



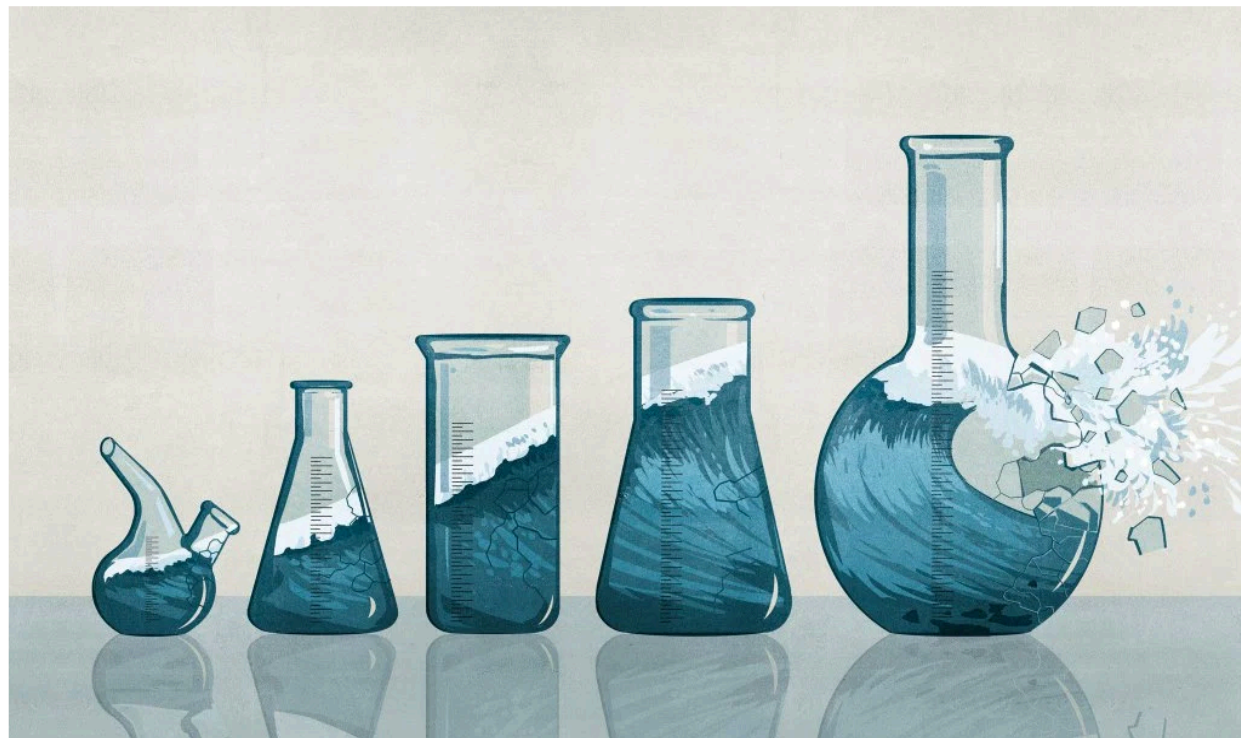
A data repository system for sharing and archiving research data

*A Solution for Publishing FAIR research data:  
Findable, Accessible, Interoperable, Reusable*

 [@dataverseorg](https://twitter.com/dataverseorg) [@mercecrosas](https://twitter.com/mercecrosas)

Science isn't broken  
 "It's just a hell a lot  
 harder than we give it  
 credit for"

But it's self-correcting



THE SCIENTIFIC METHOD | 7:00 AM | AUG 19, 2015

## Science Isn't Broken

It's just a hell of a lot harder than we give it credit for.

