# BCO-DMO
Biological & Chemical Oceanography Data Management Office

# The Frictionless Data Package:
## Data Containerization for Automated Scientific Workflows

**Adam Shepherd**, Woods Hole Oceanographic Institution | **Douglas Fils**, Consortium for Ocean Leadership | **Danie Kinkade**, Woods Hole Oceanographic Institution | **Mak Saito**, Woods Hole Oceanographic Institution
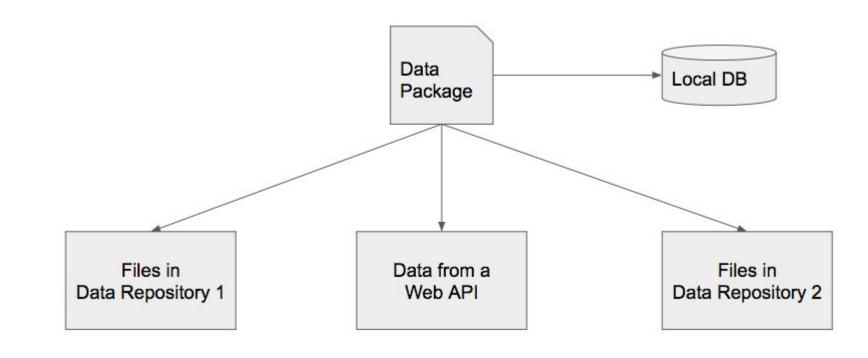
## Abstract

As cross-disciplinary geoscience research increasingly relies on machines to discover and access data, one of the critical questions facing data repositories is how data and supporting materials should be packaged for consumption. Traditionally, data repositories have relied on a human's involvement throughout discovery and access workflows. This human could assess fitness for purpose by reading loosely coupled, unstructured information from web pages and documentation. In attempts to shorten the time to science and access data resources across may disciplines, expectations for machines to mediate the process of discovery and access is challenging data repository infrastructure. This challenge is to find ways to deliver data and information in ways that enable machines to make better decisions by enabling them to understand the data and metadata of many data types. Additionally, once machines have recommended a data resource as relevant to an investigator's needs, the data resource should be easy to integrate into that investigator's toolkits for analysis and visualization.

The Biological and Chemical Oceanography Data Management Office (BCO-DMO) supports NSF-funded OCE and PLR investigators with their project's data management needs. These needs involve a number of varying data types some of which require multiple files with differing formats. Presently, BCO-DMO has described these data types and the important relationships between the type's data files through human-readable documentation on web pages. For machines directly accessing data files from BCO-DMO, this documentation could be overlooked and lead to misinterpreting the data. Instead, BCO-DMO is exploring the idea of data containerization, or packaging data and related information for easier transport, interpretation, and use. In researching the landscape of data containerization, the Frictionlessdata Data Package (http://frictionlessdata.io/) provides a number of valuable advantages over similar solutions. This presentation will focus on these advantages and how the Frictionlessdata Data Package addresses a number of real-world use cases faced for data discovery, access, analysis and visualization.
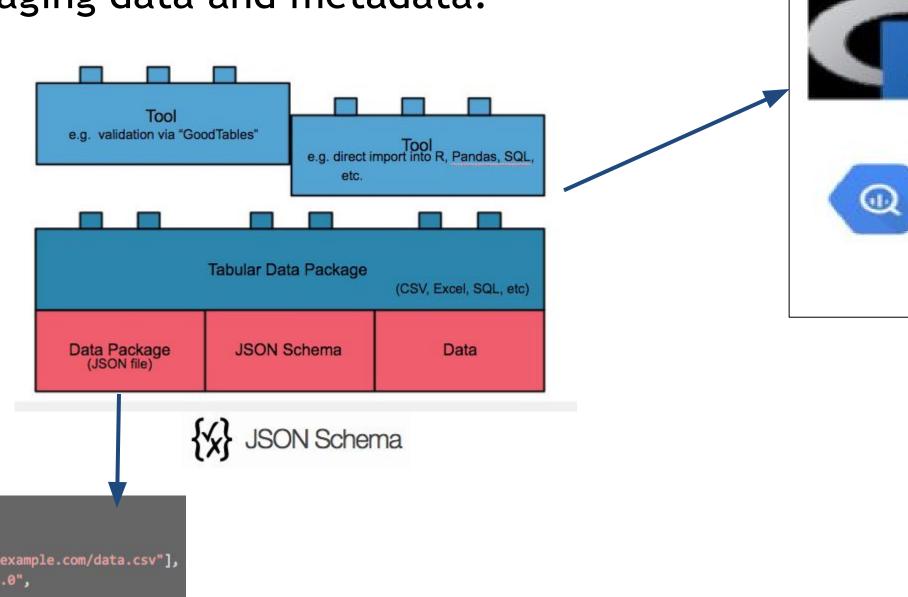
## What is Data Containerization?

