Mercè Crosas, Ph.D.
Chief Data Science and Technology Officer
Institute for Quantitative Social Science (IQSS)
Harvard University

@mercecrosas mercecrosas.com

# Data should be Findable, Accessible, Interoperable, Reusable (FAIR) by machines

Comment | OPEN

## The FAIR Guiding Principles for scientific data management and stewardship

Mark D. Wilkinson, Michel Dumontier, IJsbrand Jan Aalbersberg, Gabrielle Appleton, Myles Axton, Arie Baak, Niklas Blomberg, Jan-Willem Boiten, Luiz Bonino da Silva Santos, Philip E. Bourne, Jildau Bouwman, Anthony J. Brookes, Tim Clark, Mercè Crosas, Ingrid Dillo, Olivier Dumon, Scott Edmunds, Chris T. Evelo, Richard Finkers, Alejandra Gonzalez-Beltran, Alasdair J.G. Gray, Paul Groth, Carole Goble, Jeffrey S. Grethe, Jaap Heringa, Peter A.C 't Hoen, Rob Hooft, Tobias Kuhn, Ruben Kok, Joost Kok, Scott J. Lusher, Maryann E. Martone, Albert Mons, Abel L. Packer, Bengt Persson, Philippe Rocca-Serra, Marco Roos, Rene van Schaik, Susanna-Assunta Sansone, Erik Schultes, Thierry Sengstag, Ted Slater, George Strawn, Morris A. Swertz, Mark Thompson, Johan van der Lei, Erik van Mulligen, Jan Velterop, Andra Waagmeester, Peter Wittenburg, Katherine Wolstencroft, Jun Zhao & Barend Mons ✉ - Show fewer authors

Wilkinson et al, 'The FAIR Guiding Principles scientific data management and stewardship," Nature Scientific Data, 2016; NIH Data Commons Principles;  Joint Declaration of Data Citation Principles (Force11)

"**FAIR Principles** put specific emphasis on enhancing the ability of machines to automatically find and use the data, in addition to supporting its reuse by individuals."

"**Good data management** is not a goal in itself, but rather is the key conduit leading to knowledge discovery and innovation, and to subsequent data and knowledge integration and reuse by the community after the data publication process."

# FAIR Data Principles in Brief

- **To be Findable:**
  - (meta)data are assigned a globally unique and persistent identifier
  - data are described with rich metadata
  - metadata clearly and explicitly include the identifier of the data it describes
  - (meta)data are registered or indexed in a searchable resource

- **To be Accessible:**
  - (meta)data are retrievable by their identifier using a standardized communications protocol
  - the protocol is open, free, and universally implementable
  - the protocol allows for an authentication and authorization procedure, where necessary
  - metadata are accessible, even when the data are no longer available

- **To be Interoperable:**
  - (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.
  - (meta)data use vocabularies that follow FAIR principles
  - (meta)data include qualified references to other (meta)data

- **To be Reusable:**
  - meta(data) are richly described with a plurality of accurate and relevant attributes
  - (meta)data are released with a clear and accessible data usage license
  - (meta)data are associated with detailed provenance
  - (meta)data meet domain-relevant community standards

# We built Dataverse to incentivize data sharing, with "good data management" in mind

- An **open-source** platform to share and archive data

- Developed at **Harvard's Institute for Quantitative Social Science** since 2006

- Gives **credit and control** to data authors & producers

- Builds a **community** to:

  - define new standards and best practices

  - foster new research in data sharing and reproducibility

- Has brought **data publishing** into the hands of data authors

The **Dataverse** Project

# Dataverse is now a widely used repository platform

**21 installations around the world**
**Used by researchers from > 500 institutions**
**60,000 datasets in Harvard Dataverse repository**
**http://dataverse.org**



The **Dataverse** Project

# Dataverse has a growing, engaged community of developers and users

**38**
GitHub contributors

**332**
members in the community list

**23**
community calls
so far with **239**
participants from
**8** countries

Annual Community Meeting,
with **200** attendees

# Dataverse implements FAIR Data Principles

- **Data Citation with global persistent IDs:**
  - generate DOI automatically
  - attribution to data authors and repository
  - registration to DataCite
- **Rich Metadata:**
  - citation metadata
  - domain-specific descriptive metadata
  - variable and file metadata (extracted automatically)
- **Access and usage controls:**
  - open data as default, with CC0 waiver
  - custom terms of use and licenses, when needed
  - data can be restricted, but citation & metadata always publicly accessible
- **APIs and standards:**
  - SWORD, OAI-PMH, native API to search and get data and metadata
  - Dublin Core and DDI metadata standards
  - PROV ontology standard to capture provenance of a dataset (coming soon)

