# Transparency in Data Management
# Capturing Provenance of Data Curation at BCO-DMO

## BCO-DMO
Biological & Chemical Oceanography Data Management Office

**Adam Shepherd**, Woods Hole Oceanographic Institution | **Amber York**, Woods Hole Oceanographic Institution | **Conrad Schloer**, Woods Hole Oceanographic Institution | **Karen Soenen**, Woods Hole Oceanographic Institution

## BCO-DMO, a domain-specific repository

The Biological and Chemical Oceanography Data Management Office (BCO-DMO) staff work closely with investigators to serve data and information online from research projects funded by the U.S. National Science Foundation (NSF) Biological and Chemical Oceanography Programs, and the Office of Polar Programs Antarctic Organisms & Ecosystems Program.

The goal of this partnership is to effectively curate marine ecosystem data and accompanying documentation, facilitating efficient data discovery and re-use. Throughout the process, BCO-DMO provides services that support specific phases of the data life cycle.

The result is a rich database of research-ready data spanning the full range of marine ecosystem related measurements including in-situ and remotely sensed observations, experimental and model results, and synthesis products. The BCO-DMO system provides access to more than 9000 data sets from more than 900 projects and 2500 researchers.

## Why Capturing Provenance?

- An audit-trail explaining what has changed from one version to another

- Evidence of why domain-specific data management are needed for FAIR-*ness*

- Transparency about how archived data differs from originally submitted version

## Why ad-hoc changes are insufficient?

1. Hard to Maintain
   - code over time has dependencies
     - software & library updates
   - a single DM is the expert; hard for others to help
   - is documentation as text good enough?

2. Duplication of effort
   - same need written multiple times across multiple technologies
   - different DMs write same code for different projects
   - difficult to foster/enforce reuse of code bits
   - are all implementations executing the same procedures? same order?

## Declarative Workflows (DWs) to Capture Provenance

\* Declarative workflows focus on '**what**' to do.
\* Software focuses on '**how**' best to do it.

**Declarative workflows** as a tool for **Standardised Collaboration**:
- *Consistently* express DMs intents
- *Approachable* for non-coders: focus on 'what' should be done over 'how' to do it
- A *shared language*: DMs *understand* each others pipelines

**Declarative workflows** serve as **Workable Data**:
- Configuration data: code interpreting DWs can be changed/swapped without impacting DM intent
- .yaml files because declarative workflows are stored in a data format they can be automatically converted into a provenance record/ machine actionable provenance record.

## Frictionless Data - Data Package Pipelines

Framework to build stream-processing tabular data workflows.

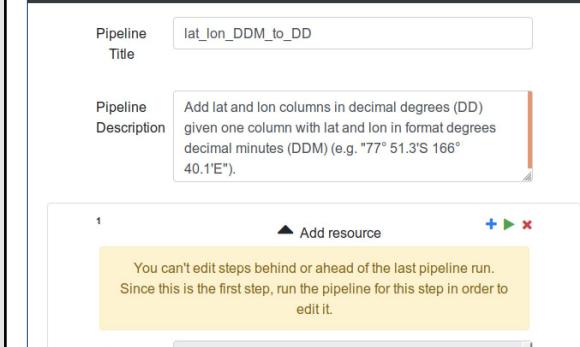- Processing data packages using pipelines/workflows
- Modular components: BCO-DMO built domain-specific processors
- A pipeline is a list of processing steps for a data package.
- Processing steps are defined in a declarative way using YAML.
- 20 processors in the standard library
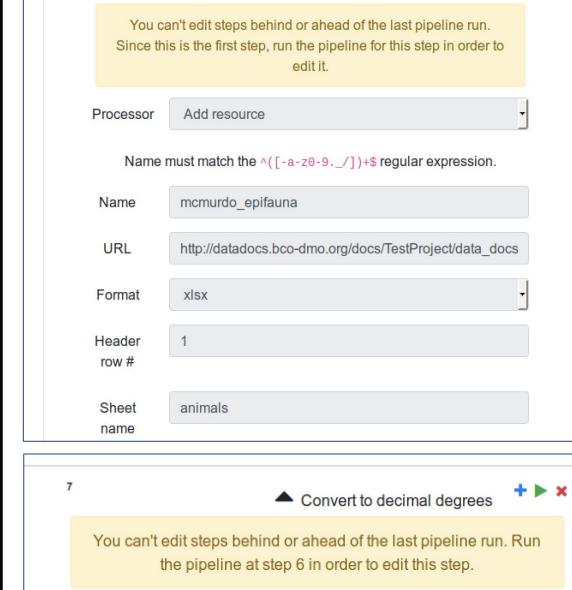- BCO-DMO extended with 18 additional processors