**Problem:** Data can be distributed across multiple locations - databases, files on a server, files on the web, or available from web APIs, etc.

The Frictionless Data Package is a set of extendible, lightweight formats for packaging data and metadata.
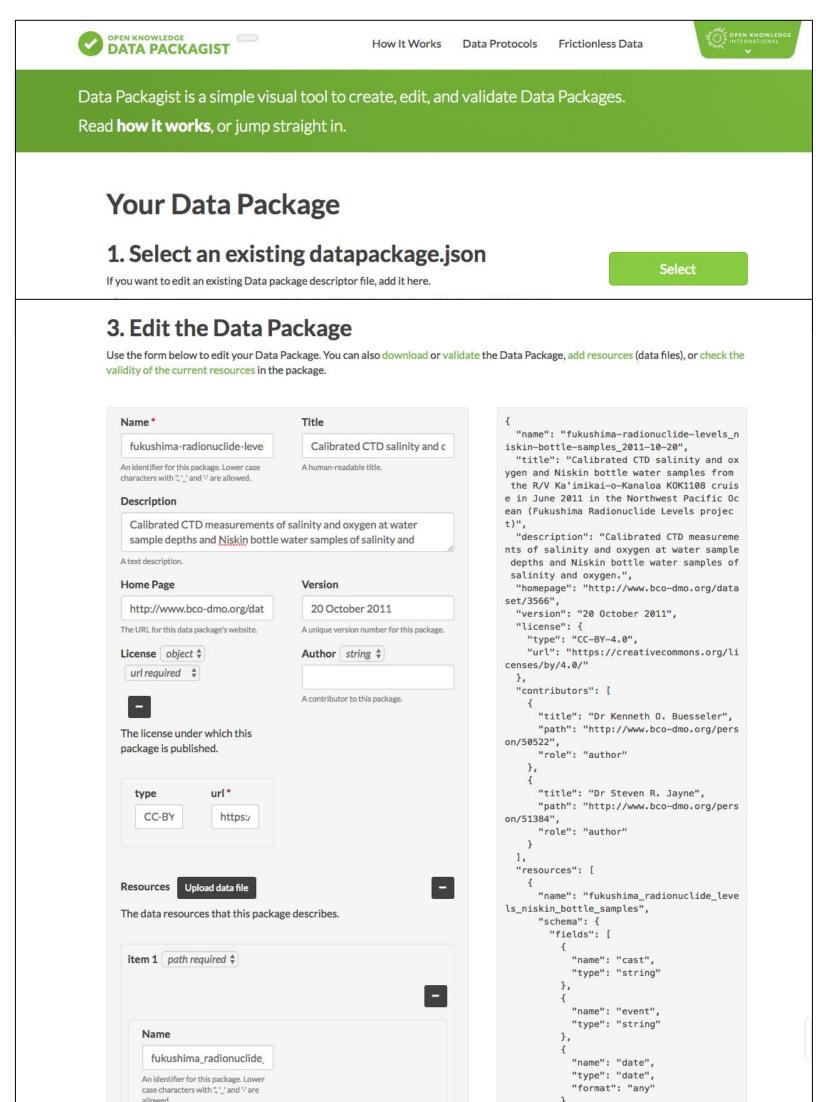


JSON Schema

**BCO-DMO wants:**

- Simpler, expedited data ingest for submitters
- Data transformation that captures provenance
- Continuous integration testing of data holdings

