

# DATA SHARING FOR BETTER SCIENCE

## THE DATAVERSE PROJECT

Mercè Crosas, Institute for Quantitative Social Science, Harvard University

 @mercecrosas

MAX PLANCK INSTITUTE FOR RADIOASTRONOMY, SEPTEMBER 12, 2017

# THIS TALK

- **Importance of Data Sharing**
  - **Reproducibility** to verify science
  - **Reuse** to advance science and evidence-based policy
- **Enabling Data Sharing**
  - **Data Policies** from journals and funding agencies
  - **Data Citation** to find datasets, give credit to data authors
  - **Data Repositories** as publishers of data

# DATA SHARING, DATA PUBLISHING

Data sharing is "the release of research data, associated metadata, accompanying documentation, and software code for re-use and analysis in such a manner that they can be discovered on the Web and referred to in a unique and persistent way."



Data Publishing Group, 2015 RESEARCH DATA ALLIANCE

Since the Beginning of Modern Science ...

# **NULLIUS IN VERBA**

**"TAKE NOBODY'S WORD FOR IT"**

(motto of the Royal Society, founded in 1660,  
launched first scientific journal in 1665)



**AUG 22 2017**  
**LEAVE A COMMENT**

BY JOHN BORGH  
BEST PRACTICES

## WHAT WE TALK ABOUT WHEN WE TALK ABOUT REPRODUCIBILITY



*Campbell's Soup Cans* (1962) by Andy Warhol. Created by replicating an existing object and then reproducing the process at least 32 times.

University of California Curation Center, DataPub blog, August 2017

## REPRODUCIBILITY AND REPLICATION (BY THE NATIONAL SCIENCE FOUNDATION):

The ability of a researcher to duplicate the results of a prior study ... using the same materials and procedures used by the original investigator. **(reproducibility)**

... if the same procedures are followed but new data are collected. **(replication)**

## EMPIRICAL, COMPUTATIONAL, AND STATISTICAL REPRODUCIBILITY (STODDEN, 2014):

**Empirical:** data and collection details are made freely available

**Computational:** code, software, hardware and implementations details are provided

**Statistical:** details on choice of statistics tests, model parameters are provided

**REPRODUCIBILITY CRISIS?**

# TRUST, BUT VERIFY

The Economist

World politics

Business & finance

Economics

Science & technology

Culture

Problems with scientific research

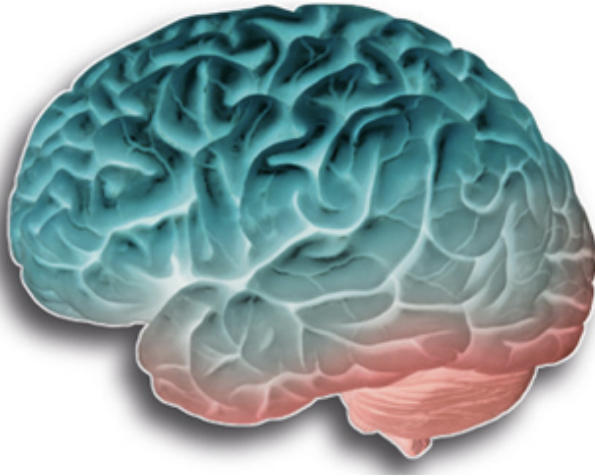
## How science goes wrong

Scientific research has changed the world. Now it needs to change itself

Oct 19th 2013 | From the print edition

Like 15k

Tweet 1,120



**6 (11%) out of 53 landmark cancer biology studies could be reproduced.**

# Science

AAAS

Home

News

Journals

Topics

Careers

Science

Science Advances

Science Immunology

Science Robotics

Science Signaling

Science Translational Medicine

SHARE

RESEARCH ARTICLE



0



19

## Estimating the reproducibility of psychological science

Open Science Collaboration<sup>†</sup>  
+ See all authors and affiliations

Science 28 Aug 2015;  
Vol. 349, Issue 6251, aac4716  
DOI: 10.1126/science.aac4716

Article

Figures & Data

Info & Metrics

eLetters

PDF

You are currently viewing the abstract.

View Full Text



**39 out of 100 psychology studies could be reproduced.**

**The Washington Post**

*Democracy Dies in Darkness*

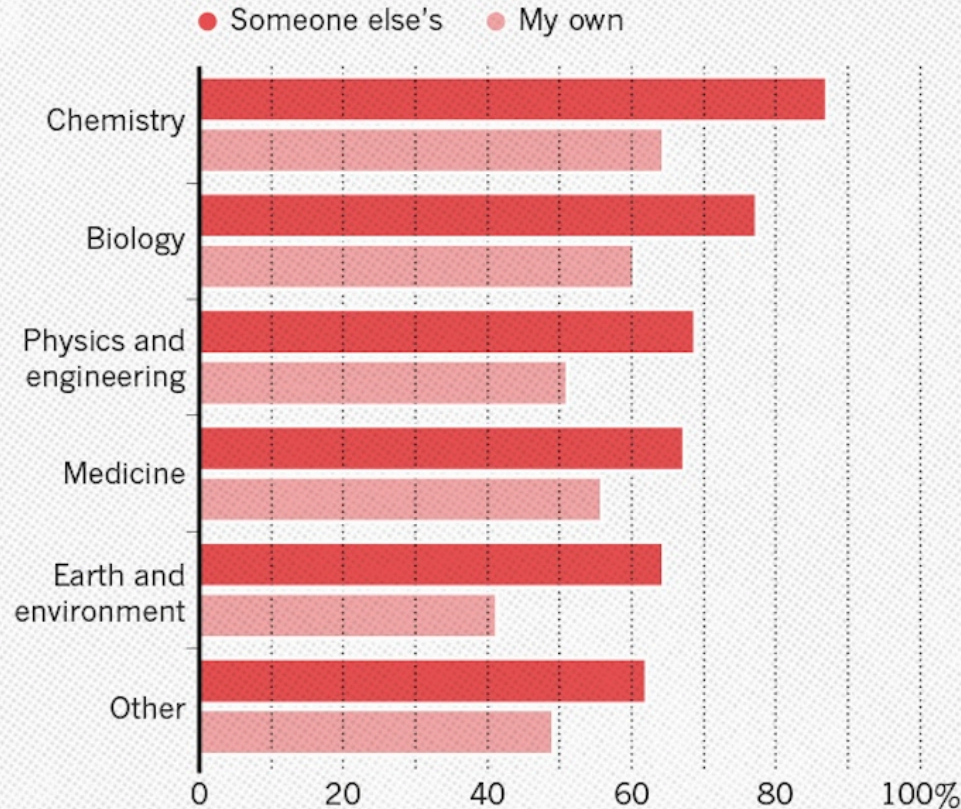
**Speaking of Science**

# No, science's reproducibility problem is not limited to psychology

Washington Post, Joel Achenbach, August 28, 2015

## HAVE YOU FAILED TO REPRODUCE AN EXPERIMENT?

Most scientists have experienced failure to reproduce results.



