# FAIR and Responsible Data Sharing with Dataverse

Beilstein Open Science Symposium

October 27, 2020

Mercè Crosas, Ph.D., Harvard University
University Research Data Management Officer
Chief Data Science and Technology Officer, IQSS
scholar.harvard.edu/mercecrosas  @mercecrosas

# Dataverse Software Platform Used Worldwide

A Network of Dataverse repositories openly sharing



- Open-source

- 63 installations

- In 6 continents

- 7,350 dataverses

- 135K datasets

- Metadata shared across the Network of Dataverses

dataverse.org

# A Growing And Engaged Community

"The **community** is the one thing that makes this work!"

Sherry Lake, Scholarly Repository Librarian, University of Virginia



- **Core development** at Harvard

- 117 GitHub contributors

- Regular, popular Community meetings, calls, working groups

- **Fast adoption (e.g., in Europe)** w/ multiple language support

- Governance and assistance from the **Global Dataverse Community Consortium**

# Dataverse Repositories Examples

# HARVARD DATAVERSE:

- Open to all researchers and research fields
- 100K searchable datasets
- 800K files
- 20M file downloads
- 92 journal dataverses
- New fee-based data curation services

dataverse.harvard.edu

---



HARVARD
Dataverse

Add Data ▾   Search ▾   About   User Guide   Support   Sign Up   Log In

**Deposit and share your data. Get academic credit.**

Harvard Dataverse is a repository for research data. Deposit data and code here.

[ Add a dataset ✚ ]

**Organize datasets and gather metrics in your own repository.**

A dataverse is a container for all your datasets, files, and metadata.

[ Add a dataverse ✚ ]

**Publishing your data is easy on Harvard Dataverse!**

Learn about getting started creating your own dataverse repository here.

[ **Getting started** ⧉ ]

Find data across research fields, preview metadata, and download files

[ Search over 100,700 datasets... ]   [ 🔍 Find ]   VIEW ALL DATA ❯

Featured   **COVID-19 Data Collection**
A curated collection of COVID-19 data deposited in the Harvard Dataverse repository.

Browse by subject

Agricultural Sciences 3,460          Computer and Information Science 1,506      Medicine, Health and Life Sciences 4,629
Arts and Humanities 2,038            Earth and Environmental Sciences 3,351      Physics 1,211
Astronomy and Astrophysics 915       Engineering 787                             Social Sciences 43,565
Business and Management 675          Law 382
Chemistry 363                        Mathematical Sciences 337

ALL SUBJECTS ❯

Recent datasets

From journal dataverses

Replication data for: Time-Varying Risk Aversion? Evidence from Near-Miss Accidents
Review of Economics and Statistics Dataverse Oct 26, 2020

Replication Data for: Trade Shocks, Firm Hierarchies, and Wage Inequality
Review of Economics and Statistics Dataverse Oct 26, 2020

Replication Data for: Mission and the Bottom Line: Performance Incentives in a Multigoal Organization
Review of Economics and Statistics Dataverse Oct 26, 2020

From other dataverses

Replication Data for: Divided Armies
Project Mars Oct 26, 2020

Replication Data for: "Polarization in America: The Relationship between Affective Polarization and Issue Positions"
Matt Levendusky Dataverse Oct 26, 2020

Occurrence of rust disease caused by Tranzschelia discolor on peach and plum in Brazil
Marlon Henrique Hahn Dataverse Oct 26, 2020

# QUALITATIVE DATA REPOSITORY:
## Dedicated to qualitative and multi-method research

**QDR**
QUALITATIVE DATA
REPOSITORY

CORE
TRUST
SEAL ✓

**Data for: Carbon captured: How business and labor control climate politics**

Version 1.0

Mildenberger, Matto. 2020. "Data for: Carbon captured: How business and labor control climate politics". Qualitative Data Repository. https://doi.org/10.5064/F6GYLSON. QDR Main Collection. V1

☰ Cite Data Project ▾

Learn about Data Citation Standards.

