# Sharing Sensitive Data with Confidence

Mercè Crosas, Ph.D.
Chief Data Science and Technology Officer
Institute for Quantitative Social Science
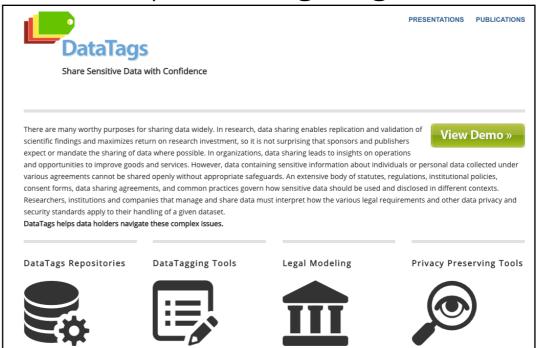Harvard University
Twitter: @mercecrosas  Web: mercecrosas.com

RDA 8th Plenary, Denver, September 16, 2016

http://privacytools.seas.harvard.edu

http://datatags.org



http://dataverse.org

**DataTags:**

- is part of the Harvard University Privacy Tools Project
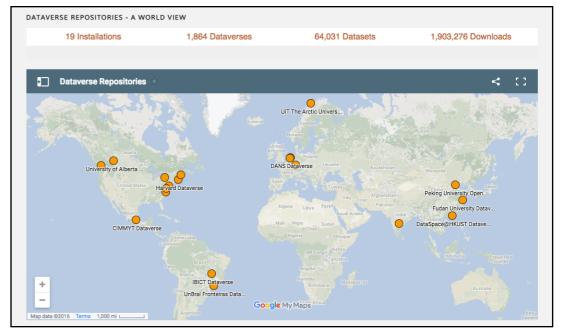- now is being integrated with Dataverse

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

# Requirements for a DataTags Repository

1. Supports more than one datatag

# Requirements for a DataTags Repository

1. Supports more than one datatag
2. Each file in the repository must have one and only one datatag

# 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

# 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

# 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

# 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,

# 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,

# 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,
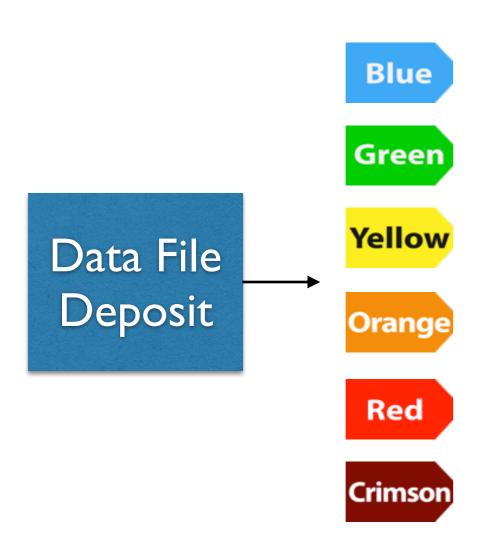    - and agree to any terms of use required to acquire the file.

# 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.
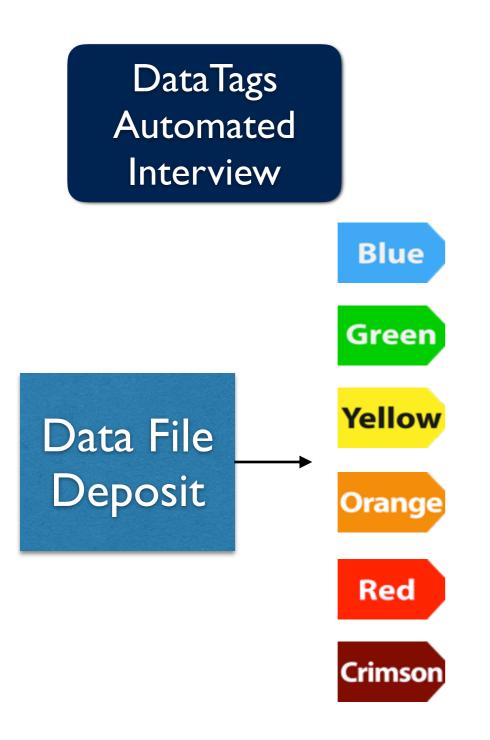4. Provides technological guarantees for requirements 1, 2 and 3.