By CHRISTIE ASCHWANDEN  
 Graphics by RITCHIE KING

If you follow the headlines, your confidence in science may have taken a hit lately. Peer review? More like self-review. An investigation in November uncovered a scam in which [researchers were rubber-stamping their own work](#), circumventing peer review at five high-profile

# INSIGHTS

Design principles for  
 synthetic ecology p. 1425 ▶

Whacking hydrogen  
 into metal p. 1429



### SCIENTIFIC INTEGRITY

## Self-correction in science at work

Improve incentives to support research integrity

By Bruce Alberts,<sup>1</sup> Ralph J. Cicerone,<sup>2</sup> Stephen E. Fienberg,<sup>3</sup> Alexander Kamb,<sup>4</sup> Marcia McNutt,<sup>5\*</sup> Robert M. Nerem,<sup>6</sup> Randy Schekman,<sup>7</sup> Richard Shiffrin,<sup>8</sup> Victoria Stodden,<sup>9</sup> Subra Suresh,<sup>10</sup> Maria T. Zuber,<sup>11</sup> Barbara Kline Pope,<sup>12</sup> Kathleen Hall Jamieson<sup>13,14</sup>

**W**EEK after week, news outlets carry word of new scientific discoveries, but the media sometimes give suspect science equal play with substantive discoveries. Careful qualifications about what is known are lost in categorical headlines. Rare instances of misconduct or instances of irreproducibility are translated into concerns that science is broken. The Octo-

ber 2013 *Economist* headline proclaimed "Trouble at the lab: Scientists like to think of science as self-correcting. To an alarming degree, it is not" (1). Yet, that article is also rich with instances of science both policing itself, which is how the problems came to *The Economist's* attention in the first place, and addressing discovered lapses and irreproducibility concerns. In light of such issues and efforts, the U.S. National Academy of Sciences (NAS) and the Annenberg Retreat at Sunnylands convened our group to examine ways to remove some of the current disincentives to high standards of integrity in science.

Like all human endeavors, science is imperfect. However, as Robert Merton noted more than half a century ago "the

activities of scientists are subject to rigorous policing, to a degree perhaps unparalleled in any other field of activity" (2). As a result, as Popper argued, "science is one of the very few human activities—perhaps the only one—in which errors are systematically criticized and fairly often, in time, corrected" (3). Instances in which scientists detect and address flaws in work constitute evidence of success, not failure, because they demonstrate the underlying protective mechanisms of science at work.

Still, as in any human venture, science writ large does not always live up to its ideals. Although attempts to replicate the 1998 Wakefield study alleging an association between autism and the MMR (measles,

# Self-correction in science?

“trust, but verify”

Data (and code) should be  
shared for other  
researchers to inspect, test  
and reuse



The Commons supports  
biomedical discovery by  
enabling  
object

## Findable Accessible Interoperable and Reusable (FAIR)

### Digital Object Compliance for the *Commons*

Data sharing is a key objective of the *Commons*, and cloud computing provides the tools, as well as to make data (including the processed data that is generated) accessible. Access control can easily be implemented in the cloud so that data and tools can be appropriately and securely shared amongst groups authorized to use them, including the appropriate protections and access for human subjects data. Thus, while the cloud provides a computing environment to share data and tools in order to be able to effectively use these digital objects, they must have attributes that make them **Findable, Accessible, Interoperable and Reusable (FAIR)**. The *Commons* is intended to be a system that will do so.

A set of Digital Object Compliance principles that supports FAIR is currently under development. The Digital Object Compliance principles are expected to evolve over time as the ability to make digital objects meet the FAIR criteria increases, thereby improving the ability of digital objects to be shared and used more easily and effectively in the *Commons*.

To meet the most basic level of compliance, it is expected that digital objects would have the following elements:

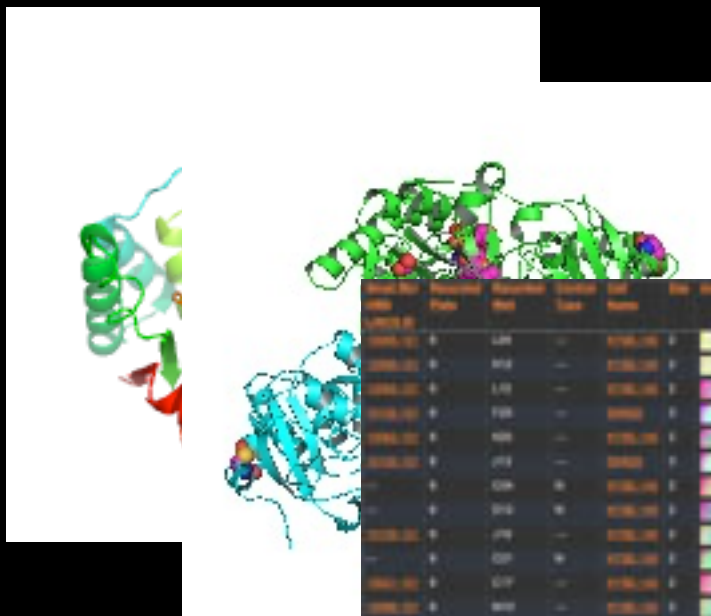
- Unique digital object identifiers
- A minimal set of searchable metadata
- Physical availability through a cloud-based *Commons* provider
- Clear access rules and controls (especially important for human subjects data)
- An entry (with metadata) in one or more indices

[Home](#) ▶ [Comment](#) ▶ [Data Descriptor](#)SCIENTIFIC DATA | 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](#)





Accession	Release Date	Resolution	Method	Model	Seq. Length	Seq. ID	Gene	Organism	Systematic Name	Source	Strain	Substrate	Temperature	pH	Time	Author	Journal	Year	Volume	Page	PMID	DOI
1G88	2001	2.5	X-ray	1	300	1	1	Homo sapiens	Proteinase 3	Human			25	7.5	10	Wang et al.	Nature	2001	408	371	11557111	10.1038/35425a
1G89	2001	2.5	X-ray	1	300	1	1	Homo sapiens	Proteinase 3	Human			25	7.5	10	Wang et al.	Nature	2001	408	371	11557111	10.1038/35425a
1G90	2001	2.5	X-ray	1	300	1	1	Homo sapiens	Proteinase 3	Human			25	7.5	10	Wang et al.	Nature	2001	408	371	11557111	10.1038/35425a
1G91	2001	2.5	X-ray	1	300	1	1	Homo sapiens	Proteinase 3	Human			25	7.5	10	Wang et al.	Nature	2001	408	371	11557111	10.1038/35425a
1G92	2001	2.5	X-ray	1	300	1	1	Homo sapiens	Proteinase 3	Human			25	7.5	10	Wang et al.	Nature	2001	408	371	11557111	10.1038/35425a
1G93	2001	2.5	X-ray	1	300	1	1	Homo sapiens	Proteinase 3	Human			25	7.5	10	Wang et al.	Nature	2001	408	371	11557111	10.1038/35425a
1G94	2001	2.5	X-ray	1	300	1	1	Homo sapiens	Proteinase 3	Human			25	7.5	10	Wang et al.	Nature	2001	408	371	11557111	10.1038/35425a
1G95	2001	2.5	X-ray	1	300	1	1	Homo sapiens	Proteinase 3	Human			25	7.5	10	Wang et al.	Nature	2001	408	371	11557111	10.1038/35425a
1G96	2001	2.5	X-ray	1	300	1	1	Homo sapiens	Proteinase 3	Human			25	7.5	10	Wang et al.	Nature	2001	408	371	11557111	10.1038/35425a
1G97	2001	2.5	X-ray	1	300	1	1	Homo sapiens	Proteinase 3	Human			25	7.5	10	Wang et al.	Nature	2001	408	371	11557111	10.1038/35425a
1G98	2001	2.5	X-ray	1	300	1	1	Homo sapiens	Proteinase 3	Human			25	7.5	10	Wang et al.	Nature	2001	408	371	11557111	10.1038/35425a
1G99	2001	2.5	X-ray	1	300	1	1	Homo sapiens	Proteinase 3	Human			25	7.5	10	Wang et al.	Nature	2001	408	371	11557111	10.1038/35425a
1G00	2001	2.5	X-ray	1	300	1	1	Homo sapiens	Proteinase 3	Human			25	7.5	10	Wang et al.	Nature	2001	408	371	11557111	10.1038/35425a

Research Data Lifecycle  
(collect, process, analyze, compute)