**Data Project Metrics** ❓

1 Download ❓

**Description** ❓

This data project has been published in parallel to Mildenberger, M. 2020. *Carbon Captured: How Business and Labor Control Climate Politics* (Cambridge, MA: MIT Press). In this book, I advance a new theory to explain cross-national differences in the timing and content of climate reforms across advanced industrial economies.

This book seeks to make the sources of its inferences as transparent as possible. In doing so it draws from best practices outlined by the American Political Science Association in its 2012 Council Statement on Openness in Political Science and in a 2013 Guideline for Data Access and Research Transparency for Qualitative Research in Political Science. This data project and its associated Transparency Appendix (TRAX) are key components of this effort. Where permitted by copyright, the data project provides digital copies of grey literatures, policy documents, and media reports referenced in the book or in the book's

# TEXAS DATA REPOSITORY:
## A Consortium of 11 Universities

# DATAVERSENO:
## 9 Universities in Norway

# Key Dataverse Features

# Organization of a Dataverse Repository



**Dataverse**

Collection of datasets
Own administration
Own branding (& can be embedded in your site)

**dataset**

Citation
Metadata
Versioning
Terms/permissions
Collection of Files

**File**

Citation
Preview/Explore
Metadata
Versioning
Permissions

# Features for FAIR and Responsible Data Sharing

✔ Data Citation with DOI for datasets and files

✔ Credit to data authors

✔ Link from data to related article

✔ Standards-based and custom metadata

    ✔ DDI, Schema.org, DataCite, Dublin Core, OAI-ORE, OpenAire, JSON-LD

✔ Access controls (open vs  guestbook vs restricted ) with licenses and terms of use

✔ Versioning and provenance

✔ Descriptive Statistics generated from variables in tabular data files

✔ Conversion to multiple formats of tabular data files

✔ Flexible upload of large data files and hierarchical folder structure: Web UI, API, Standalone Client

✔ Integration with external tools through extensive API

✔ Make Data Count (coming soon)

# FINDABILITY:
## Full, standard data citation automatically generated



Data Citation, with **DataCite DOI**, fully compliant with Force11 **Joint Declaration of Data Citation Principles**

# FINDABILITY AND REUSABILITY:
## Support for multiple metadata standards

Files | Metadata | Terms | Versions

⬆ Export Metadata ▾

**Citation Metadata** ▲

Dublin Core
DDI
DataCite
DDI HTML Codebook
JSON
OAI_ORE
OpenAIRE
Schema.org JSON-LD

| | |
|---|---|
| **Dataset Persistent ID** ❓ | doi:10.7910/DVN/S1EUAF |
| **Publication Date** ❓ | 2020-10-14 |
| **Title** ❓ | Replication Data for: To Emerge? Breadwinning, Motherhood, and Women's Decisions to Run for Office |
| **Author** ❓ | Teele, Dawn (University of Pennsylvania) - ORCID: 0000-0003-3079-3083<br>Bernhard, Rachel (UC Davis)<br>Shames, Shauna |
| **Contact** ❓ | Use email butt...<br>Teele, Dawn (Uni... |
| **Description** ❓ | This is anonymiz... |
| **Subject** ❓ | Social Sciences |
| **Keyword** ❓ | Candidate Emergence, American Politics, Gender and Elections, Campaign Training |
| **Related Publication** ❓ | 202x. Bernhard, Rachel, Shauna Shames, and Dawn Teele. "To Emerge? Breadwinning and Income in Women's Decisions to Run." Forthcoming: American Political Science Review. |
| **Depositor** ❓ | Teele, Dawn |
| **Deposit Date** ❓ | 2020-09-28 |

Rich support for Metadata Standards in **human- and machine-readable formats.**

# ACCESSIBILITY:
## Access terms available for restricted data

# ACCESSIBILITY:
## Metadata always available

# INTEROPERABLE:

## Rich statistics metadata derived from each variable in data file



**Replication Data for: "Climate Amenities, Climate Change, and American Quality of Life" Journal of the Association of Environmental and Resource Economists 3, no. 1 (March 2016): 205-246.**