The **Dataverse** Project

# Dataset Landing Page

# Dataset Landing Page

# Dataset Landing Page

# Dataset Landing Page

# Standard file formats and automatic metadata extraction allow data exploration

| Var1 | Var2 | Var3 | Var4 |
|------|------|------|------|
|      |      |      |      |
|      |      |      |      |
|      |      |      |      |
|      |      |      |      |

TwoRavens: summary stats & analysis



geospatial variable

| Var1 | Var2 | Var3 | Var4 |
|------|------|------|------|
|      |      |      |      |
|      |      |      |      |
|      |      |      |      |
|      |      |      |      |

WorldMap: geospatial exploration



The Dataverse Project

Led by Boston University, the MOC is a collaborative effort among BU, Harvard, UMass Amherst, MIT, and Northeastern University, as well as the Massachusetts Green High-Performance Computing Center (MGHPCC) and Oak Ridge National Laboratory (ORNL)
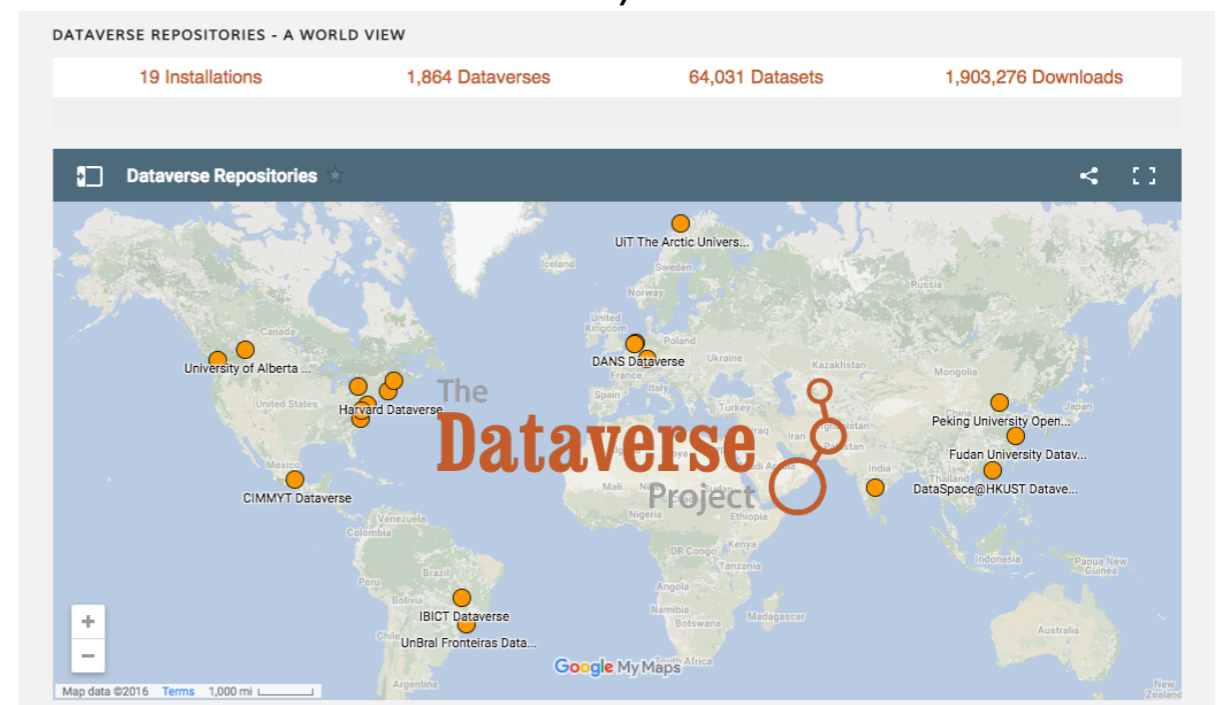
# Cloud Dataverse:



- is a collaboration between **Massachusetts Open Cloud (MOC)** and **Dataverse**

- will allow replication of data in multiple storage locations, and access to cloud computing

Dataverse led by Harvard's IQSS

# Dataverse Now          with Cloud Dataverse

MOC

Data depositor

Data users

**Access** object in Swift +
**Compute** with Sahara/Hadoop

Publish dataset

download

The Dataverse Project
**Repository**

The Dataverse Project

**Data
Replication**

Metadata

Data files

Data + metadata

Swift
Object
Store

# Future Dataset Landing Page: Access to Compute

Each Dataverse repository can choose to enable the Cloud Dataverse option

http://privacytools.seas.harvard.edu



http://datatags.org



**DataTags:**

- is part of the Harvard University Privacy Tools Project

- will be integrated with Dataverse

http://dataverse.org



Open source research data repository software

A **datatag** is a set of security features and access requirements for file handling.