Publish and Archive Data

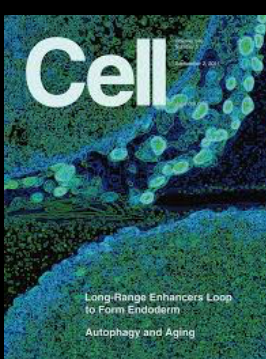


Data Citation

Metadata

Access Controls

Publish Research Results



17 Installations

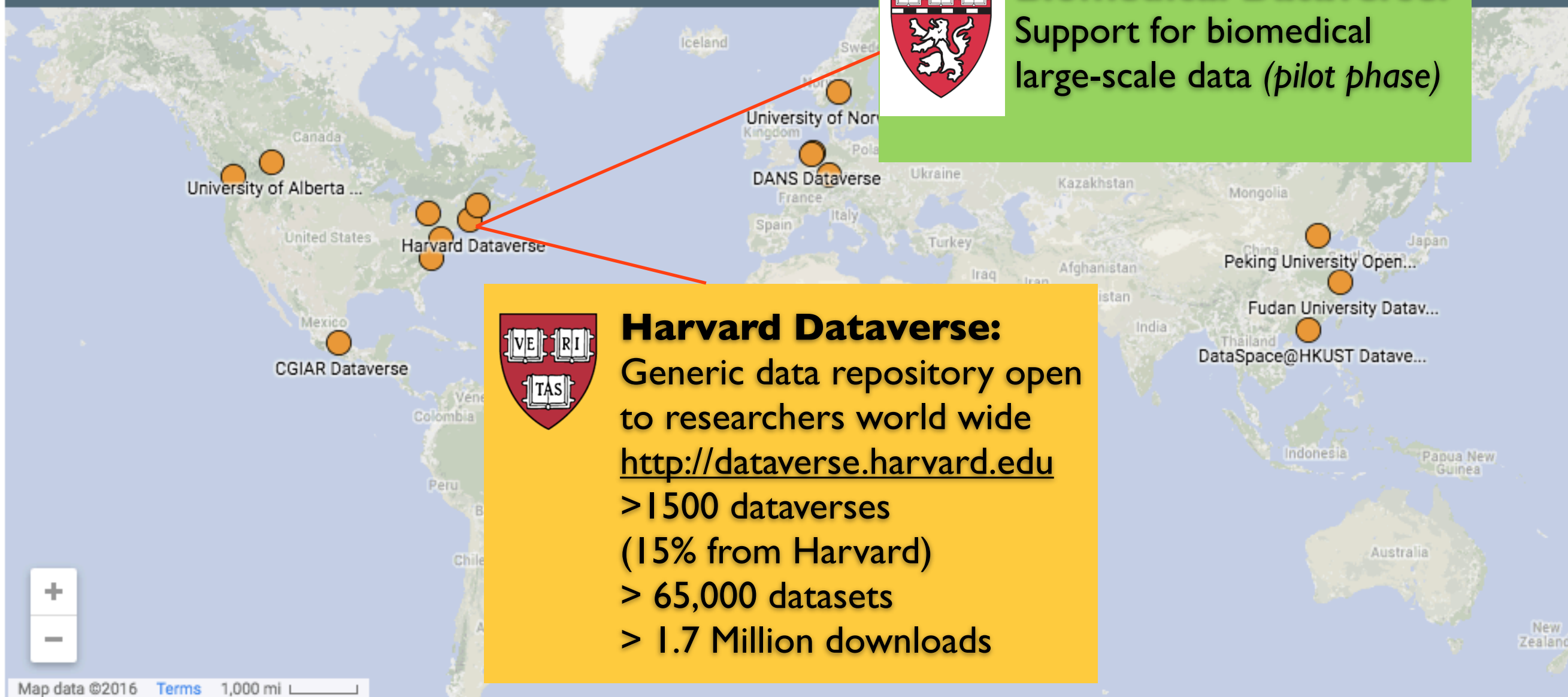
1,500+ Dataverses

65,000+ Datasets

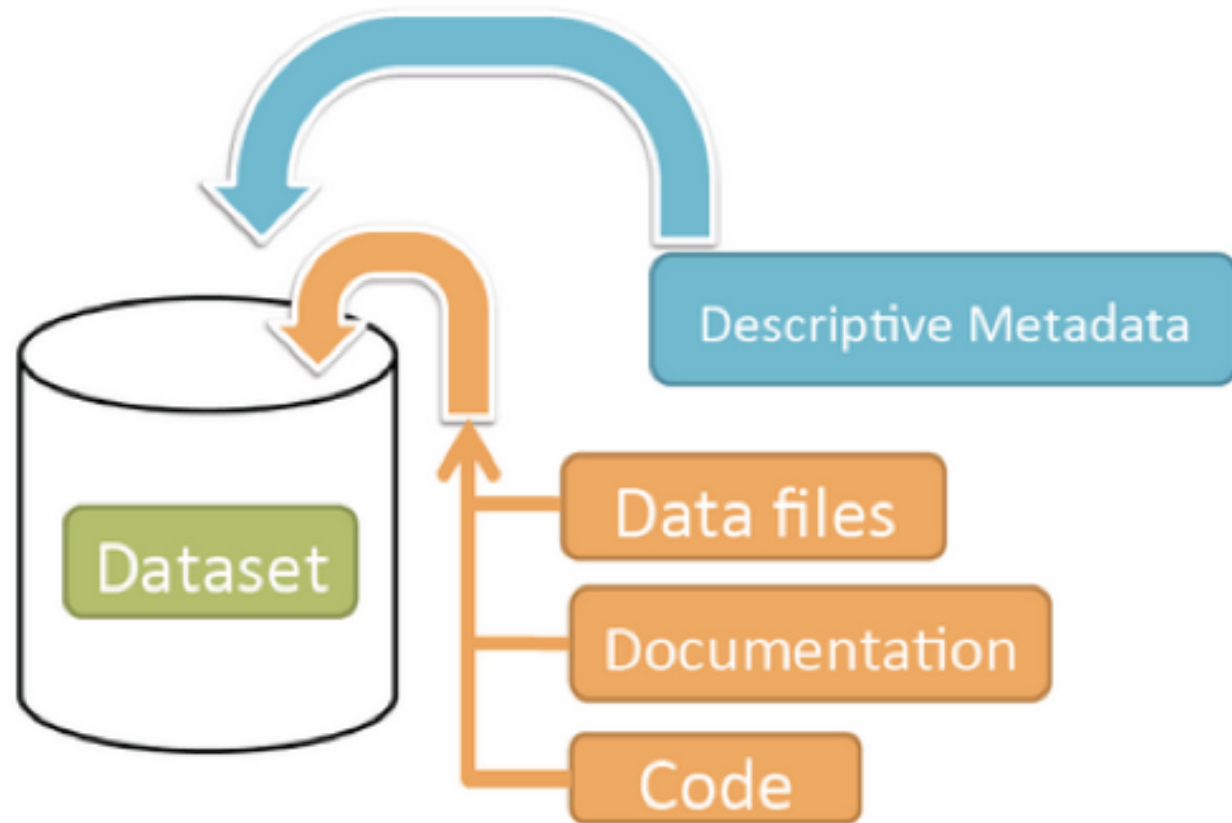
1,700,000+ Downloads



Dataverse Repositories ☆

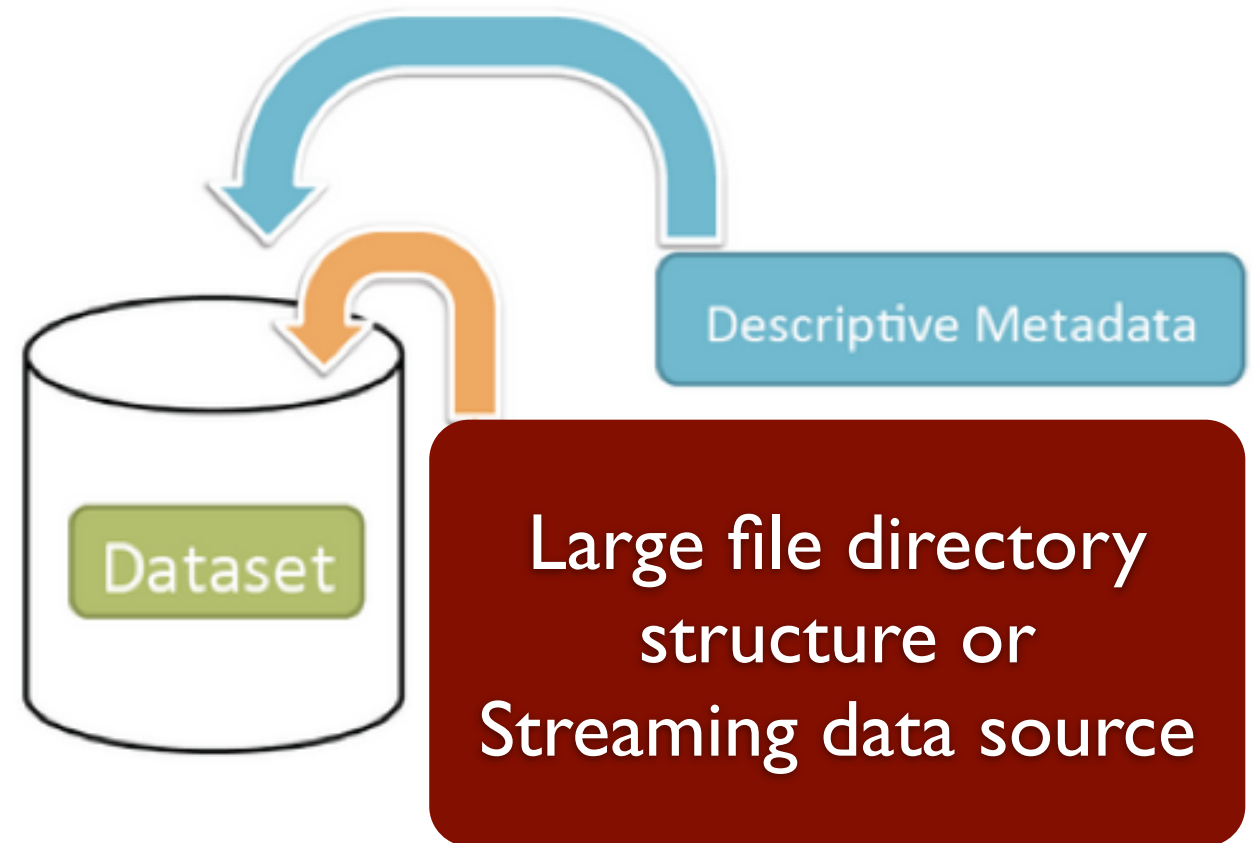