**ClimateRegressionData_150327.tab**

Albouy, David, Graf, Walter, Kellogg, Ryan, and Wolff, Hendrik, 2018, "Replication Data for: "Climate Amenities, Climate Change, and American Quality of Life" Journal of the Association of Environmental and Resource Economists 3, no. 1 (March 2016): 205-246.", https://doi.org/10.7910/DVN/QCE1XY, Harvard Dataverse, V1, UNF:6:CBIOoHJrG5/T6i+XjwBVwg== [fileUNF]

Extensive variable metadata (descriptive statistics) **automatically** derived from tabular data file in DDI format

# INTEROPERABLE:
## Enable use of standard vocabularies



FAIR Controlled Vocabularies can be added in **metadata template**

# REUSABLE:
## Licenses, Terms, and Tiered-Access to Data



HARVARD
Dataverse

Add Data ⌄   Search ⌄   About   User Guide   Support   Sign Up   Log In

## Replication Data for: To Emerge? Breadwinning, Motherhood, and Women's Decisions to Run for Office

Version 1.0

Teele, Dawn; Bernhard, Rachel; Shames, Shauna, 2020, "Replication Data for: To Emerge? Breadwinning, Motherhood, and Women's Decisions to Run for Office", https://doi.org/10.7910/DVN/S1EUAF, Harvard Dataverse, V1, UNF:6:gAlQI8fH9OpP/AdvZlo/1A== [fileUNF]

Cite Dataset ⌄      Learn about Data Citation Standards.

Access Dataset ⌄

Contact Owner     Share

Dataset Metrics ❓

4 Downloads ❓

**Description** ❓          This is anonymized replication data. (2020-09-28)
**Subject** ❓              Social Sciences
**Keyword** ❓              Candidate Emergence, American Politics, Gender and Elections, Campaign Training
**Related Publication** ❓  202x. Bernhard, Rachel, Shauna Shames, and Dawn Teele. "To Emerge? Breadwinning Women's Decisions to Run." Forthcoming: American Political Science Review.

- CC0 as default waiver for open data;
- Optional Licenses and custom Terms;
- Tiered-access for restricted data

| Files | Metadata | Terms | Versions |

Terms of Use ⌃

**Waiver** ❓         Our Community Norms as well as good scientific practices expect that proper credit is given via citation. Please use the data citation above, generated by the Dataverse.

CC0 - "Public Domain Dedication"   [PUBLIC DOMAIN]

Guestbook ⌃

**Guestbook** ❓      No guestbook is assigned to this dataset, you will not be prompted to provide any information on file download.

# What's Next

# Sensitive Data

# Support for Sensitive Data in Dataverse

## Non-Sensitive DataTags (TODAY)

- Data uploaded to Dataverse via one of the current options
- Stored locally

| | |
|---|---|
| Blue | Publicly open, no barriers |
| Green | Publicly open, but need to register to access |
| Yellow | Restricted, need to be granted permissions, but non-sensitive |

## Sensitive DataTags (FUTURE RELEASE)

- Stored in a **Trusted Remote Storage Agent,** accessed through notary service
- Metadata published in Dataverse

| | |
|---|---|
| Orange | Requires Data Use Agreement (DUA); requires data enclave *(moderate sensitivity)* |
| Red | Requires DUA; stricter security requirements and audits *(high sensitivity)* |
| Crimson | Only metadata and no link to data; data stored outside network *(maximum sensitivity)* |

