# DataONE:
# Current Services, New Tools and Future Developments

**Amber Budden**
Director for Community Engagement and Outreach

**Dave Vieglais**
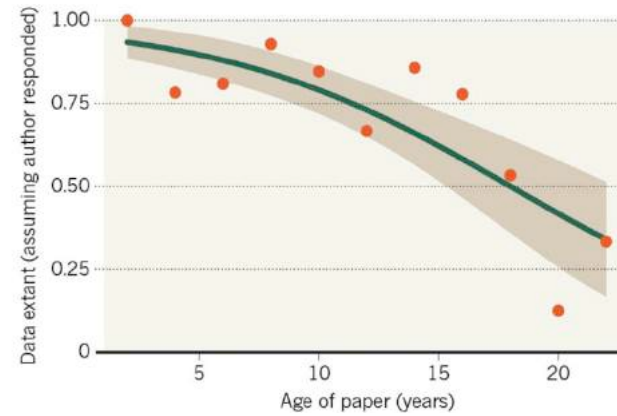Director for Development and Operations

NSF

dataone.org

# Science and Data Challenges





**MISSING DATA**

As research articles age, the odds of their raw data being extant drop dramatically.

Open Government Initiative

My Administration is committed to creating an unprecedented level of openness in Government. We will work together to ensure the public trust and establish a system of transparency, public participation, and collaboration. Openness will strengthen our democracy and promote efficiency and effectiveness in Government.

— PRESIDENT OBAMA, 01/21/09



**Data Sharing by Scientists: Practices and Perceptions**

Metadata standards

| DIF | DwC | DC | EML | FGDC | Open GIS | ISO | My Lab | none |
|-----|-----|-----|-----|------|----------|-----|--------|------|
| 12  | 21  | 26  | 95  | 95   | 96       | 97  | 266    | 676  |

# DataONE
# Vision and Mission

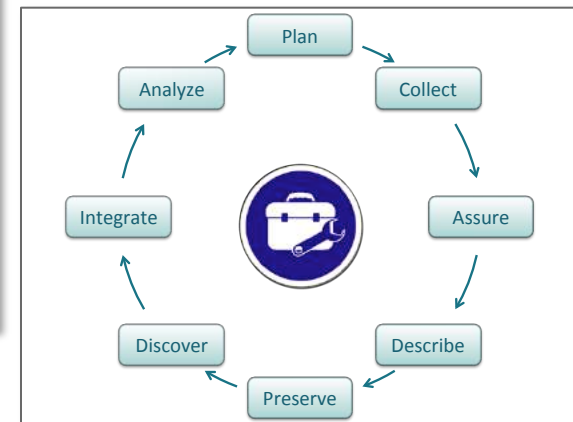*Providing universal access to data about life on earth and the environment that sustains it*

