

DATA SHARING FOR BETTER SCIENCE AND BETTER HEALTH: THE DATAVERSE PROJECT

Mercè Crosas, Institute for Quantitative Social Science, Harvard University

 @mercecrosas

XXXVII JORNADAS DE ECONOMÍA DE LA SALUD, BARCELONA, 6-8 SEPTIEMBRE

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

Regulating the internet giants

The world's most valuable resource is no longer oil, but data

The data economy demands a new approach to antitrust rules

 [Print edition](#) | [Leaders >](#)

May 6th 2017



DATA SCIENCE, BIG DATA

"Every two years, the amount of digitized data is equal to all of the data ever collected before. The world's knowledge is at our fingertips, and **data science** allows us to effectively and efficiently make use of that knowledge. This is facilitating a societal shift as big as the Industrial Revolution. "

UVAToday Q&A, August 21, 2017



Phil Bourne

Data Science Director, UVA

Former Associate Director for

Data Science, NIH

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 BORGHI
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

A microscopic image showing a cluster of cells. The cells are stained with a blue dye, likely DAPI, which highlights the nuclei. The overall structure is irregular and dense, with many small, round, red-stained structures interspersed among the blue-stained cells. The background is dark blue.

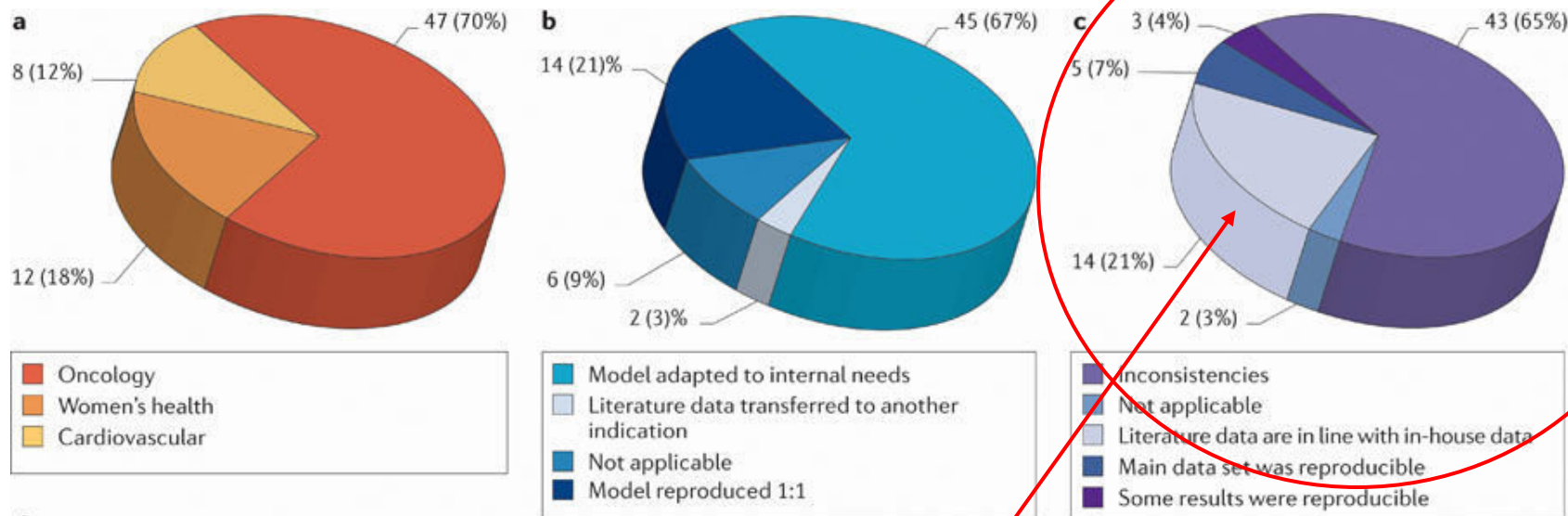
REPRODUCIBILITY IN CANCER STUDIES

**ONLY 6 (11%) OUT OF 53
LANDMARK STUDIES
THAT CLAIM TO TREAT
CANCER COULD BE
REPRODUCED**

Begley & Ellis, Nature, 2012

(from Scientists at Amgen biotechnology)