## **Nature's survey of 1,576 researchers:**

703 Biology

106 Chemistry

95 Earth and Environmental

203 Medicine

**236 Physics and Engineering**

233 Other

Nature, 2016, "1,500 scientists lift the lid on reproducibility", vol 533, Issue 734

# Digital Science: reproducibility and visibility in Astronomy

J.E. Ruiz<sup>1</sup>, L. Verdes-Montenegro<sup>1</sup>, S. Sánchez<sup>1</sup>, J.D. Santander-Vela<sup>1</sup>, and J. Garrido<sup>1</sup>

<sup>1</sup> Instituto de Astrofísica de Andalucía – CSIC

*Highlights of Spanish Astrophysics VII, Proceedings of the X Scientific Meeting of the Spanish Astronomical Society held on July 9 - 13, 2012, in Valencia, Spain. J. C. Guirado, L.M. Lara, V. Quilis, and J. Gorgas (eds.)*

**"In the Wf4Ever project we propose to improve the quality of science with metrics based on reproducibility and reuse, preserving decomposable thoroughly curated digital artefacts that enhances reproducibility and visibility of the experiment, as well as allowing more accurate mechanisms for credit attribution."**

**SHARING DATA, CODE, AND  
WORKFLOWS FACILITATES  
REPRODUCIBILITY AND  
REDUCES DUPLICATION**




**BUT DATA SHARING IS MORE  
THAN POSTING YOUR DATA IN  
YOUR WEBSITE**

 OPEN ACCESS  PEER-REVIEWED

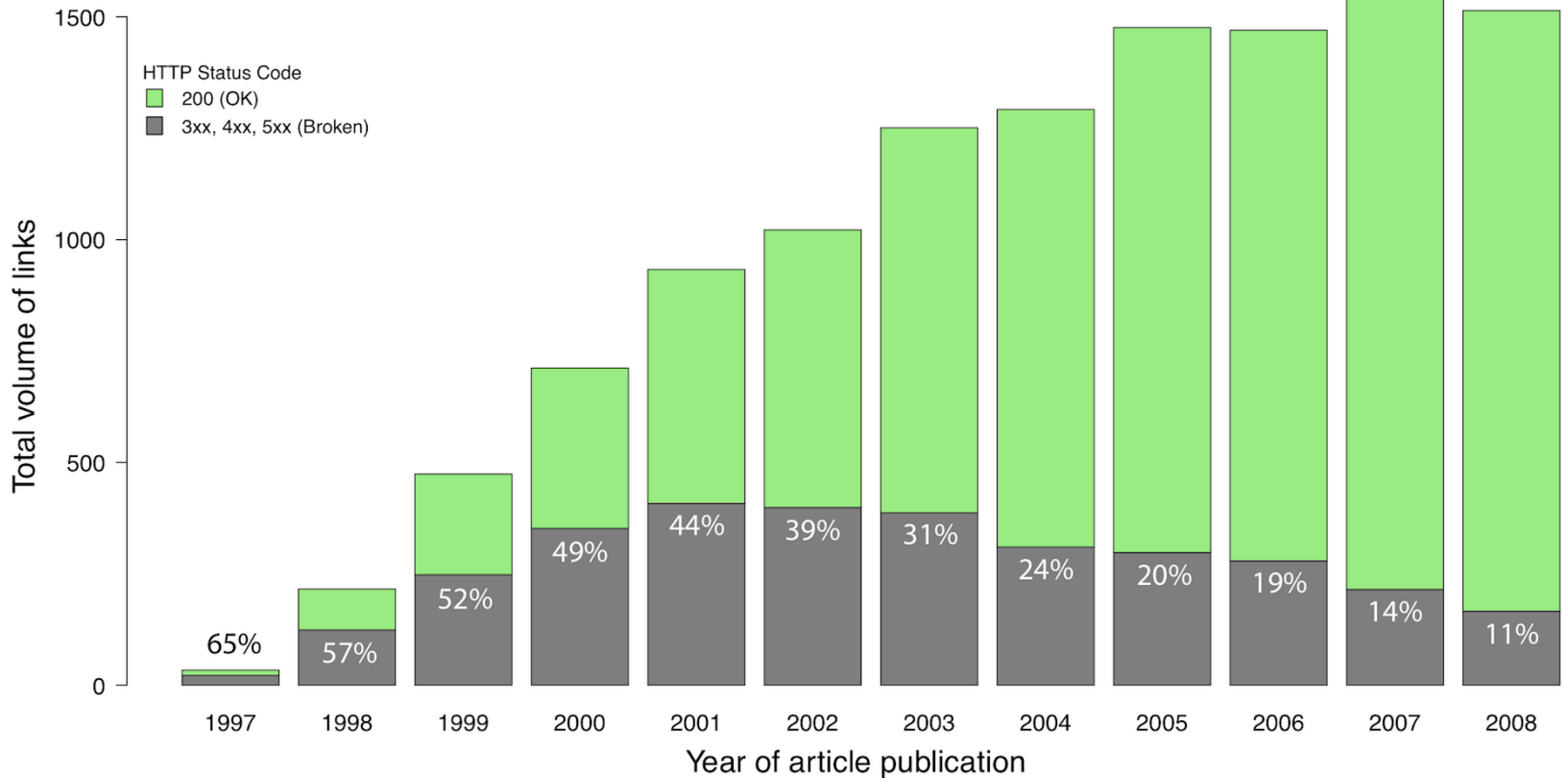
RESEARCH ARTICLE

# How Do Astronomers Share Data? Reliability and Persistence of Datasets Linked in AAS Publications and a Qualitative Study of Data Practices among US Astronomers