1. Building community

2. Developing sustainable data discovery and interoperability solutions

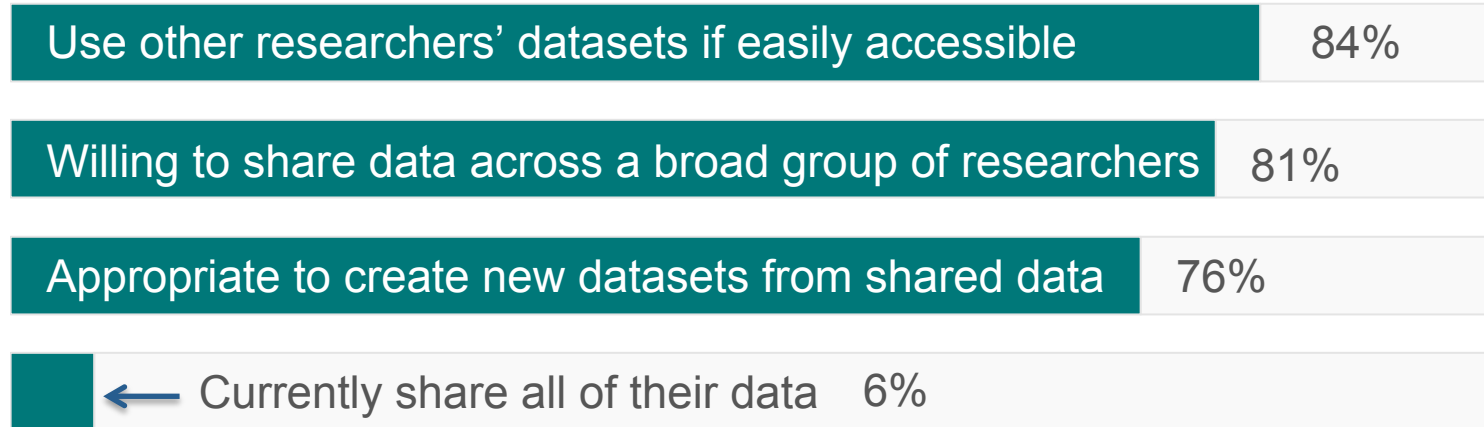