# Dataverse + Impact + OpenDP

**Public Repository**



Blue  Green  Yellow

Dataset and file metadata

Data Use Agreement (DUA)
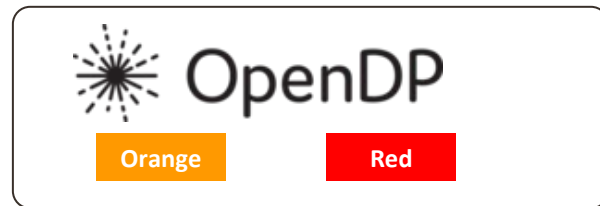Set by Data Owner
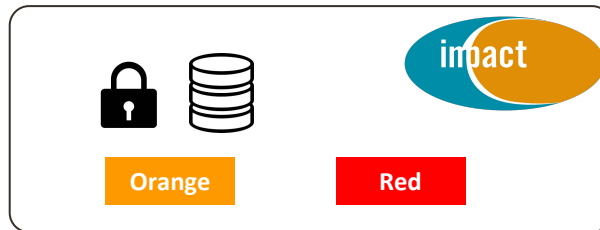
*Find sensitive dataset*

**Data User**

**Notary Service**



Differentially private statistical release of the data

**Secure Compute Environment to run Differentially Private statistics**



Orange  Red

**Trusted Remote Storage Agents (TRSA) or data enclaves**



Orange  Red

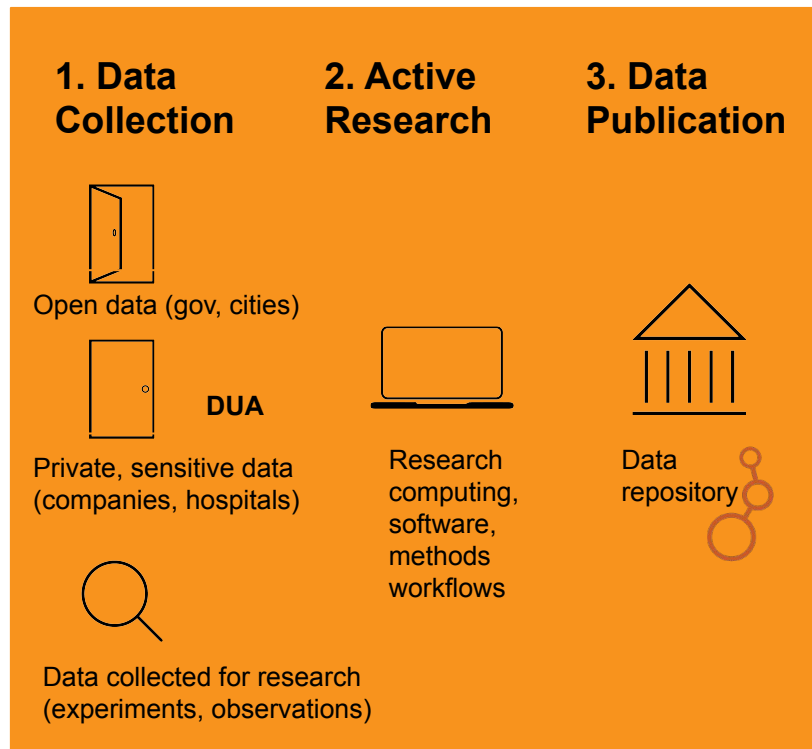Sensitive Data Files
Large Data Files

# Towards a Data Commons

"a **data commons** brings together (or co-locates) **data with cloud computing** infrastructure and **commonly used** software services, tools & applications for **managing, analyzing and sharing data** to create an **interoperable** resource for a research community."