Alberto Pepe , Alyssa Goodman, August Muench, Merce Crosas, Christopher Erdmann

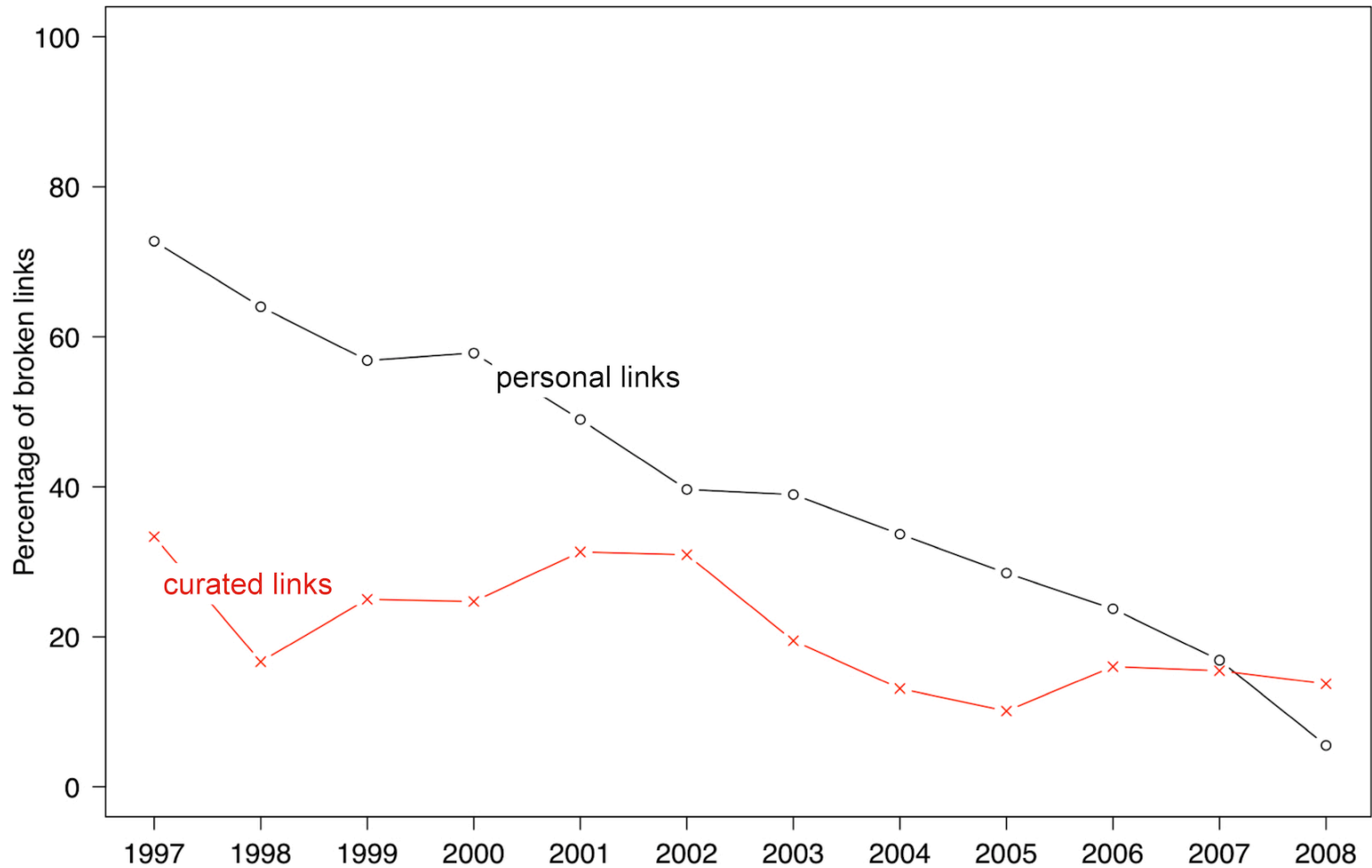
Published: August 28, 2014 • <https://doi.org/10.1371/journal.pone.0104798>

# MORE THAN HALF OF LINKS TO DATA IN ARTICLES FROM 15 YEARS AGO ARE BROKEN



External links in all articles published between 1997 and 2008 in the four main astronomy journals published by the American Astronomical Society.

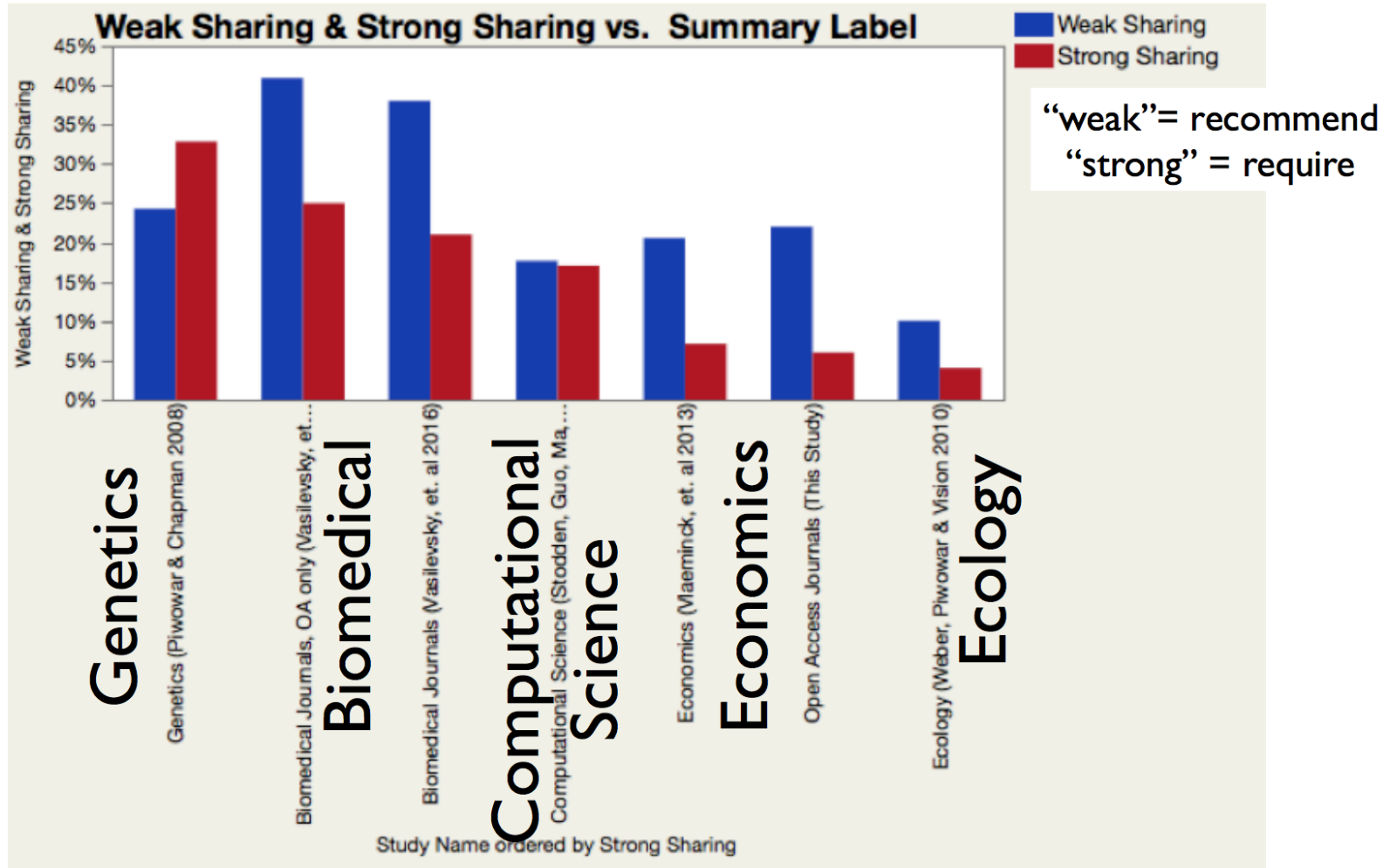
# 70% OF LINKS TO PERSONAL WEBSITES FROM ARTICLES PUBLISHED IN 1997 ARE BROKEN



# HOW CAN WE IMPROVE DATA SHARING?

- New Norms
- New Incentives
- New Technology

# FORMAL DATA-SHARING POLICIES ARE APPLIED IN JOURNALS ACROSS DISCIPLINES



Castro, Crosas, Garnett, Sheridan, Altman, 2017, Journal of Scholarly Publishing

# MANY FUNDERS REQUIRE DATA SHARING & OPEN DATA

## PRIVATE RESEARCH FUNDERS

- Bill and Melinda Gates Foundation Information Sharing Approach
- Sloan Foundation Data Sharing Policy
- Wellcome Trust Data Sharing Policy
- Arnold Foundation
- Moore Foundation
- Robert Wood Johnson Foundation
- HHMI Policy on the Sharing of Publication-Related Materials, Data and Software

