# NM EPSCoR Data Management and Project Data Repositories

Karl Benedict          Jonathan Wheeler

February 4, 2020

## A High-Level Integration of the Research Lifecycle and a Project Data Preservation and Publication Workflow

The EPSCoR SMART Grid Center includes diverse research activities that are being performed on a variety of systems ranging from desktop workstations and laptop computers to the high-performance computer systems at NMSU and UNM. This work is progressing through a standard research lifecycle that includes the production, management, and analysis of data; the sharing of data with project collaborators to enable shared research activities; and the publication of results. It is in the context of the publication of results that the project's data preservation and publication platforms come into play - providing the capability to preserve the data products generated by the project, and enabling FAIR (Findable, Accessible, Interoperable, and Reusable) access to project data products, including the generation of persistent identifiers (DOIs) for the published data. Meeting this objective meets the requirements that are imposed by both our project sponsor [NSF] and a growing number of publishers who require public access to data associated with published papers. Meeting this need requires the generation of reusable data products, the collection and compilation of documentation (metadata) associated with those data, and the placement of those data in an appropriate repository where they can be discovered and accessed.

In support of the EPSCoR SMART Grid Center project we have a number of current and planned capabilities to enable the efficient management, preservation, and sharing of FAIR data. The following figure highlights the workflow that we currently support and future workflow components that will be developed over the course of the project.

## Publishing Data in Dryad

The Dryad data repository is maintained by the California Digital Library. The University of New Mexico is a member institution through the NM ESPCoR SMART Grid Center. There are two separate paths to publication for SMART Grid researchers:

1. **Everyone** may use mediated upload services, supported by the UNM Libraries' Research Data Services. This service is available to both UNM and non-UNM affiliated researchers. Note that currently at least one author must be from UNM, and that person will have to initiate the dataset upload process in Dryad.
2. **UNM affiliated researchers** can directly access Dryad by logging with their ORCiD and UNM credentials.

Regardless of institutional affiliation, the first author listed on a dataset **must have an ORCiD**. This is a Dryad requirement, not an EPSCoR requirement.
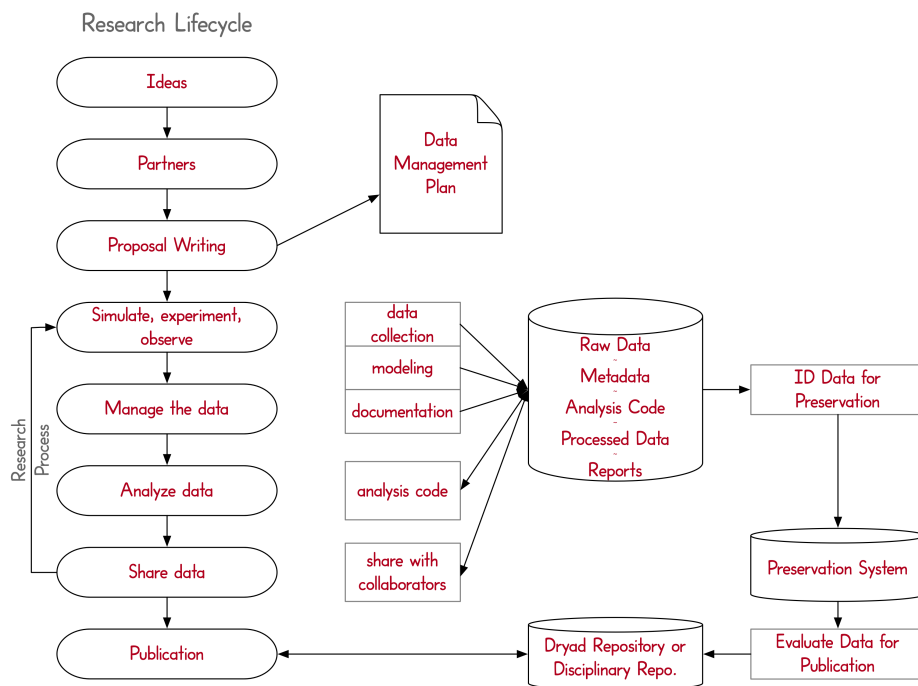
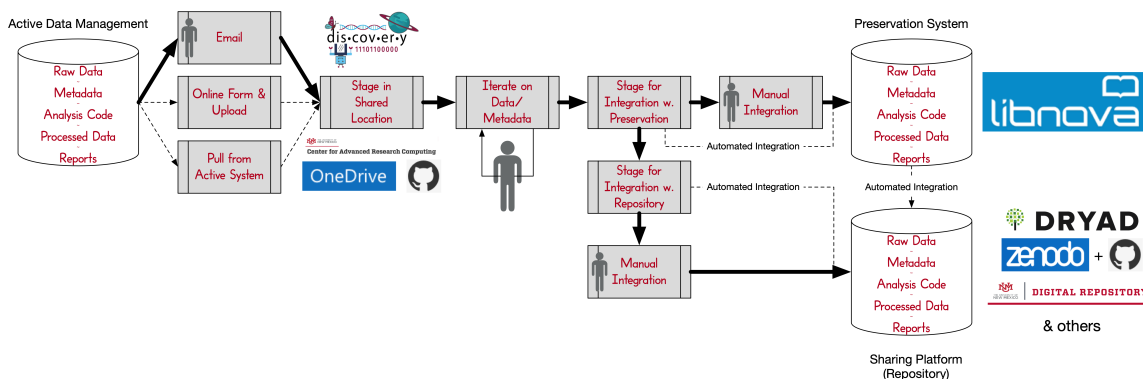Figure 1: Research lifecycle with high-level data workflow components