**Problem:** Data can be large so that traditional packages (TAR, ZIP, BagIt) are inefficient transports.

**Q: How do we package data to handle these various locations and sizes?**

datapackage.json
- can be added to traditional packaging formats (BagIt, TAR, ZIP) for describing local files
- can point to data resources that aren't local to the package or file system
- for data accessible by URL, can be the only file needed to be passed in transport
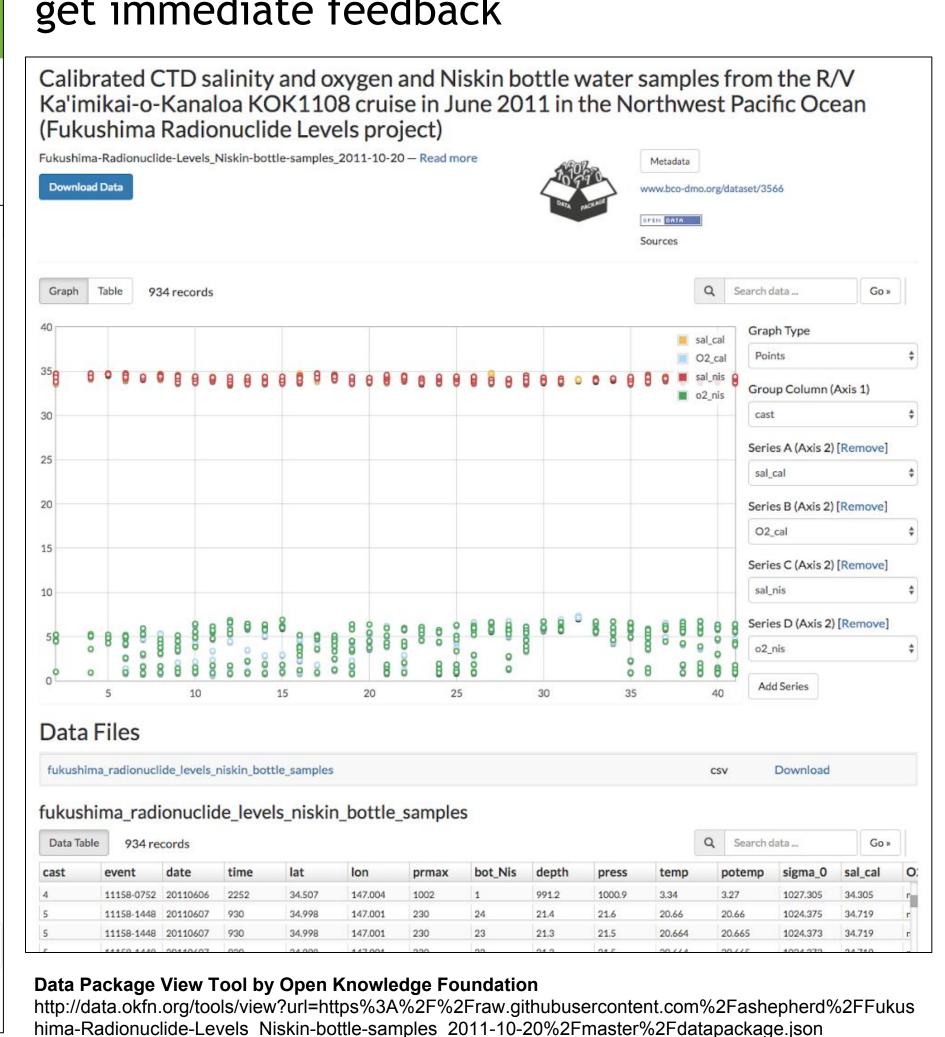
## Simpler, expedited data ingest for submitters

**DataPackagist** - A web service for creating Data Packages.

https://github.com/frictionlessdata/datapackagist

**Data Submitters:** login, describe submission, get immediate feedback

Data Package View Tool by Open Knowledge Foundation
http://data.okfn.org/tools/view?url=https%3A%2F%2Fraw.githubusercontent.com%2Fashepherd%2FFukushima-Radionuclide-Levels_Niskin-bottle-samples_2011-10-20%2Fmaster%2Fdatapackage.json
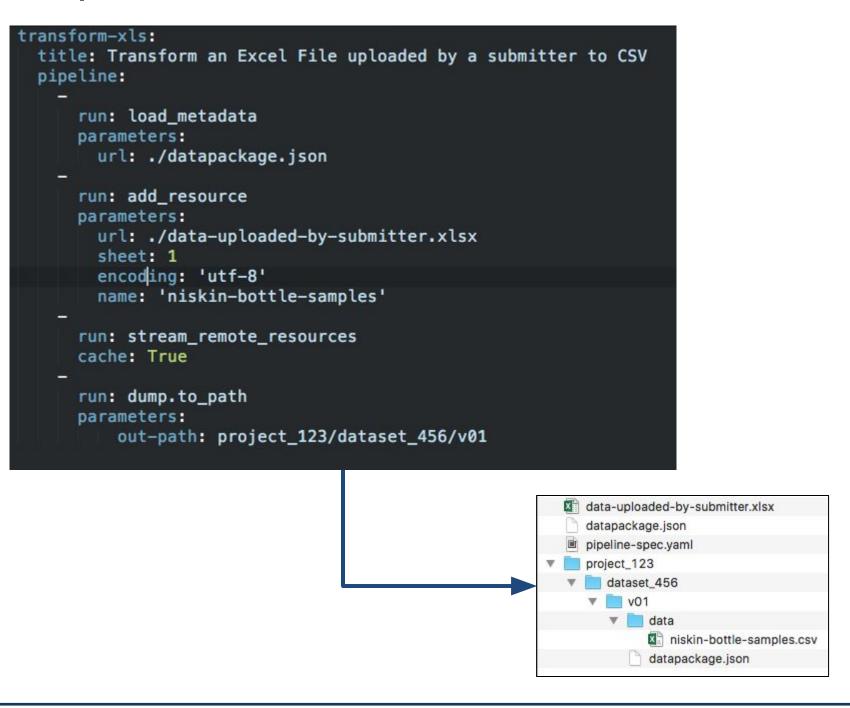
## What about the work of a data manager?