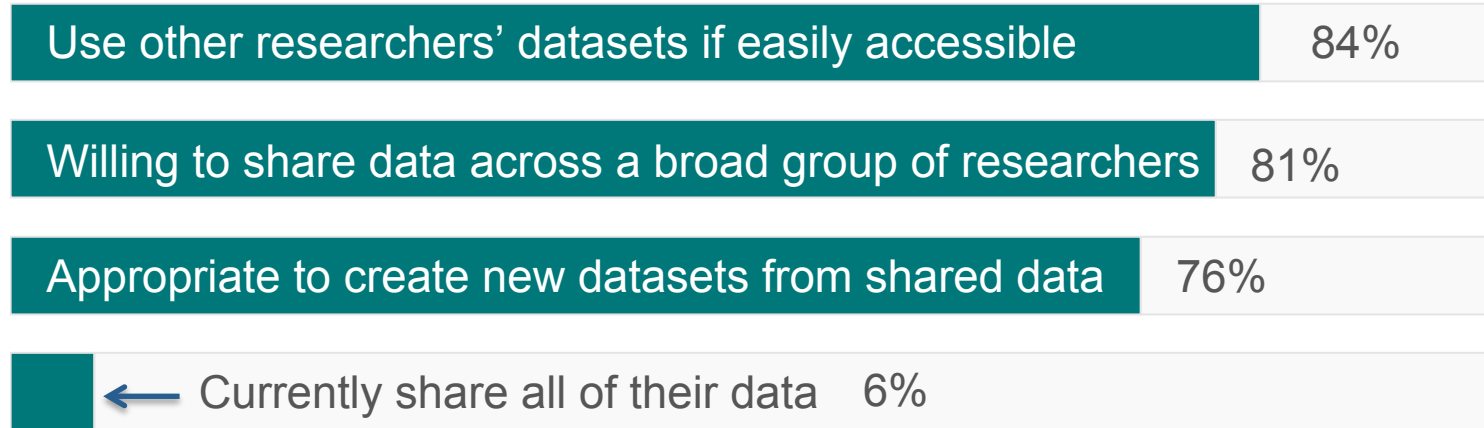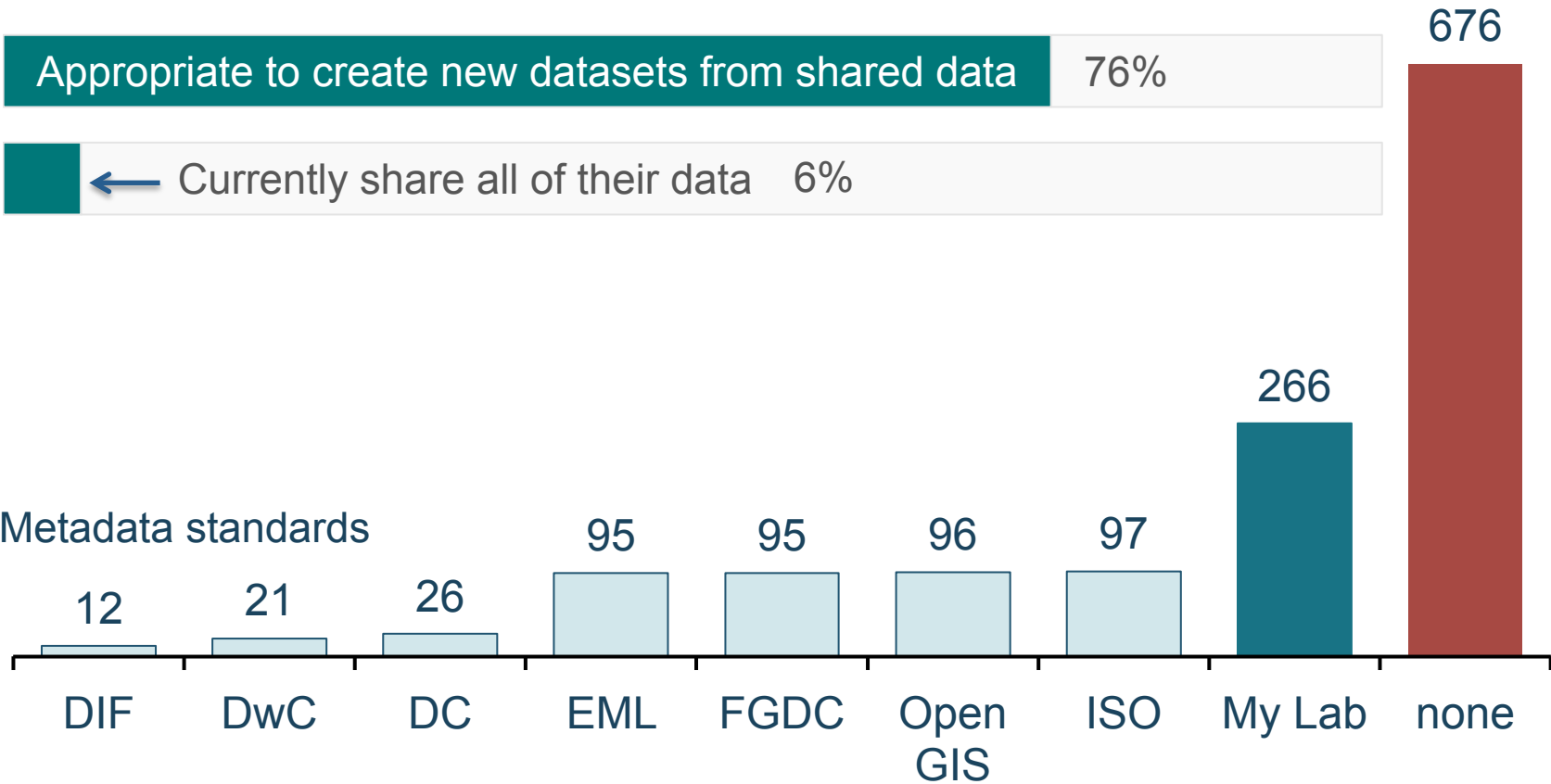3. Supporting researcher tools and services