Figure 2: EPSCoR SMART Grid Center data preservation and publication workflow

**Getting an ORCiD**

ORCiD is a unique, persistent identifier for researchers. Publishers increasingly require ORCiDs as part of the manuscript submission process, but anyone who needs to register with ORCiD can

1. Go to https://orcid.org/. Click on *SIGN IN/REGISTER* at the top right of the page.
2. UNM and NMSU researchers can click on *Institutional account.* Researchers affiliated with other institutions can click on *Personal account.*
3. Researchers can fill out their profile as needed, but for Dryad purposes the ORCiD iD is generated by the system right away.



Figure 3: ORCID sign in



Figure 4: ORCID iD

**Submitting Data to Research Data Services**

As noted above, the UNM Libraries' Research Data Services (RDS) unit can mediate all deposits to Dryad on behalf of researchers. An SMART Grid specific online interface for submitting data to RDS is in development. Until that is ready researchers should:

1. Use the provided CSV metadata template to prepare metadata and documentation. The fields in the template correspond to required fields in Dryad and align with the DataCite metadata schema used by Dryad and other data repositories.
2. Contact RDS at rds@unm.edu to request access to a data upload directory. At this time, the default upload directory will be a shared OneDrive directory but RDS will work with researcher preferences. Alternative methods for submitting data are available.
3. The data upload directory will contain a copy of the CSV metadata template. Researchers can edit this or upload a local copy.
4. One UNM affiliated author will log into Dryad to begin the submission process. RDS will coordinate with that author to complete the submission.
5. RDS will consult with researchers as needed for additional metadata or documentaiton and will submit the dataset to Dryad on the researcher's behalf. Datasets will be published under a "review" status to allow for feedback, updates, and corrections from researchers prior to final publication.

| term | notes | value |
|---|---|---|
| title | required field | |
| author_1_name | required field | |
| author_1_email | required field | |
| author_1_institution | required field | |
| author_1_orcid | required field | |
| abstract | required field; a short description of the dataset | |
| author_2_name | required field if there is a second author | |
| author_2_email | required field if there is a second author | |
| author_2_institution | required field if there is a second author | |
| author_3_name | required field if there is a third author | |
| author_3_email | required field if there is a third author | |
| author_3_institution | required field if there is a third author | |
| *insert additional author info as needed* | | |
| keywords | recommended; controlled vocabulary terms recommended | |
| methods | recommended; upload relevant documents as appropriate | |
| usage_notes | recommended | |
| funding_information | recommended | NSF OIA-1757207 |
| *insert additional funding info as needed* | | |
| related_works | recommended; include publications citing these data, analysis code, other datasets | |
| *insert additional relation info as needed* | | |
| location_information | recommended where appropriate; use lat/long coordinates | |

Figure 5: metadata template

## Overview of the Dryad Submission Process

**Dataset Size Limits**

Dryad provides two options for uploading data. Both options have different size limits:

- When uploading directly from a computer the limit is 10GB per DOI.
- When uploading from a server or shared files service (Google Drive, Box, etc.) the limit is 300GB per DOI.