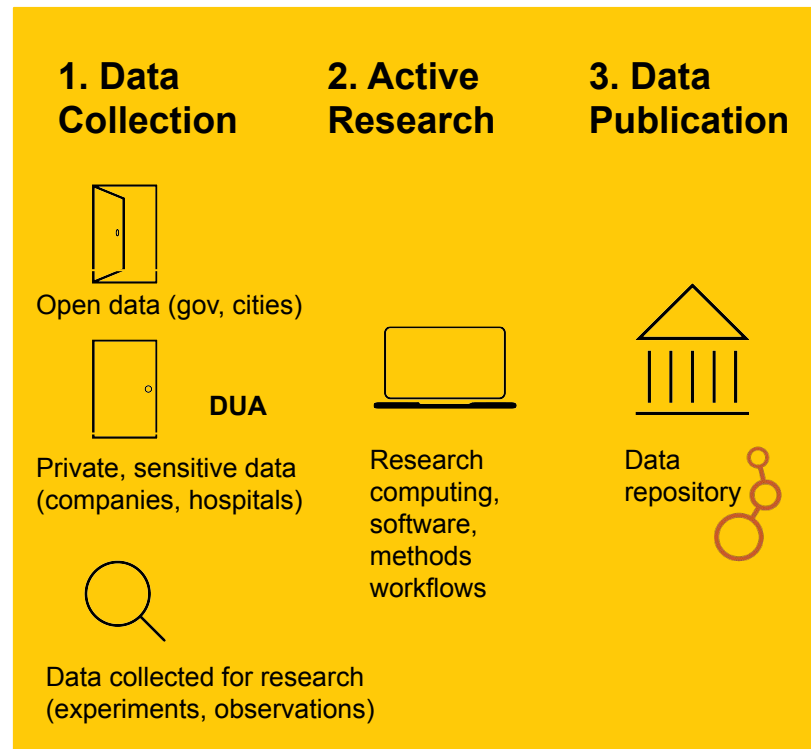[Robert Grossman, on the NIH Data Commons Consortium initiative]

# The Problem

- During Active Research, researcher A doesn't know about data from Researcher B; **not easy to collaborate**
- Research lifecycle steps 1, 2, & 3 are **duplicative** and **not connected**; not easy to manage data throughout lifecycle and publish research output
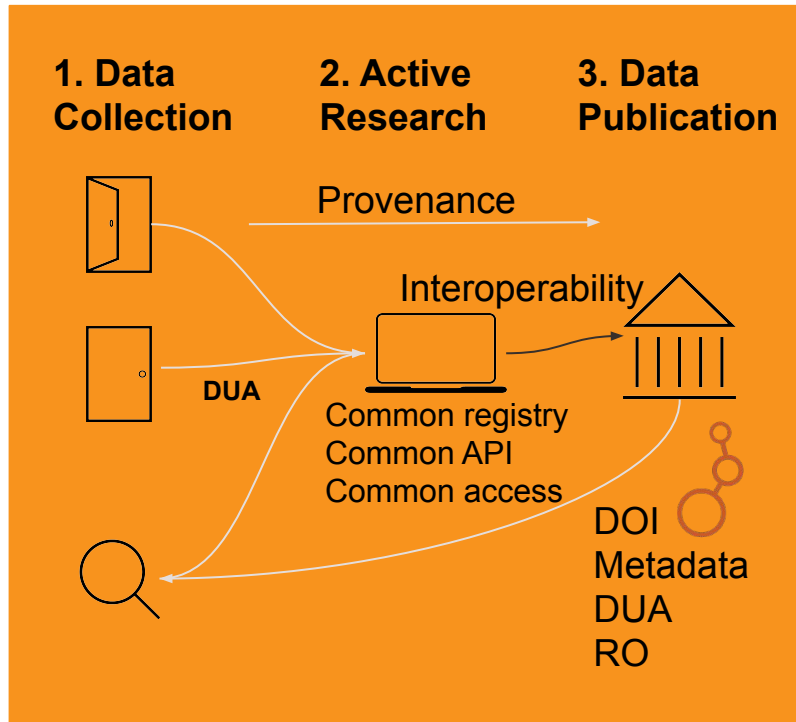
# A Solution for an Institution
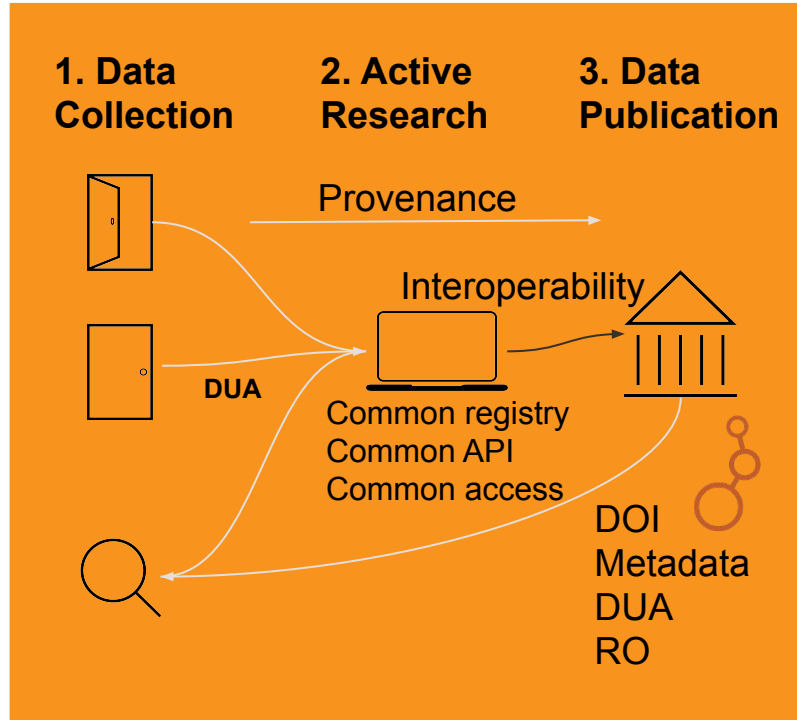
### Institution A