# Dataverse Now



Up to 100 files per dataset  
Up to ~ Gb per file

# Dataverse Big



> 1000s files per dataset  
> Gb per file



# Dataverse Projects supporting Big Data

## Biomedical Big Data

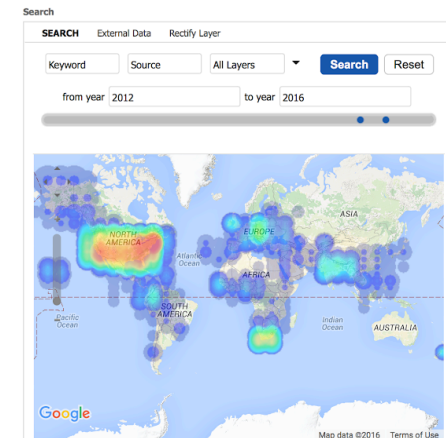
- 1000s image files per data set
- Non-http batch upload
- Data files in directory structure



Library of Integrated  
Network-based  
Cellular Signatures



## Social Science Big Data



A Billion  
streaming  
Geotweets



Alfred P. Sloan  
FOUNDATION

## Cloud Dataverse



The  
**Dataverse**  
Project

+ Swift Storage  
(scalable, access to  
computing)

## Integration with Data Management Systems

STARFISH +

The  
**Dataverse**  
Project  
at HMS

iRODS +

The  
**Dataverse**  
Project  
at UNC