## PUBLIC RESEARCH FUNDERS

- Department of Agriculture
- Department of Commerce
- Department of Defense
- Department of Education
- Department of Energy
- Department of Health and Human Services
  - Agency for Healthcare Research and Quality (AHRQ)
  - Assistant Secretary for Preparedness and Response (ASPR)
  - Center for Disease Control and Prevention (CDC)
  - Food and Drug Administration (FDA)
  - National Institutes of Health (NIH)
- Department of Homeland Security
- Department of Housing and Urban Development
- Department of Interior
- Department of Labor
- Department of Transportation
- Department of Veterans Affairs
- Environmental Protection Agency (EPA)

# Are reproducibility and open science starting to matter in tenure and promotion review?

July 14th, 2017, Brian Nosek

Tags: [openscience](#), [reproducibility](#)







# The NEW ENGLAND JOURNAL of MEDICINE

[HOME](#)[ARTICLES & MULTIMEDIA ▾](#)[ISSUES ▾](#)[SPECIALTIES & TOPICS ▾](#)[FOR AUTHORS ▾](#)[CME >](#)

## SOUNDING BOARD

### Data Authorship as an Incentive to Data Sharing

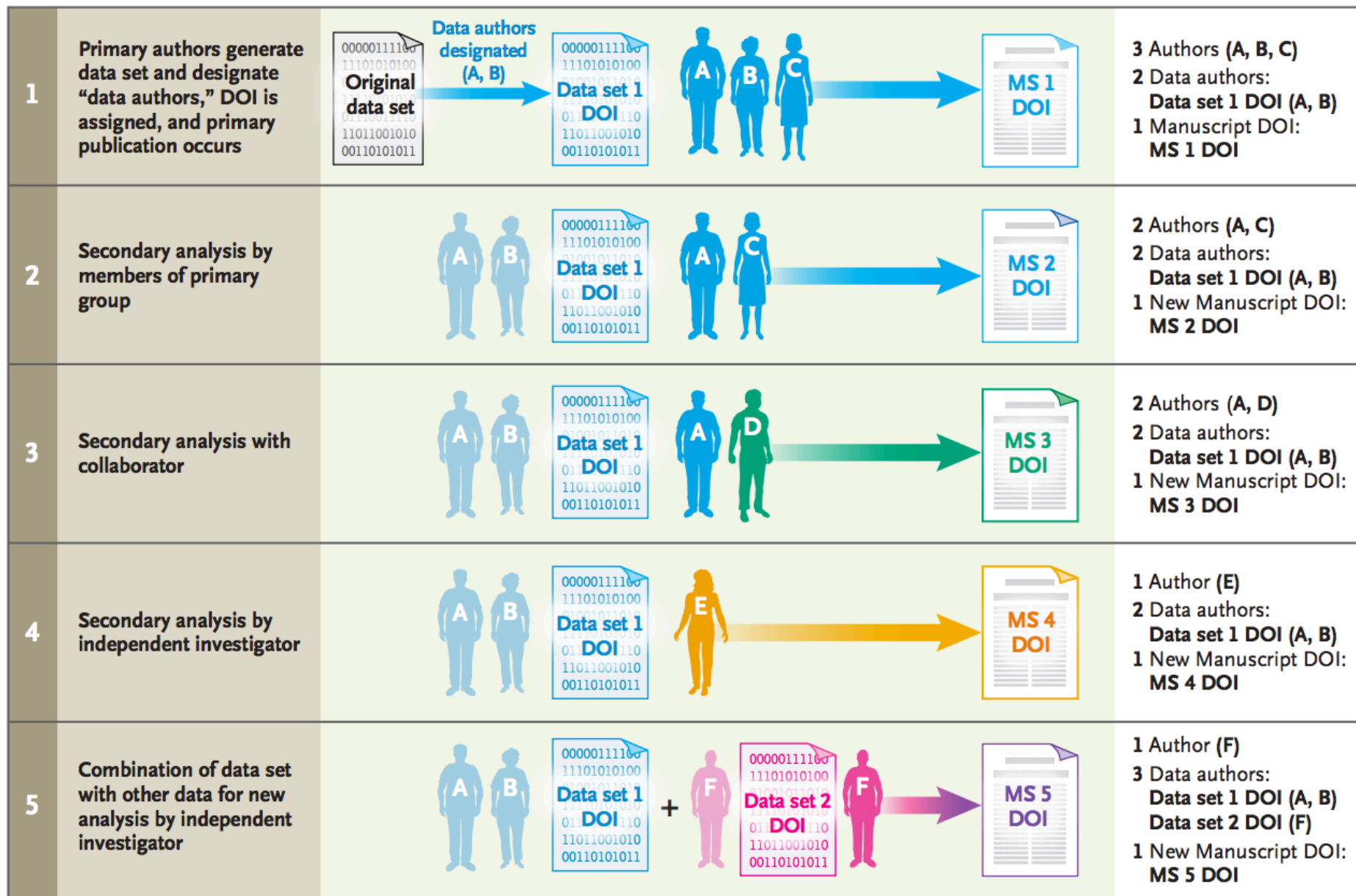
Barbara E. Bierer, M.D., Mercè Crosas, Ph.D., and Heather H. Pierce, J.D., M.P.H.

N Engl J Med 2017; 376:1684-1687 | [April 27, 2017](#) | DOI: 10.1056/NEJMSb1616595

Share: [f](#) [t](#) [g+](#) [in](#) [+](#)

[Article](#)[References](#)[Citing Articles \(3\)](#)[Letters](#)[Metrics](#)