# Repositories today do not have an easy, standard way to support sensitive data
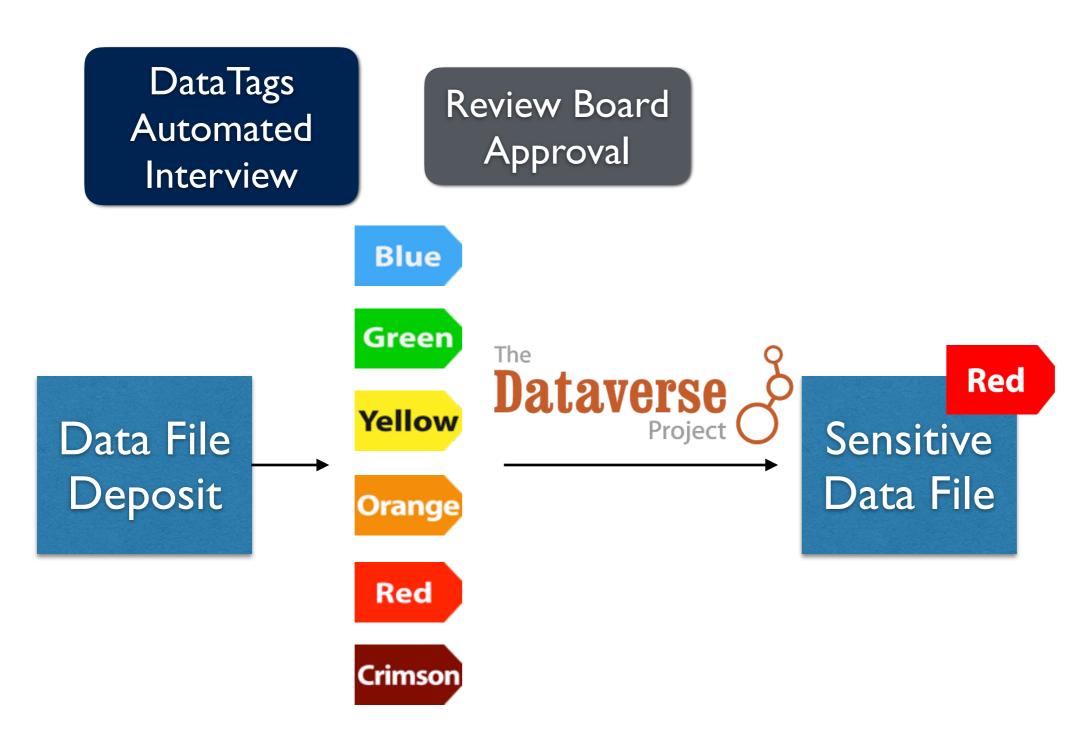
**Dataverse** Project Terms of Use

"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."
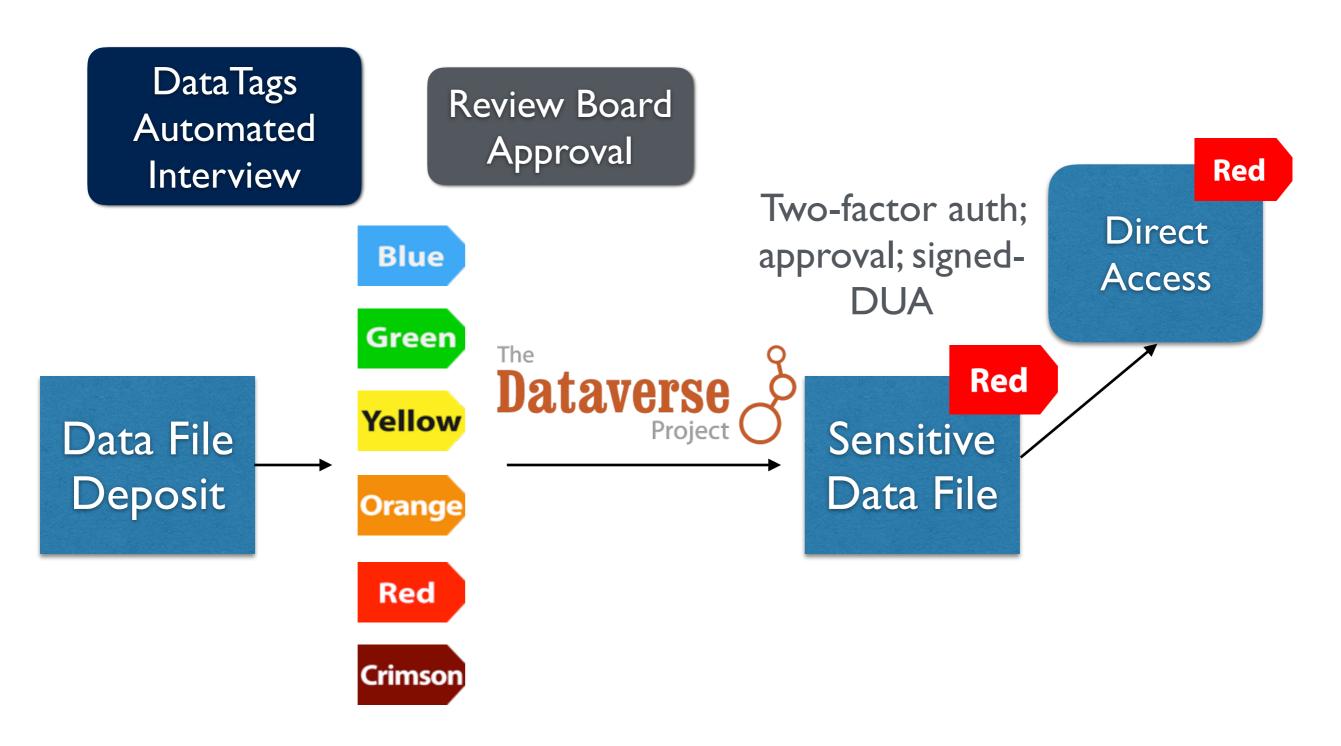
# Making Dataverse a DataTags Repository to share sensitive data

