Microsoft® Research
# Faculty Summit 2010

# The Dataverse Network – Combining Preservation with Scholarly Recognition

Mercè Crosas, Ph.D.
Director of Product Development
Institute for Quantitative Social Science (IQSS)
Harvard University

# Acknowledgements

This work is done with contribution from:

- The Dataverse Network software development team at IQSS: Gustavo Durand, Ellen Kraffmiller, Kevin Condon, Leonid Andrev, Bob Treacy, Michael Heppler, Steve Kraffmiller

- Gary King, Albert J. Weatherhead III University Professor, Harvard University

- Micah Altman, Senior Research Scientist, IQSS, Harvard University

# The Problem with Data

**Professional archives** focus on long term access by the wider community

- Persistent identifiers
- Fixity
- Backups and recovery
- Metadata standards
- Conversion standards
- Preservation standards
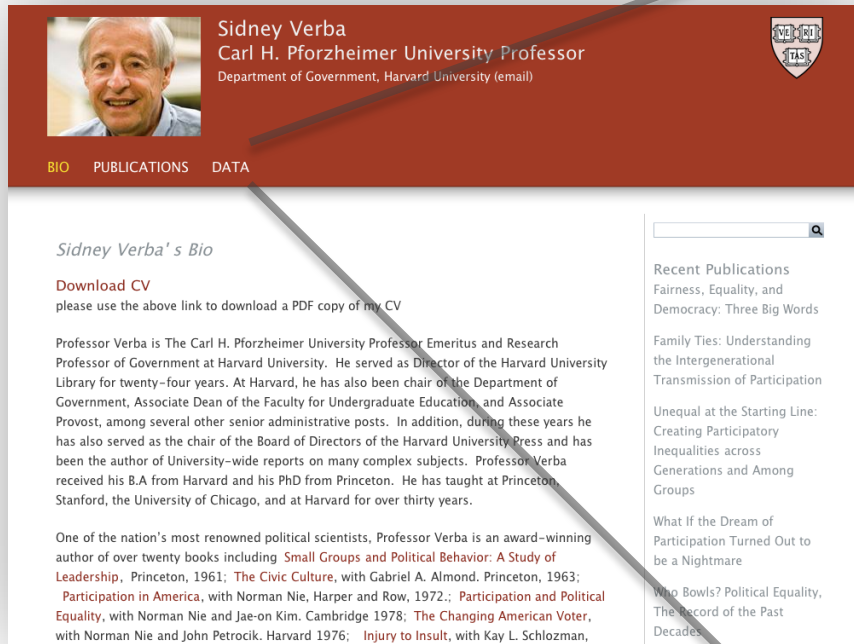
$\neq$

**Data owners** focus on recognition and control

- Branding and visibility
- Data discovery
- Ease of use
- Scholarly citation
- Control over updates
- Terms of access and use

*… but do not attract maximum contribution from data owners*

*… but do not assure long-term use as would a professional archive*

# A Solution through technology

**Centralized** Archival Infrastructure

**Distributed** Recognition and Control

- Persistent identifiers
- Fixity
- Backups and recovery
- Metadata standards
- Conversion standards
- Preservation standards

**+**

- Branding and visibility
- Data discovery
- Ease of use
- Scholarly citation
- Control over updates
- Terms of access and use

**The Dataverse Network**

*Enables people to solve a political problem through technology – Combining centralized archiving with distributed ownership*

# Use Case 1: Dataverse for Researchers – Web Visibility and Branding



**Your website**

**Your dataverse**

# Use Case 1: Dataverse for Researchers – Recognition through Citatition

- The researcher is the administrator of his/her dataverse.
- Manages data curation workflow.
- Sets permissions and terms of use to datasets.