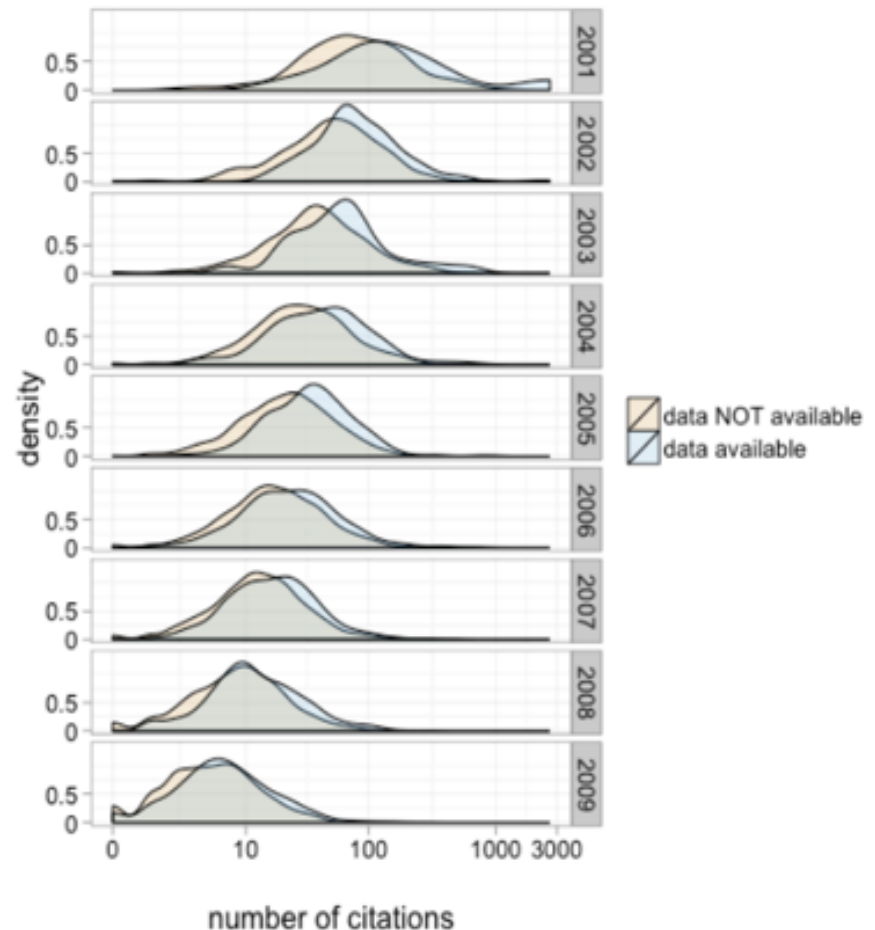
"We believe that both as a matter of fairness and as a matter of providing an incentive for data sharing, the persons who initially gathered the data should receive appropriate and standardized credit that can be used for academic advancement, for grant applications, and in broader situations."



# DATA SHARING INCREASES CITATIONS

From 10,555 studies with gene expression microarray data:

- Studies that shared data received 9% more citations
- Data reuse by other researchers continued for 6 years



Piowar and Vision (2013), Data reuse and the open data citation advantage. PeerJ 1:e175; DOI 10.7717/peerj.175





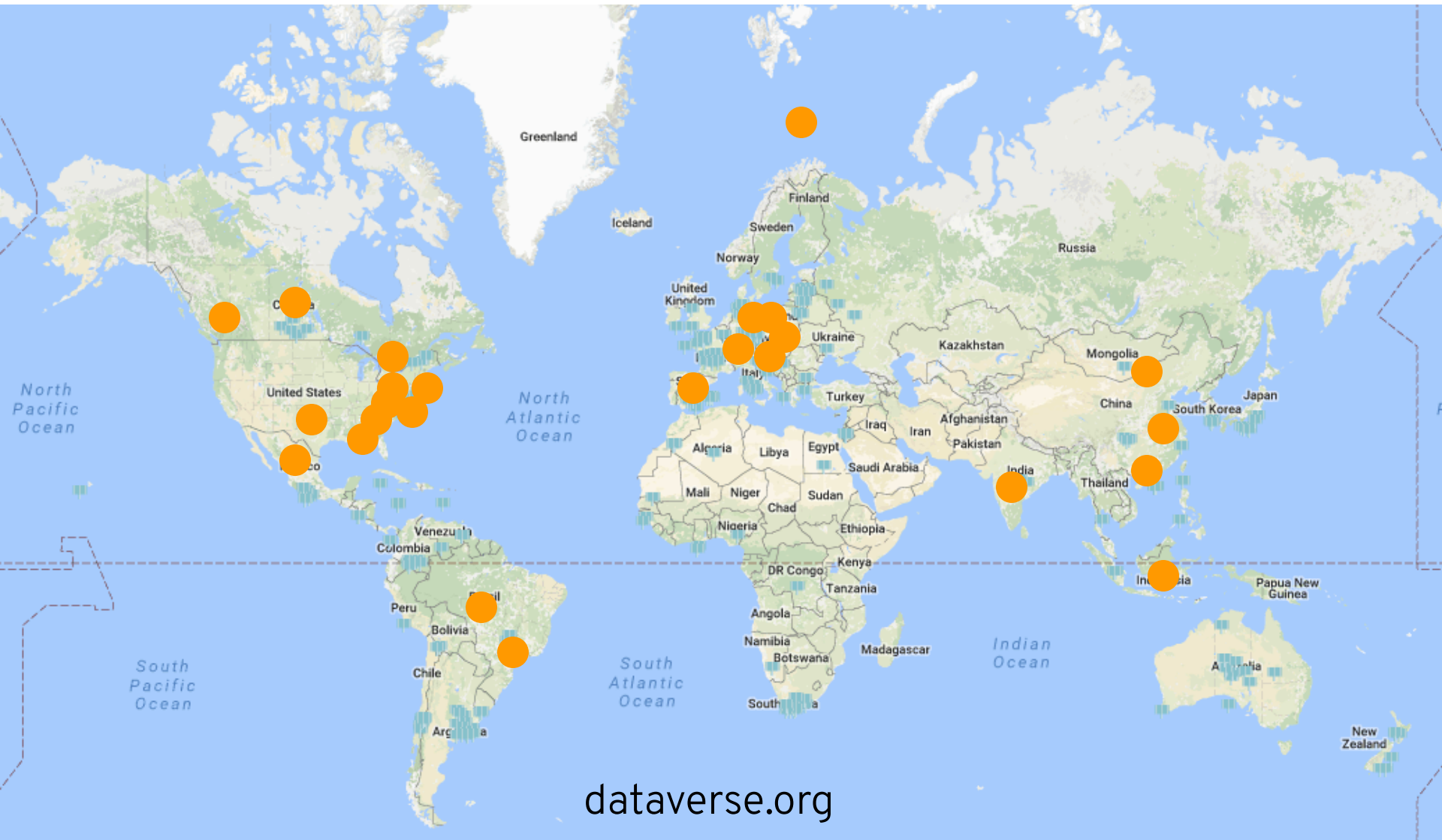
**OUR INSTITUTE PROVIDES A  
TECHNOLOGY SOLUTION TO  
DATA SHARING**



An open-source software to share, cite, and find data.  
Developed at Harvard's Institute for Quantitative Social Science

2006 (we started)