# Scientists want to share data

| | |
|---|---|
| Use other researchers' datasets if easily accessible | 84% |
| Willing to share data across a broad group of researchers | 81% |
| Appropriate to create new datasets from shared data | 76% |
| ← Currently share all of their data | 6% |

DataONE

# Scientists want to share data

| | |
|---|---|
| Use other researchers' datasets if easily accessible | 84% |
| Willing to share data across a broad group of researchers | 81% |
| Appropriate to create new datasets from shared data | 76% |
| ← Currently share all of their data | 6% |

Metadata standards

| DIF | DwC | DC | EML | FGDC | Open GIS | ISO | My Lab | none |
|---|---|---|---|---|---|---|---|---|
| 12 | 21 | 26 | 95 | 95 | 96 | 97 | 266 | 676 |

DataONE

# Libraries not yet providing data services



| | | |
|---|---|---|
| Metadata creation | 67% | |
| Conversion of data/datasets for ingest | 75% | |
| Selection of data/datasets for ingest | 70% | |
| Selection of data/datasets for repository | 66% | |

never    occasionally    monthly    weekly    daily

# DataONE Cyberinfrastructure
## Coordinating Nodes

Components for a flexible, scalable, sustainable network

**Coordinating Nodes**
- retain complete metadata catalog
- indexing for search
- network-wide services
- ensure content availability (preservation)
- replication services

THE UNIVERSITY of NEW MEXICO

OAK RIDGE National Laboratory

NCEAS

UCSB

UT

www.dataone.org/coordinating-nodes

# DataONE Cyberinfrastructure
## Member Nodes

Components for a flexible, scalable, sustainable network



**Coordinating Nodes**

**Member Nodes**
- diverse institutions
- serve local community
- provide resources for managing their data
- retain copies of data

www.dataone.org/member-nodes

# DataONE Cyberinfrastructure
## Investigator Toolkit

Components for a flexible, scalable, sustainable network



**Coordinating Nodes**

**Member Nodes**

**Investigator Toolkit**

>> command line interface

VT · python · MENDELEY

dash DATA SHARING MADE EASY · zotero

Kepler · ONE R

DataONESearch

DMPTool

www.dataone.org/investigator-toolkit

DataONE

# DataONE Member Nodes
## Current and Upcoming



Upcoming Member Nodes

# Data Holdings

## Uploads

209,359 metadata    393,837 data

The number of individual metadata and data files uploaded over time. Only the first version of each file is counted.

files uploaded: 0, 50,000, 100,000, 150,000, 200,000, 250,000, 300,000, 350,000

2013   2014   2015   2016

## File formats

We breakdown the types of metadata and data files uploaded. Only the most recent version of each file is included.

**209,359 metadata files**

- FGDC 1999 / Other
- EML 2.0.1 — 23,124
- EML 2.1.1 — 26,093
- 005/gmd-noaa — 29,276
- EML 2.1.0 — 31,008
- Dryad 3.1 — 58,445
- FGDC 1998 — 33,106

**393,837 data files**

- image/jpeg / Other
- application/pdf — 33,134
- text/xml — 47,246
- text/plain — 57,626
- text/csv — 116,201
- Application file — 99,574

www.search.dataone.org/#profile

# Three Components

# Three Components
## Integration



DataONE
Service
Specifications

HTTPS REST
XML Messages

ITK

MNs

CNs

DataONE

# purl.dataone.org/architecture

# Shown in a High Level Design

# Objects in DataONE
## Objects and Packaging



**CiTO**
The Citation Typing Ontology

Sys MD

Science Metadata

documents

isDocumentedBy

∞

Data

Sys MD

∞

aggregatedBy

aggregatedBy

**System Metadata**
Identifier
Checksum
Object Type
Size
Access Control
Replication Policy

Resource Map

**OAI-ORE**
Open Archives Initiative Object Reuse and Exchange

Sys MD

# Objects in DataONE
## Everything on Member Nodes

# Objects in DataONE
## Synchronization to Coordinating Nodes

# Coordinating Node Processing
## Indexing for Discovery

# Data Life Cycle

Benjamin Halpern, Melanie Frazier, John Potapenko, Kenneth Casey, Kellee Koenig, Catherine Longo, Julia Lowndes, Cotton Rockwood, Elizabeth Selig, Kimberly Selkoe, and Shaun Walbridge. 2015. **Cumulative human impacts: raw stressor data (2008 and 2013).** KNB Data Repository. doi:10.5063/F1S180FS.