STUDY FROM BAYER: ONLY 14(21%) OF 67 PUBLICATIONS COULD BE REPRODUCED



d

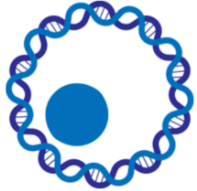
	Model reproduced 1:1	Model adapted to internal needs (cell line, assays)	Literature data transferred to another indication	Not applicable
In-house data in line with published results	1 (7%)	12 (86%)	0	1 (7%)
Inconsistencies that led to project termination	11 (26%)	26 (60%)	2 (5%)	4 (9%)

21% of Literature Data are in line with in-house data

Nature Reviews | Drug Discovery

Prinz, Schlange, Asadullah, 2011, Nature Reviews Drug Discovery

Reproducibility Project: Cancer Biology (RP:CB) Overview



The Reproducibility Project: Cancer Biology (RP:CB) is an initiative to conduct direct replications of 50 high-impact cancer biology studies. The project anticipates learning more about predictors of reproducibility, common obstacles to conducting replications, and how the current scientific incentive structure affects research practices by estimating the rate of reproducibility in a sample of published cancer biology literature. The RP:CB is a collaborative effort between the Center for Open Science and the Center for Cancer Research at the University of Michigan.

Through independent direct replication studies, the project aims to identify best practices that maximize reproducibility and to provide a more accurate picture of the reliability of published findings and the factors that influence them.

Additionally, we expect to learn about:

- The overall rate of reproducibility in a sample of the published cancer biology literature.
- Obstacles that arise in conducting direct replication of original studies.
- Predictors of replication success such as the journal in which the original finding was published, the citation impact of the original report, the number of direct replications that have been published elsewhere, the transparency of materials and methods included with the publication, and adherence to publishing checklists and guidelines.

Independently replicating a subset of experimental results from 50 high-profile papers in the field of cancer biology published between 2010-2012