github.com/frictionlessdata/datapackage-pipelines    github.com/BCODMO/bcodmo_frictionless

## BCO-DMO UI for Frictionlessdata Datapackage Pipelines

**Watch Demo !!!**
(YouTube)

- **Encourages consistency** across commonly occurring processing tasks

- **Ensures proper data validation** occurs before data made publicly available

- *Example* pipeline-spec.yaml
- *Example* datapackage.json

Read more: blog.bco-dmo.org/2020/02/09/frictionless-data-pipelines-for-ocean-science

## Elevating pipelines: From Pipeline to PROV

In the BCO-DMO Knowledge Graph, this provenance enables:

1) transparency across data curation process

2) speed, consistency across DM staff

3) query & quality control
   *ex: List all datasets that used the "convert_to_decimal_degrees" processor*

See a provenance record (RDF-Turtle)

Click the AWS lambda to see Python code converts Pipeline to PROV at: github.com/BCODMO/dpp2prov

## Machine readable provenance

### PROV Data Model

*Courtesy of https://www.w3.org/TR/prov-p/*

### PROV-O Extended Data Model

*Courtesy of https://www.w3.org/TR/prov-o/*

## Ontology for Provenance & Plans (P-Plan)

*Courtesy of http://purl.org/net/p-plan#*

## Open Knowledge Foundation

*With support from* Alfred P. Sloan Foundation  Google.org

*for specs:* **frictionlessdata.io/**
*for tooling:* **github.com/frictionlessdata/**

## WOODS HOLE OCEANOGRAPHIC INSTITUTION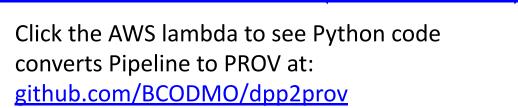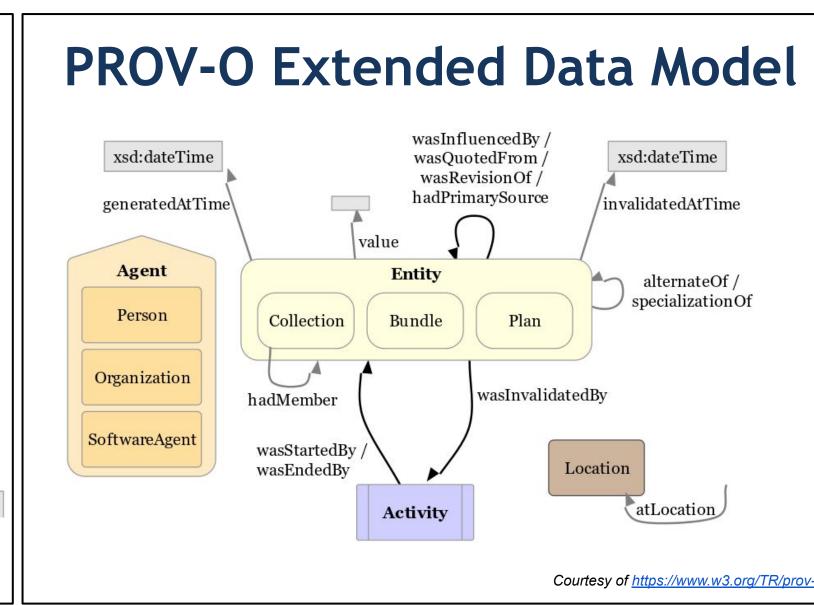