**Copy Citation**

| | Files in this dataset   Package: urn:uuid:57d4a0c5-cffe-4b57-b863-faec725153fa | | | | | **Download all ☁** |
|---|---|---|---|---|---|---|
| 🗀 | Name | | File type | Size | Downloads | |
| 🗎 | Metadata: Cumulative human impacts: raw stressor data (2008 and 2013) | | .xml (EML) | 30 KB | 2092 views | **Download ☁** |
| ⊞ | raw_2008_inorganic_mol.zip | More info ❶ | ZIP folder | 77 MB | 213 downloads | **Download ☁** |
| ⊞ | raw_2013_demersal_nondest_low_bycatch_mol.zip | More info ❶ | ZIP folder | 215 MB | 208 downloads | **Download ☁** |
| ⊞ | raw_2008_artisanal_fishing_mol.zip | More info ❶ | ZIP folder | 46 MB | 218 downloads | **Download ☁** |

▸ **Show 34 more items in this data set**

## General

**Identifier**

raw_2013_uv_mol_20150714095238

**Abstract**

This is a portion of the data used to calculate 2008 and 2013 cumulative human impacts in: Halpern et al. 2015. Spatial and temporal changes in cumulative human impacts on the world's ocean. Seven data packages are available for this project: (1) supplementary data (habitat data and other files); (2) raw stressor data (2008 and 2013); (3) stressor data rescaled by one time period (2008 and 2013, scaled from 0-1); (4) stressor data rescaled by two time periods (2008 and 2013, scaled from 0-1); (5) pressure and cumulative impacts data (2013, all pressures); (6) pressure and cumulative impacts data (2008 and 2013, subset of pressures updated for both time periods); (7) change in pressures and cumulative impact (2008 to 2013). All raster files are .tif format and coordinate reference system is mollweide wgs84. Here is an overview of the calculations: Raw stressor data -> rescaled stressor data (values between 0-1) -> pressure data (stressor data after adjusting for habitat/pressure vulnerability) -> cumulative impact (sum of pressure data) -> difference between 2008 and 2013 pressure and cumulative impact data. This data package includes 2008 and 2013 raw stressor data. The 2008 data includes 18 raster files (preceeded by raw_2008_). The 2013 data includes 19 raster files (preceeded by raw_2013_). There is no sea level rise data for 2008.

**Publication Date**

2015-07-14

# Member Node Profiles

# Use Provenance for Transparency, Reproducibility

What ***input data*** went into this study?

What ***methods*** were used?

… with what ***parameter*** settings, ***calibrations***, …?

*Can we **trust** the data and methods?*

- **Provenance** (*lineage*): track **origin** and **processing history** of data ➔ trust, data quality ~ audit trail for attribution, credit

- **Discovery** of data, methodologies, experiments

# Dataset Provenance



*Record contains provenance information*

# Provenance
## … of Figures

# Provenance
## … of Data

bit.ly/DWS_01_04

# The Problem: Enabling researchers to effectively find data in DataONE

**DataONE:**

**209,300 Metadata Records**
*describing over* **393,000 Data Objects**
*from* **31 Member Nodes**

*... and growing*

# Semantics
## For greater clarity and consistency

Litter?

DataONE

# Displaying semantics of attribute labels

# Manual annotation UI

# Semantic search

## Community Need
## Data Use Metrics

**Challenge**: Data citation and usage reporting are rare, difficult to find, but highly valuable

**Goal**: Index the science literature to provide citation and usage metrics for data and software in DataONE

How interested would you be to know each of the following about the impact of your data?



Kratz and Strasser (2015) doi:10.1038/sdata.2015.39

# Data Use Metrics
## Approach

- Leverage 'Making Data Count' prototype
- Index usage and citation in papers and open access sources
- Powerful reports for users, repositories, and funders

# Data Use Metrics
## Outputs

### For users and repositories:
- Citation and usage services
  - with DataCite
  - interactive displays, reports
- Notification services
  - when cited, by whom…

### For funders:
- Per-award reports
- Program-wide reports
- Impact assessments



graphic from Kratz and Strasser (2015). doi: 10.1038/sdata.2015.39

# Data Use Metrics
## Outcomes

- **Enable Greater Attribution**
  - Article level
  - micro-citation

- **Enhance Resource Discovery**
  - Greater motivation to share
  - More resources to explore