**Data Package Pipelines** - Framework for processing data packages in pipelines of modular components.

https://github.com/frictionlessdata/datapackage-pipelines

- A pipeline has a list of processing steps, and it generates a single data package as its output.
- A pipeline is defined in a declarative way, not in code, stored in a file named *pipeline-spec.yaml*.
- Data Package Pipelines define some common processors, custom processors can be created.

Data Submissions may need to be converted to open formats

Data Managers can extend the datapackage.json to add semantic markup.

The resulting datapackage.json can then be used to populate repository metadata catalog.
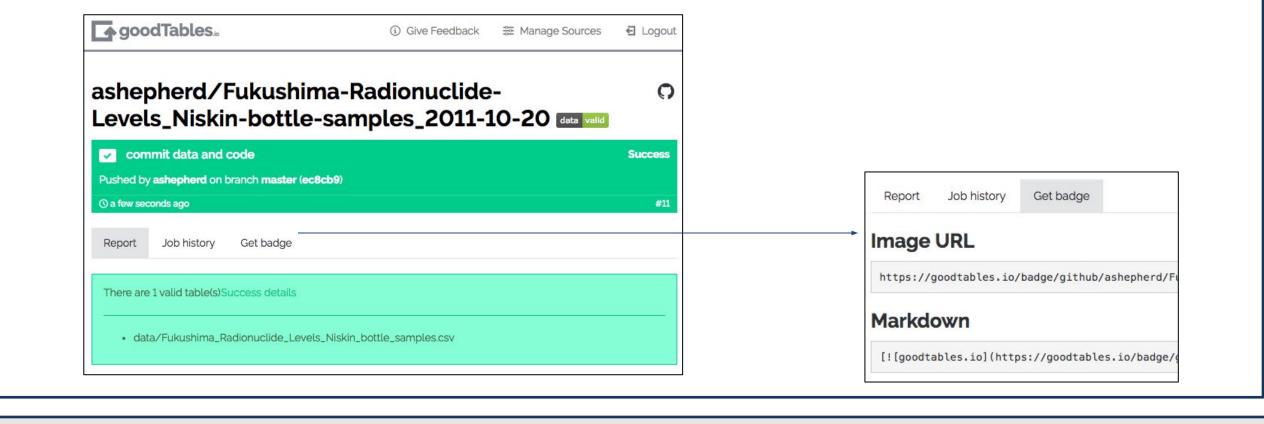
## Continuous Integration Testing for Data

**Goodtables.io** - Continuous data validation as a service.
https://github.com/frictionlessdata/goodtables.io

Because Data Package Pipelines are declarative, *pipeline-spec.yaml* files are **provenance** records.
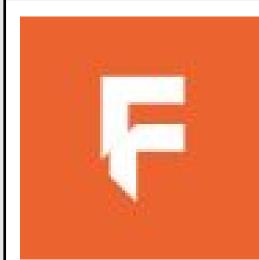
- Pipelines can be re-run to verify that the workflow is reproducible
- Data Packages can be validated
  - Is the datapackage.json well-formed ?
  - Does it meet the JSON schema specification ?
- Tabular Data Packages have deeper validation
  - Column header checking vs. datapackage.json defined fields
  - Methods for handling missing data values
  - Regex processing of a data cell for conforming to a pattern

for specs: **https://frictionlessdata.io/**
for tooling: **https://frictionlessdata.io/software/**

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

for more information on data in context through self-publishing metadata:

**IN33B-0116: Open Core Data approaches to exposing facility data to support FAIR principles**
Wed Dec 13, 1:40PM - 6:00PM
Poster Hall D-F