# Making Dataverse a DataTags Repository to share sensitive data

Data File Deposit

# Making Dataverse a DataTags Repository to share sensitive data

# Making Dataverse a DataTags Repository to share sensitive data

# Making Dataverse a DataTags Repository to share sensitive data



DataTags Automated Interview

Review Board Approval

Blue

Green

Yellow

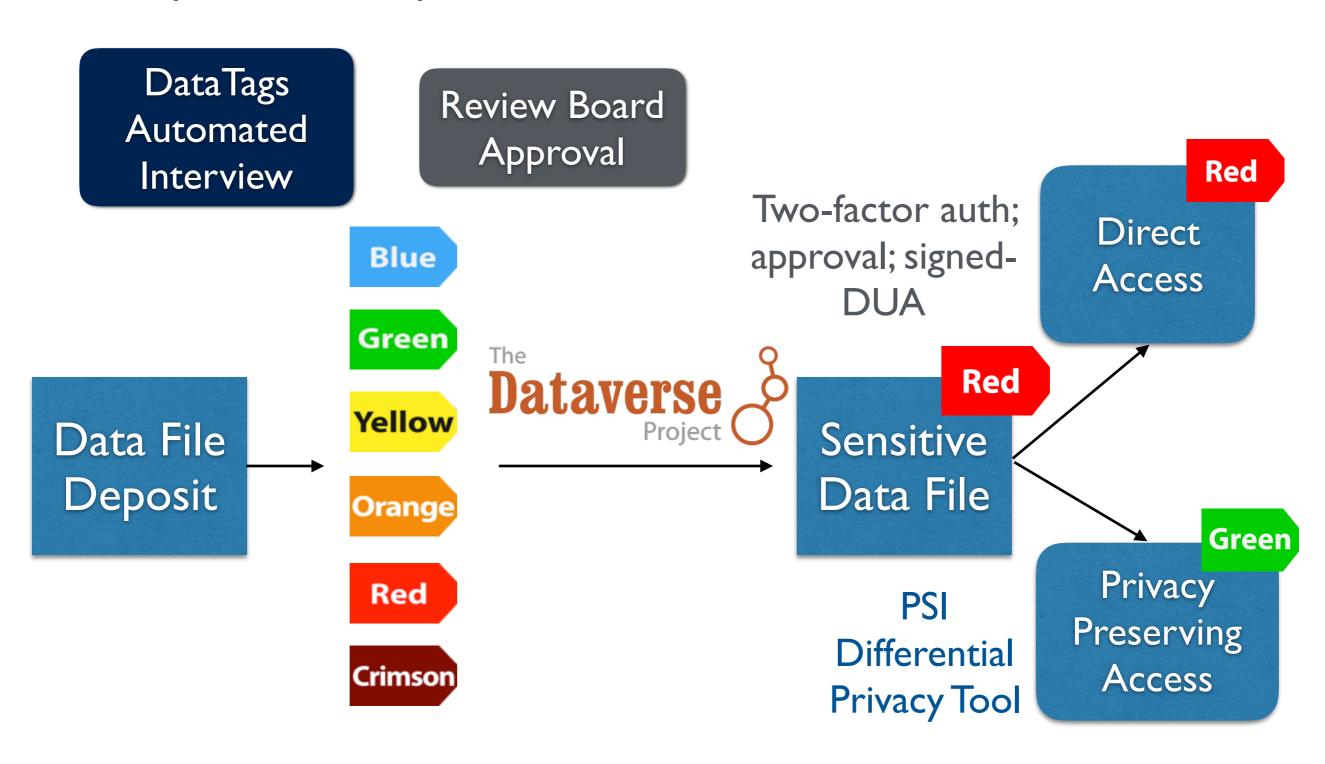Data File Deposit

Orange

Red

Crimson

# Making Dataverse a DataTags Repository to share sensitive data



DataTags Automated Interview

Review Board Approval

Blue
Green
Yellow
Orange
Red
Crimson

The Dataverse Project

Data File Deposit

Sensitive Data File

Red

# Making Dataverse a DataTags Repository to share sensitive data



DataTags Automated Interview

Review Board Approval

Blue

Green

Yellow

Orange

Red

Crimson

The Dataverse Project

Data File Deposit

Two-factor auth; approval; signed-DUA

Sensitive Data File — Red

Direct Access — Red