**A Commons vision with Dataverse:**

- **Common registry** with **metadata**, so Researcher A can find and access data from Researcher B during Active Research

- **Common API** to integrate with research tools and computing to access the data seamlessly

- **DUA and controlled access** tracked and throughout lifecycle and shared through Dataverse

- **Provenance** tracked throughout the lifecycle to produce reproducible research outputs – package data, code, workflows in a Research Object(RO)
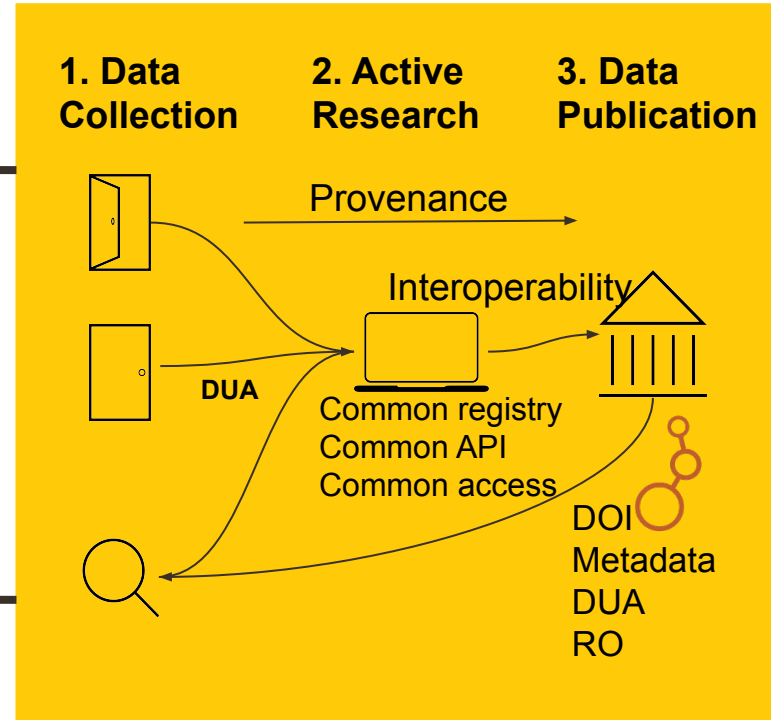
# THANKS

Mercè Crosas, Ph.D., Harvard University
scholar.harvard.edu/mercecrosas  @mercecrosas
dataverse.org