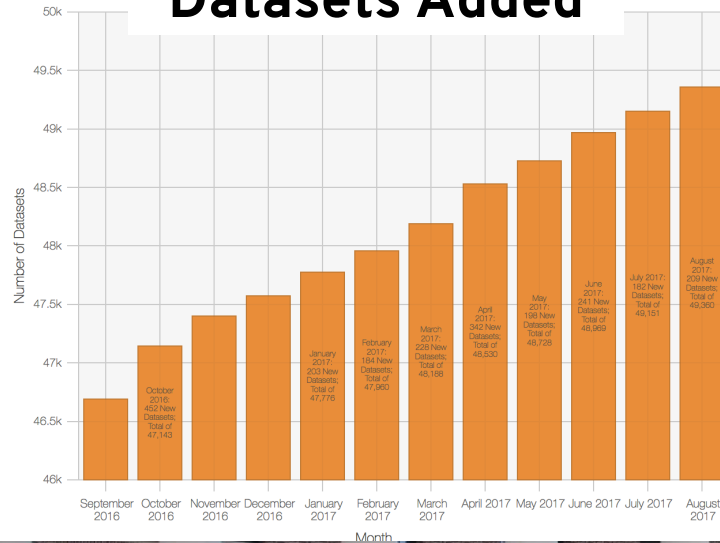
2017



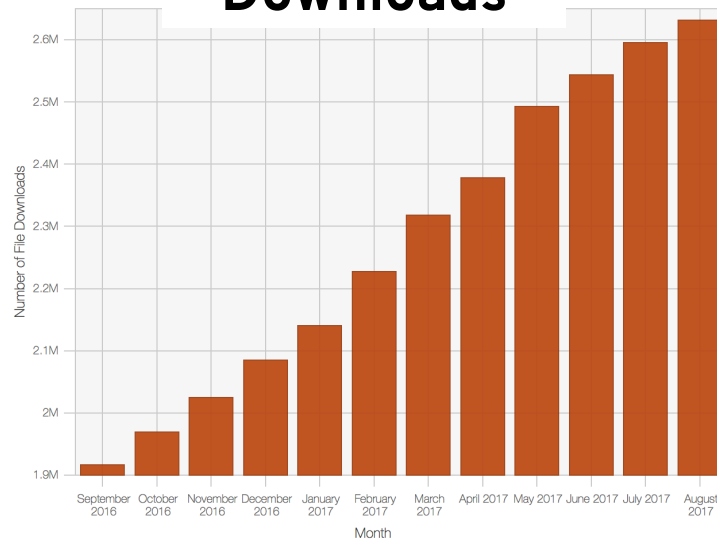


# HOW RESEARCHERS SHARE & USE DATA WITH DATAVERSE

## Datasets Added



## Downloads



## Harvard Dataverse Repository

- > 70,000 datasets total
- > 49,000 datasets uploaded to Harvard Dataverse repository
- 200 datasets/month

- > 340,000 files
- 4,000 files/month

- > 2.5 M downloads
- 60,000 downloads/month

[dataverse.harvard.edu](http://dataverse.harvard.edu)

# OUR CONTRIBUTIONS TO ENHANCE DATA SHARING

## **King, 1995, Replication, Replication**

## **2014, Joint Declaration of Data Citation Principles**

## **Wilkinson et al, 2016, The FAIR Guiding Principles for Scientific Data Management and Stewardship**

Altman et al, 2001, A Digital Library for the Dissemination and Replication of Quantitative Social Science

Pepe et al, 2014, How Do Astronomers Share Data?

Bierer, Crosas, Pierce, 2017, Data Authorship as an Incentive to Data Sharing

Altman and King, 2007, A Proposed Standard for the Scholarly Citation of Quantitative Data

Goodman et al, 2014, Ten Simple Rules for the Care and Feeding of Scientific Data

King, 2007, An Introduction to the Dataverse Network as an Infrastructure for Data Sharing

Crosas, Honaker, King, Sweeney, 2015, Automating Open Science for Big Data

Crosas, 2012, The Dataverse Network: an open source application for sharing, discovering, and preserving research data

Castro et al, 2015, Achieving Human and Machine Accessibility of Cited Data

Crosas, 2013, A Data Sharing Story

Sweeney, Crosas, Bar-Sinai, 2015, Sharing Sensitive Data with Confidence: The DataTags System

Altman and Crosas, 2013, The Evolution to Data Citation: from principles to implementation

Meyer et al. 2016, Data Publication with the Structural Biology Data Grid Supports Live Analysis



2017



Data should be ...

**FINDABLE**

**ACCESSIBLE**

**INTERPOPERABLE**

**REUSABLE**

Wilkinson et al., 2016, "The FAIR Guiding Principles for Scientific Data Management and Stewardship"

Nature Scientific Data

# FAIR DATA IN DATAVERSE

Data Citation  
with Persistent  
Identifier (DOI)

Data Files

Metadata

Data Licenses,  
User Agreements

Dataset Versions

The screenshot shows the Harvard Dataverse interface for the dataset "L1688 Dust Temperature and Opacity" (Version 1.0) by Goodman, Alyssa (2015). The interface includes a search bar, navigation tabs (Files, Metadata, Terms, Versions), and a list of files. Arrows from the left text labels point to specific features: "Data Citation with Persistent Identifier (DOI)" points to the DOI link; "Data Files" points to the "Files" tab; "Metadata" points to the "Metadata" tab; "Data Licenses, User Agreements" points to the "Terms" tab; and "Dataset Versions" points to the "Versions" tab.

**Dataverse** Search About User Guide Support Sign Up Log In

**L1688 Dust Temperature and Opacity** Version 1.0

Goodman, Alyssa, 2015, "L1688 Dust Temperature and Opacity", doi:10.7910/DVN/OWVFCE, Harvard Dataverse, V1 [Cite Dataset](#)

[Learn about Data Citation Standards.](#)

**Description** data reduction by Aaron Meisner and Hope Chen

**Subject** Astronomy and Astrophysics

**Files** Metadata Terms Versions

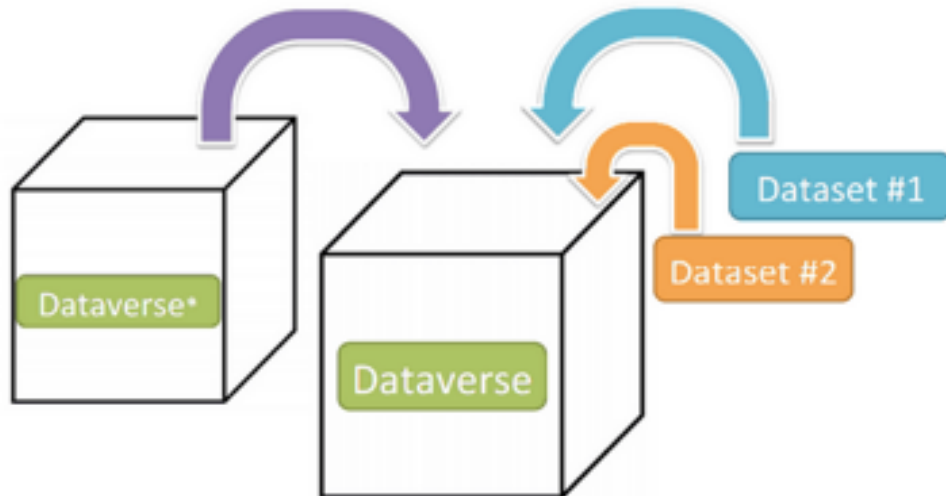
Search this dataset Find

**2 Files**

File Name	Size	Date	Downloads	MD5	Metadata	Download
<a href="#">tr_350um_L1688_meisner.fits</a>	FITS - 163.1 KB	Jul 14, 2015	2 Downloads	MD5: e6e94a6215c5ba1fa5e0074b4bb33056	This is a FITS file with 1 (primary) HDU. The following recognized metadata keys have been found in the FITS file: NAXIS0; NAXIS1; CRVAL2; NAXIS; CD1_1; CRVAL1;	<a href="#">Download</a>
<a href="#">T_L1688_meisner.fits</a>	FITS - 163.1 KB	Jul 14, 2015	4 Downloads	MD5: c1d3daba39d2f29517e3eb232bed413a	This is a FITS file with 1 (primary) HDU. The following recognized metadata keys have been found in the FITS file: NAXIS0; NAXIS1; CRVAL2; NAXIS; CD1_1; CRVAL1;	<a href="#">Download</a>