# Making Dataverse a DataTags Repository to share sensitive data



DataTags Automated Interview

Review Board Approval

Blue

Green

Yellow

Orange

Red

Crimson

The Dataverse Project

Data File Deposit

Sensitive Data File — Red

Two-factor auth; approval; signed-DUA

Direct Access — Red

PSI Differential Privacy Tool

Privacy Preserving Access — Green
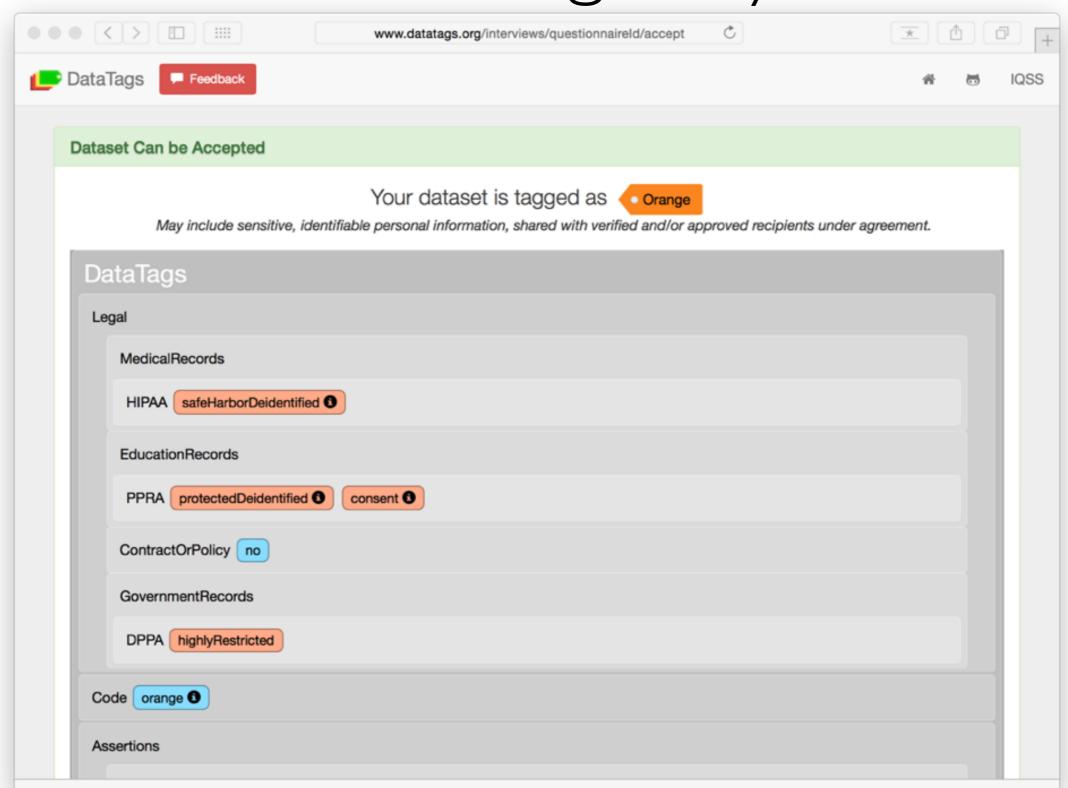
# The DataTags automated interview ...

# The DataTags automated interview ...

# The DataTags automated interview ...

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

**References:**

Sweeney L, Crosas M, Bar-Sinai M. Sharing Sensitive Data with Confidence: The Datatags System. Technology Science. 2015. Technology Science

Bar-Sinai M, Sweeney L, Crosas M. DataTags, Data Handling Policy Spaces, and the Tags Language. Proceedings of the International Workshop on Privacy Engineering. 2016. IEEE Symposium on Security and Privacy

http://datatags.org (Demo of DataTags automated interview)

http://privacytools.seas.harvrd.edu

http://dataverse.org