A **datatags repository** is one that stores and shares data files in accordance with a standardized and ordered levels of security and access requirements

# DataTags Levels

| Tag Type | Description | Security Features | Access Credentials |
|---|---|---|---|
| Blue | Public | Clear storage, Clear transmit | Open |
| Green | Controlled public | Clear storage, Clear transmit | Email- or OAuth Verified Registration |
| Yellow | Accountable | Clear storage, Encrypted transmit | Password, Registered, Approval, Click-through DUA |
| Orange | More accountable | Encrypted storage, Encrypted transmit | Password, Registered, Approval, Signed DUA |
| Red | Fully accountable | Encrypted storage, Encrypted transmit | Two-factor authentication, Approval, Signed DUA |
| Crimson | Maximally restricted | Multi-encrypted storage, Encrypted transmit | Two-factor authentication, Approval, Signed DUA |

*DataTags and their respective policies*
Sweeney L, Crosas M, Bar-Sinai M. Sharing Sensitive Data with Confidence: The Datatags System. Technology Science. 2015.

# Requirements for a DataTags Repository

1. Supports more than one datatag
2. Each file in the repository must have one and only one datatag
   - ◉ additional requirements cannot weaken the file security
   - ◉ and cannot required the same or more security than a more restrictive datatag
3. A recipient of a file from the repository must
   - ◉ satisfy file's access requirements,
   - ◉ produce sufficient credentials as requested,
   - ◉ and agree to any terms of use required to acquire the file.
4. Provides technological guarantees for requirements 1, 2 and 3.

# Repositories today do not have an easy, standard way to support sensitive data

**Dataverse** Project Terms of Use

"User Uploads must be void of all identifiable information, such that re-identification of any subjects from the amalgamation of the information available from all of the materials (across datasets and dataverses) uploaded under any one author and/or user should not be possible."

# We are making Dataverse a DataTags Repository to share sensitive data

# The DataTags automated interview ...



www.datatags.org/interviews/FlowChartSet-1/start

**DataTags** 🚩 Feedback  🏠 ⬛ IQSS

**Question: Please select one a...**

Do the data concern living p

[ Yes ] [ No ]

**Answer Feed**

This space will contain your answer
when n...

**Engine Trace**
Useful for debugging the interview

| Id | Type | I... |
|----|------|------|
| $0 | ask | D... |

DataTags project. © 201...

---

www.datatags.org/interviews/FlowChartSet-1/q/medicalRecord...

**DataTags** 🚩 Feedback  🏠 ⬛ IQSS

**Question: Please selec...**

Do the data contain h

[ Yes ] [ No ]

**Answer Feed**

Do the data concern living...

**Engine Trace**
Useful for debugging the interv...

| Id |
|----|
| medicalRecordsCompliance... |
| $4 |
| $3 |
| $0 |

---

www.datatags.org/interviews/questionnaireId/q/MR7

**DataTags** 🚩 Feedback  🏠 ⬛ IQSS

**Question: Please select one answer**

Do the data contain information from a covered entity or business associate of a covered entity?

Terms

**Business associate**
A business associate is any person or organization, including a subcontractor, that acts on behalf of, or provides services to, a covered entity involving the use or disclosure of protected health information. This includes, but is not limited to, legal, actuarial, accounting, consulting, claim processing, data analysis, administration, utilization review, quality assurance, billing, benefit management, practice management, and re-pricing activities.

**Covered entity**
A covered entity is a health plan, health care clearinghouse, or health care provider that transmits any health information in electronic form.

- Health plans include health insurance companies, health maintenance organizations [HMOs], company health plans, and government programs that pay for health care, such as Medicare, Medicaid, and the military and veterans health care programs.
- Health care providers include doctors, clinics, psychologists, dentists, chiropractors, nursing homes, and pharmacies.
- Health care clearinghouses include entities that process nonstandard health information they receive from another entity into a standard, i.e. standard electronic format or data content, or vice versa.

[ Yes ] [ Not Sure ] [ No ]

**Answer Feed**

Do the data contain information related to substance abuse diagnosis, referral, or treatment?  [ No ] [ ↻ Revisit ]

**Current Tags**

**DataTags**

Code [ green ] 🛈

# … helps you generate a machine-readable datatag for your data

# Thanks!

**Learn more at http//<u>dataverse.org</u>**

@mercecrosas    <u>mercecrosas.com</u>