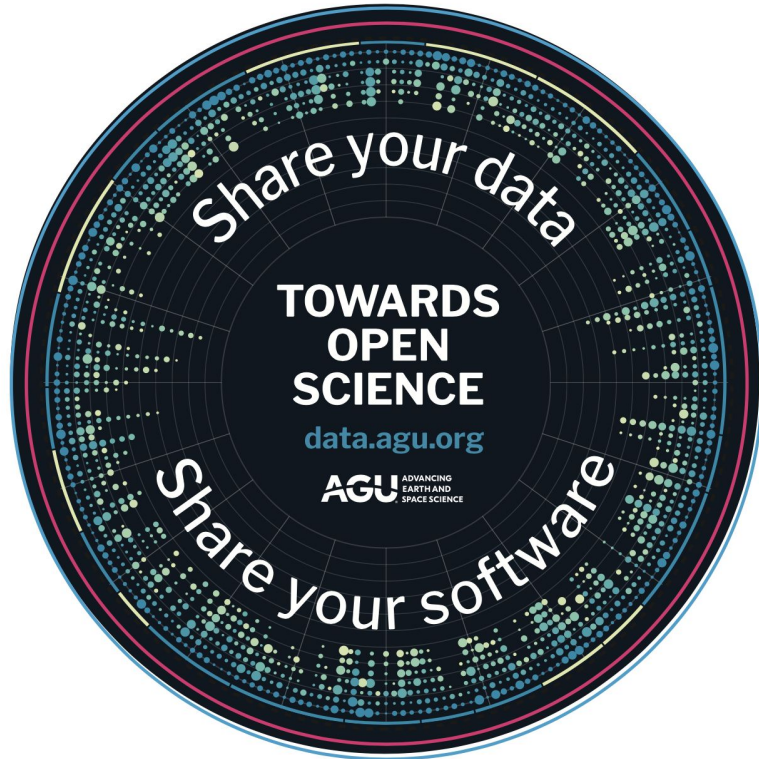- Institutions/disciplinary services need to work together to help simplify workflows for authors (e.g., help desk, repositories)
- We, the community, need to find better ways to integrate data/software best practices earlier in the research process, embed in research workflows (e.g, platforms, notebooks)
- Researchers are inundated with guidance, we need to streamline information as much as possible in combination with the point above (e.g, checklists)
- We need to demonstrate to researchers the value of sharing data/software by leveraging metadata (e.g., filtering, indexing)

# Thank you.



**Chris Erdmann**

Assistant Director, Data Leadership

American Geophysical Union

0000-0003-2554-180X

@libcce | cerdmann@agu.org