**A Researcher ...**

| Creates Dataverse | → | Adds Study | → | Sets Permissions | → | Releases Study |

Adds Study
└ Enters Cataloging
└ Uploads Data Files

Sets Permissions
└ Terms of Use

# Use Case 2: Dataverse for Journals – Replication of Published Work

## Journal, Book

## Datasets, text, images



**Formal Data Citation**

**Please use this citation**

```
Jacquelyn Boone James; Rosalind C. Barnett, 2003, "The Aspirations and Experiences of
Undergraduates, 2000", http://hdl.handle.net/1902.1/00092  UNF:3:0sjRWertkBtT9wlgcuP4Fg==
Murray Research Archive [Distributor] V1 [Version]
```

**Persistent identifier and url that never changes**

**Universal Numerical Fingerprint (UNF) to verify dataset**

**Citation For**

- A Journal for replication data can be set as an open dataverse.
- Authors upload data after registering.
- Journal reviews data before releasing it.

```
┌─────────────────┐      ┌─────────────────┐      ┌─────────────────┐      ┌─────────────────┐
│  Author Adds    │ ───> │  Sends Study    │ ───> │    Curator      │ ───> │    Releases     │
│  Study          │      │  for Review     │      │  reviews data   │      │    Study        │
└─────────────────┘      └─────────────────┘      └─────────────────┘      └─────────────────┘
         │
┌─────────────────┐
│  Enters         │
│  Cataloging     │
└─────────────────┘
         │
┌─────────────────┐
│  Uploads        │
│  Data Files     │
└─────────────────┘
```

The Dataverse provides:

- **Versioning** – Keeps all changes to data and metadata to allow citation of old versions.

- **Deaccession** – Doesn't delete studies permanently so citation is valid forever.

- **Universal Numeric Fingerprints** – Fixity that survives changes of format.

**Dataset File**

| Variable and Format Information | | | | | |
|---|---|---|---|---|---|
| 1 | 4 | 4 | 21 | ⋯ | 121 |
| 1 | 2 | 2 | 91 | ⋯ | 212 |
| 1 | 9 | 2 | 72 | ⋯ | 104 |
| 0 | 2 | 2 | 2 | ⋯ | 321 |
| 1 | 6 | 2 | 12 | ⋯ | 204 |
| 1 | 9 | 4 | 52 | ⋯ | 311 |
| 0 | 3 | 2 | 23 | ⋯ | 92 |
| 0 | 2 | 5 | 91 | ⋯ | 212 |
| 0 | 5 | 8 | 91 | ⋯ | 91 |
| 1 | 9 | 1 | 72 | ⋯ | 104 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋱ | ⋮ |
| 1 | 2 | 2 | 91 | ⋯ | 212 |

Preservation and verifiable format by:

1. Extracting variable and format metadata.

2. Converting dataset to a format independent of software package.

3. Applying cryptographic algorithm to canonical format.

4. Obtaining an alphanumeric string based on semantic content:
   UNF5:EKgHvTNfkkS86dNzABIhNw==

DDI

Dublin Core

FGDC

Marc

Export metadata in multiple formats

Harvest metadata with OAI-PMH

Dataverse Network

Dataverse Network

Archival Site

Replicate data and metadata with LOCKSS

The Dataverse facilitates:

- Multiple copies of data
- Metadata exchange with other systems.

**New: A 'wiki' dataverse is an archive to share and improve data through contributions (with a moderator to review them before their release)**
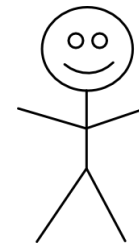
- Multiple contributors can be updating the same study.
- Changes are only visible after being reviewed.
- A new way to gather more data for a study and improve its cataloging.

**Update Metadata or Upload Data**

**Review Changes**

**Release Changes**

or

**Don't accept Changes**

Contributors

Curators

# Current Usage

- The IQSS Dataverse Network at Harvard holds:
  - 257 public dataverses from around the world
  - 36,000 studies (mostly in Social Science)
  - 640,000 files
- ... *and* there are over 10 more installations of the Dataverse Network at other Institutions.

# How to Get Started …

- Create a dataverse at:  http://dvn.iq.harvard.edu
- Or install the Dataverse Network software for your Institution.
- For more information, go to: http://thedata.org
- Contact me at: mcrosas@hmdc.harvard.edu


*THANK YOU!*

Microsoft® Research
Faculty Summit 2010