# A DATAVERSE IS A CONTAINER OF DATASETS AND A DATASET IS A CONTAINER OF DATA FILES, DOCUMENTATION, AND CODE

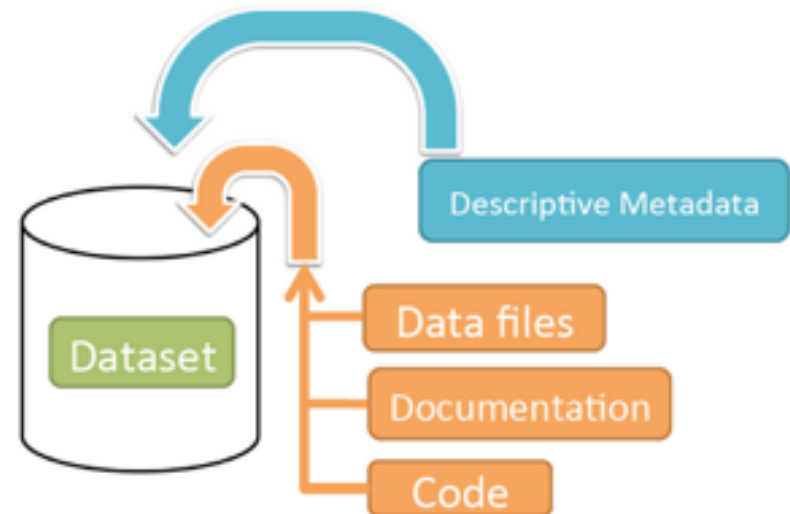
Schematic Diagram of a **Dataverse** in Dataverse 4.0



Container for your **Datasets** and/or **Dataverses**\*

\* Dataverses can now contain other Dataverses (this replaces Collections & Subnetworks)

Schematic Diagram of a **Dataset** in Dataverse 4.0



Container for your data, documentation, and code.

# DATAVERSE SUPPORTS ASTRONOMY DATA

The screenshot displays the Dataverse interface for a dataset titled "L1688 Dust Temperature and Opacity" (Version 1.0). The header includes the Dataverse logo and navigation links: About, User Guide, Support, Sign Up, and Log In. The dataset title is accompanied by a document icon and a "Version 1.0" badge. Below the title, the citation information is provided: "Goodman, Alyssa, 2015, 'L1688 Dust Temperature and Opacity', doi:10.7910/DVN/OWVFCE, Harvard Dataverse, V1". A "Cite Dataset" button and a link to "Learn about Data Citation Standards" are also present. The "Description" field states "data reduction by Aaron Meisner and Hope Chen", and the "Subject" field is "Astronomy and Astrophysics". Navigation tabs for "Files", "Metadata", "Terms", and "Versions" are shown. A search bar with the placeholder "Search this dataset..." and a "Find" button is located below the tabs. The "2 Files" section lists two FITS files: "tau350um\_L1688\_meisner.fits" and "T\_L1688\_meisner.fits". Each file entry includes its size (163.1 KB), upload date (Jul 14, 2015), download count, MD5 hash, a description of the file's structure (1 primary HDU), recognized metadata keys (NAXIS0, NAXIS1, CRVAL2, NAXIS, CD1\_1, CRVAL1), and a "Data" label. Each file has a corresponding "Download" button.

**L1688 Dust Temperature and Opacity** Version 1.0

Goodman, Alyssa, 2015, "L1688 Dust Temperature and Opacity", doi:10.7910/DVN/OWVFCE, Harvard Dataverse, V1

**Description** data reduction by Aaron Meisner and Hope Chen

**Subject** Astronomy and Astrophysics

Files Metadata Terms Versions

Search this dataset... Find


2 Files


**tau350um\_L1688\_meisner.fits**  
FITS - 163.1 KB - Jul 14, 2015 - 2 Downloads  
MD5: e6e94a6215c5ba1fa5e0074b4bb33056  
This is a FITS file with 1 (primary) HDU. The following recognized metadata keys have been found in the FITS file: NAXIS0; NAXIS1; CRVAL2; NAXIS; CD1\_1; CRVAL1;  
**Data**


**T\_L1688\_meisner.fits**  
FITS - 163.1 KB - Jul 14, 2015 - 4 Downloads  
MD5: c1d3daba39d2f29517e3eb232bed413a  
This is a FITS file with 1 (primary) HDU. The following recognized metadata keys have been found in the FITS file: NAXIS0; NAXIS1; CRVAL2; NAXIS; CD1\_1; CRVAL1;  
**Data**

- Supports default astronomy metadata fields (based on virtual observatory schema)
- Extracts header metadata from FITS files upon ingest

# DATAVERSE USED BY MAX-PLANCK INSTITUTE ...

 **Dataverse**

 [About](#) [User Guide](#) [Support](#) [Sign Up](#) [Log In](#)


 **J. Michael's Analysis Projects Dataverse** (Max-Planck Institute for Extraterrestrial Physics)


[Harvard Dataverse](#) > **J. Michael's Analysis Projects Dataverse**


[Contact](#) [Share](#)

Data products, analysis results, and code associated with my publications to be used for replication.

[Advanced Search](#)

☒  **Dataverses (0)**


☒  **Datasets (1)**

☐  **Files (133)**

**Publication Date**  
2017 (1)

**Subject**  
[Astronomy and Astrophysics \(1\)](#)

1 to 1 of 1 Result

 **Replication Data for: Is the Spectral Width of GRBs a Reliable Measure of GRB Emission Physics?**  
May 15, 2017  
Burgess, J. Michael, 2017, "Replication Data for: Is the Spectral Width of GRBs a Reliable Measure of GRB Emission Physics?", doi:10.7910/DVN/BDC2GS, Harvard Dataverse, V1

This data set includes the processed GBM PHA/BAK/RSP files used in the publication. The files are readable by XSPEC or 3ML's OGIPLike plugin. Additionally, the 3ML analysis results FITS files are included. These results include the Bayesian posteriors from the fits. These can be...

# DATAVERSE IN THE ASTRONOMY NEWS ...



WEIRD NEWS 08/31/2017 05:15 am ET

## Scientists Detect Mysterious Radio Signals From Deep Space

Strange bursts are coming from 3 billion light years away.



By Ed Mazza



More than 12,000  
downloads!

The screenshot shows the Harvard Dataverse interface for a dataset titled 'The Sound of Fast Radio Burst FRB 121102'. An orange arrow points from the text 'More than 12,000 downloads!' to the '12,520 Downloads' metric. The page includes a description of the dataset, its subject (Astronomy and Astrophysics), keywords (radio interferometry, astronomy, fast radio burst), and a list of 9 files. The first two files are audio files in mp3 format, each with a download button.