Note that the limits are per DOI. That means there is a practical limit on the number and size of revisions that may be submitted for any single dataset.

1. After logging in, click the button to *Start a New Dataset*.
2. Indicate whether the data are related to a manuscript, a published article, or "other." The section headed *Dataset: Basic Information* includes the required metadata fields. Note that the submitting author's ORCiD is included. Other authors will receive an email requesting them to add their ORCiD.
3. Additional (recommended) metadata can be provided. Please note that some usage information has been specified project-wide in the DMP.
4. Upload files using one of the two methods described above.
5. The final *Review and Submit* step includes options for enabling private peer review. Note that the default license in Dryad is public domain.
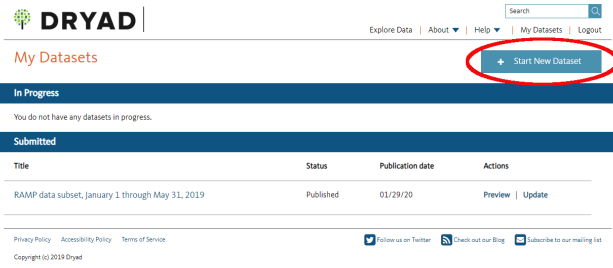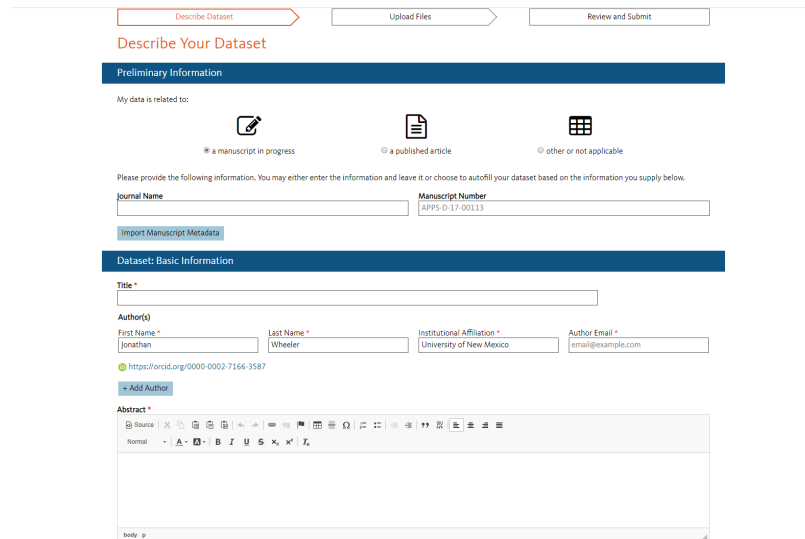
Figure 6: Dryad step 1



Figure 7: Dryad step 2

Figure 8: Dryad step 3

## About DOIs

Datasets published in Dryad are assigned a DOI immediately upon submission. No additional workflow is required for a DOI in Dryad.

DOIs are available for datasets and research products that may be published in repositories other then Dryad. Many repositories will create DOIs, but if not UNM is a DataCite member and RDS can manually generate DOIs as needed.

## Dataset Versioning

Dryad supports dataset versioning. Any dataset can be revised using a process similar to the one outlined above for submitting a new dataset. Dataset revisions in Dryad all have the same DOI, with changes to files included in the *Data Files* section.

Figure 9: Dryad step 4



Figure 10: Dryad step 5

Figure 11: Dryad step 6

Other features to note in the image above include the link to download the entire current version of the dataset, as well as the *Download Data Publication* button, which exports the metadata in PDF format. The PDF includes links to the data files.