- **Build Community Engagement**
  - Awareness of others' work

- *Promote Reproducible Science*



graphics from: Lin and Fenner, 2013 Information Standards Quarterly; doi: 10.3789/isqv25no2.2013.04

bit.ly/DWS_01_03

# Technical Resources

Architecture and API Documentation
- purl.dataone.org/architecture

Mailing List
- developers@dataone.org

IRC
- irc.ecoinformatics.org #dataone

Subversion, GitHub
- repository.dataone.org/software/cicore
- github.com/DataONEorg

Previous Webinars:
- dataone.org/previous-webinars

# Community Engagement
# Education and Outreach

# Best Practices
## Database and Primer

www.dataone.org/best-practices

# Data Management Modules



www.dataone.org/education-modules

# Screencast Tutorials



www.dataone.org/screencast-tutorials

# DataONE Webinar Series



number of individuals

350
300
250
200
150
100
50
0

64.7%

1  2  3  4  5  6  7  8  9  10  11

webinar event

5 point response scale (1 low)

5

4

3

2

1

Relevance    Speaker    Format
             Knowledge

www.dataone.org/webinars

# Librarian Outreach Kit



www.dataone.org/for-librarians

# Other communication mechanisms

# DataONE Users Group

- A self-organizing, independent group providing feedback to DataONE
- 310 members, 13 member Steering Committee, 2 Co-chairs
- Open participation and membership
- Annual summer meeting co-located with ESIP



www.dataone.org/dataone-users-group

## Save the Date:
## DataONE Users Group Meeting

Please save July 17-18, 2016 for the open DataONE Users Group meeting to be co-located with the Summer ESIP Federation Meeting at the Friday Center, Chapel Hill, North Carolina. The DataONE Users Group (DUG) meeting will be a 2-day event featuring plenary presentations, topical breakout sessions, and community-led discussions.

**There is no registration fee to attend and participate in the DUG meeting.**

Registration and hotel block will open in the spring, a few months before the meeting. Please visit https://www.dataone.org/dataone-users-group for updates and to join the DUG.

### Meeting Theme and Objectives

The 2016 Meeting theme, "**Expanding Data Networks**," will focus on the new challenges and efforts in making data accessible, discoverable, and deliverable while promoting open data policies, standards, and compliance with funders' emerging data management requirements. A strong emphasis is on data synthesis and technological progress made in data network infrastructure.

The scientific program of the 2016 meeting will invite talks and posters on the following topics:
- Leveraging research data level metrics for large data repositories and data networks
- Integrating the needs and inputs of data users to advance and improve data discoverability
- Assessing the progress, impact, and success in promoting open data policies

DataONE encourages DataONE Member Nodes, data scientists, researchers, scientists, students and others to submit abstracts for posters and talks.

### Abstract Submission for Posters and Talks

Please submit an abstract (250 words maximum) to dugchairs@dataone.org and indicate whether you prefer to present a talk or a poster. Talks will be approximately 10-20 minutes in duration, to be confirmed with development of the agenda. The poster session will be held the evening of Sunday July 17th during the reception event.

Submissions will be reviewed by the DataONE Users Group Steering Committee. Accepted abstracts will be published on the DataONE website.

### Important dates

Abstract Submission Deadline**: April 15th 2016**
Author Notification: **May 15th 2016**

DUG Steering Committee: Felimon Gayanilo (co-chair), Plato Smith (co-chair), Steven Aulenbach, Amber Budden, Debora Drucker, Rebecca Koskela, Myrica McCune, Laura Moyers, Shannon Rauch, Robert Sandusky, Stephanie Simms, Heather Soyka

# DataONE Users Group Meeting

July 17-18th 2016
Research Triangle, NC

Theme:
# Expanding Data Networks

# www.DataONE.org

@DataONEorg

facebook.com/DataONEorg

vimeo.com/DataONEorg

slideshare.net/DataONEorg

aebudden@dataone.unm.edu

vieglais@ku.edu