**Harvard Dataverse** 12,520 Downloads

**The Sound of Fast Radio Burst FRB 121102** Version 1.0

Law, Casey, 2016, "The Sound of Fast Radio Burst FRB 121102", doi:10.7910/DVN/QSWJE6, Harvard Dataverse, V1

**Description** Sound files generated from nine radio transients detected by the Very Large Array toward cosmological radio transient, FRB 121102. The sound of FRB 121102 has a "chirp" that is caused by dispersion. Original publication is Chatterjee et al (2017), "The direct localization of a fast radio burst and its host". Original data is available at https://doi.org/10.7910/DVN/TLDKXG.

**Subject** Astronomy and Astrophysics

**Keyword** radio interferometry, astronomy, fast radio burst

**Notes** These bursts were observed with the VLA from 2.5 to 3.5 GHz with a dispersion measure of approximately 570 pc/cm<sup>3</sup>.

**Files** **Metadata** **Terms** **Versions**

Search this dataset... Find

**9 Files**

File Name	Format	Size	Date	Downloads	MD5	Download
57623.mp3	audio/mp3	200.0 KB	Dec 31, 2016	6,256 Downloads	9967df426bfaa8d239110dbd3c9ffbc0	Download
57633.scan13.mp3	audio/mp3	200.0 KB	Dec 31, 2016	1,930 Downloads	0c864a33799d4f98a4e0f018f7ad59c3	Download

**WHAT ARE WE WORKING ON NOW?**

# **DATA PROVENANCE**

**TRACK THE ORIGINAL SOURCE OF A DATASET**



# SCIENTIFIC DATA

OPEN

## Comment: If these data could talk

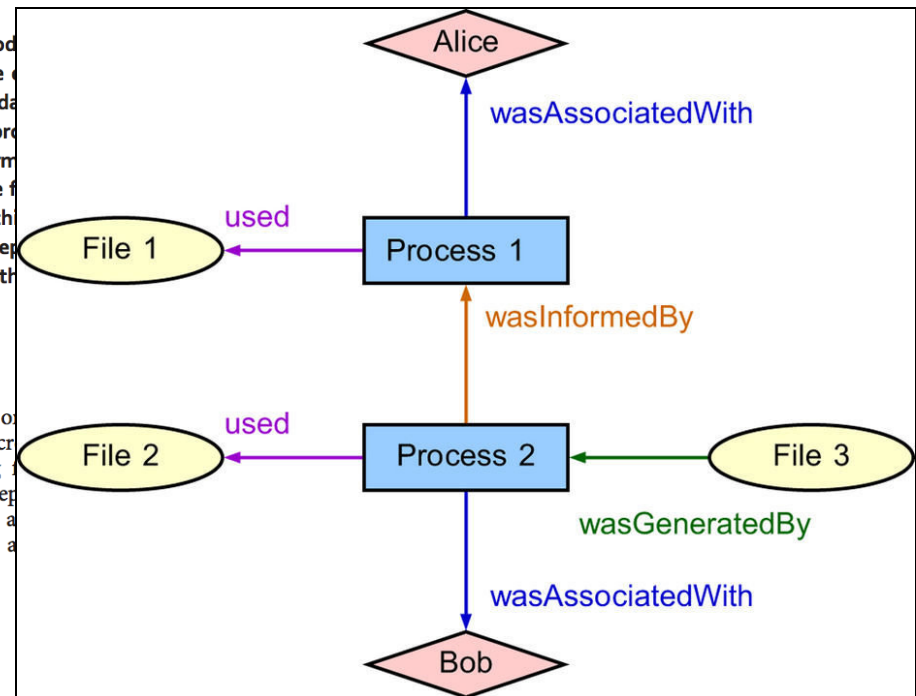
Thomas Pasquier<sup>1</sup>, Matthew K. Lau<sup>2</sup>, Ana Trisovic<sup>3,4</sup>, Emery R. Boose<sup>2</sup>, Ben Couturier<sup>3</sup>,  
Mercè Crosas<sup>5</sup>, Aaron M. Ellison<sup>2</sup>, Valerie Gibson<sup>4</sup>, Chris R. Jones<sup>4</sup> & Margo Seltzer<sup>1</sup>

Received: 12 April 2017  
Accepted: 24 July 2017  
Published: 5 September 2017

In the last few decades, data-driven methods have become essential for managing and analyzing the growing flood of data. However, in many fields exhibit distressingly low rates of reproducibility. In this issue, we believe that there is a lack of formal records from the data source to the analysis to the final publication. To make their research and data accessible, they need to report through *systematic* and *formal* records of their publications and researchers.

### Reproducibility

The success and power of science depends on the reproducibility of its results. Issues with reproducibility have surfaced across many fields, including medicine<sup>1</sup>. Although the lack of reproducibility remains a worrisome issue. This comes at a time when data is exponentially<sup>3</sup>. At the same time, the data is becoming computationally demanding.



# **CLOUD DATAVERSE**

**COMBINE DATA REPOSITORIES WITH CLOUD  
COMPUTING**

**Users, External Tools, Services**



Deposit

Access

Compute

**Software: Services & Tools**



**Data Storage**



**Cloud Computing**



+



**FAIR Cloud Dataaverse**

# **DATA PRIVACY**

**CLASSIFY AND HANDLE DATASETS BASED ON  
THEIR PRIVACY LEVEL**

## Dataverse® as a DataTags repository

### Data file deposit

Assistance to assign DataTag from:

- DataTags automated interview
- RobotLawyer auto-generated data user agreements (DUA)
- Review Board



**orange**

### Direct Access

Requires:

- User registration
- Approval needed for access
- Signed DUA

**green**

### Privacy Preserving Access

- Requires user registration
- Provides access to differentially private statistics using Private data Sharing Interface (PSI)

Harvard Data Privacy Tools Project: [privacytools.seas.harvard.edu](https://privacytools.seas.harvard.edu)

DataTags Project: [datatags.org](https://datatags.org)

# **INTEGRATION WITH TOOLS**

**DATVERSE AS PART OF THE DATA LIFECYCLE**

## Data Collection

Lab
E-Notebooks
Instruments
Surveys
...

Track Provenance

Assign DUA  
&  
metadata

Run data &  
code

Explore &  
Visualize data

Journals &  
Funders

Data  
Citation

Work with  
Sensitive Data

Cloud Computing and  
Storage

The  
**Dataverse**<sup>®</sup>  
Project







# DATAVERSE COMMUNITY



**49 SOFTWARE  
CONTRIBUTORS**

# BI-WEEKLY COMMUNITY CALLS

235 ATTENDEES

26 ORGANIZATIONS/UNIVERSITIES

11 COUNTRIES

# **ANNUAL COMMUNITY MEETING**

**NEXT: JUNE 13, 14, 15, 2018**

# THANKS

@mercecrosas

[scholar.harvard.edu/mercecrosas](https://scholar.harvard.edu/mercecrosas)

[dataverse.org](https://dataverse.org)