Aims to release results by end of 2017

DATA SHARING TO ADVANCE SCIENCE AND POLICY MAKING

- **Outbreak Data**
- **City Data**

AB



SABETILAB



Home

Research



BROAD
INSTITUTE

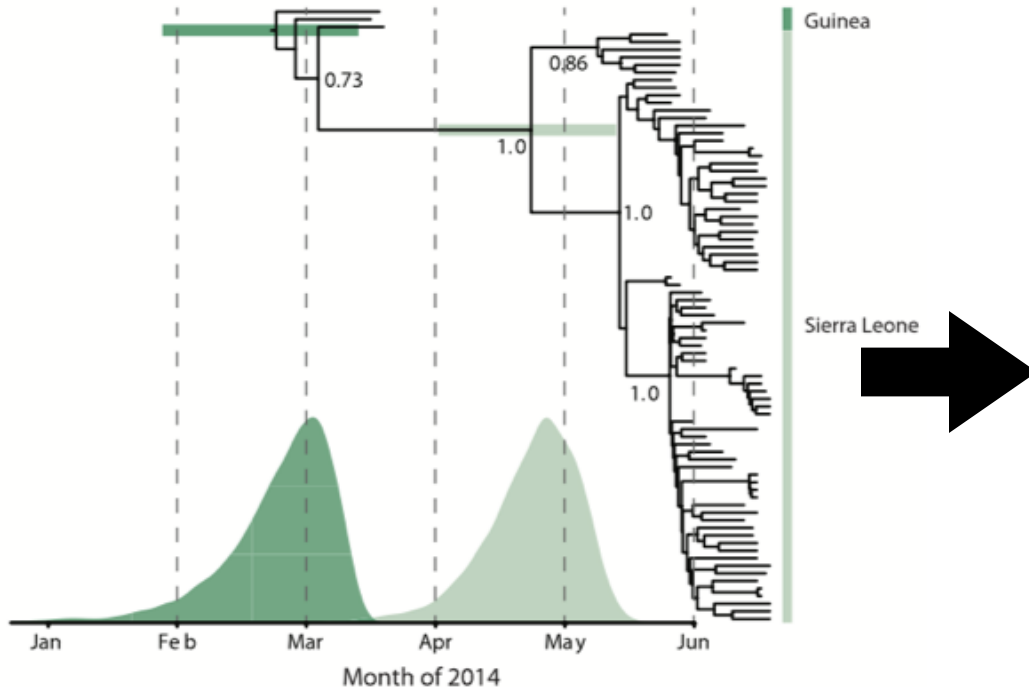
Outreach

WELCOME TO THE SABETI LAB

INFECTIOUS DISEASE RESEARCH

SABETI LAB SHARED EBOLA DATA DURING OUTBREAK

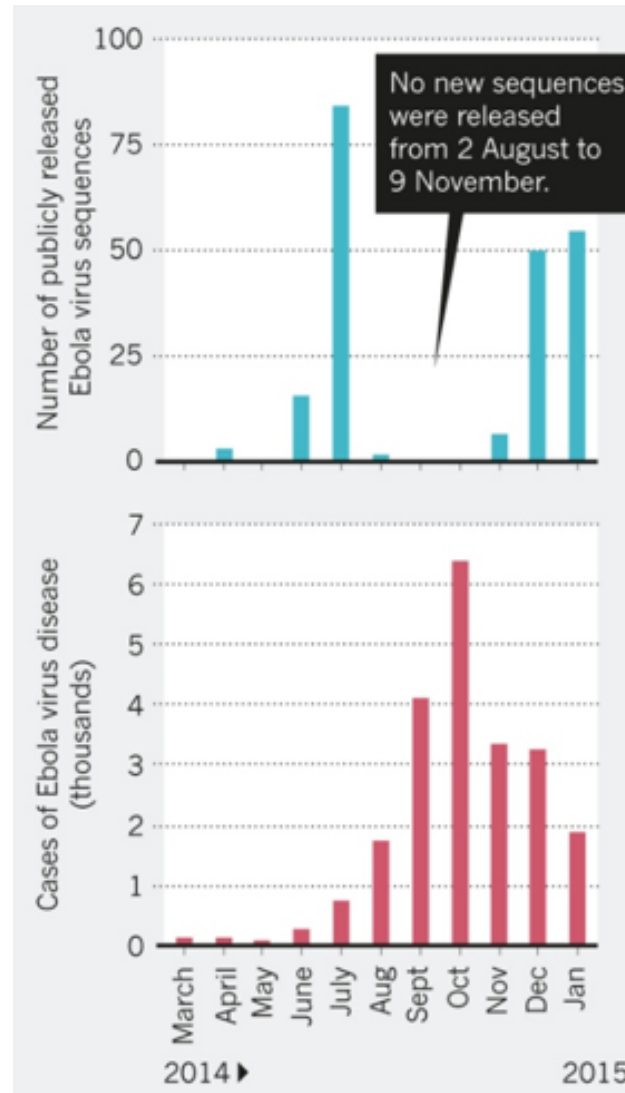
Lab released the first publicly available Ebola sequences (on GenBank), and clinical data (on Harvard Dataverse).



"We were amazed by the surge of collaboration that followed"

Yozwiak, Shaffner, Sabeti, 2015 "Make Outbreak Research Open Access" Nature

BUT, WE DON'T SHARE DATA OFTEN ENOUGH



Gaps in data sharing during the pike of the Ebola outbreak

A young woman with short dark hair and a serious expression is looking directly at the camera. She is wearing a white t-shirt and small hoop earrings. In her hands, she holds a black smartphone and several papers, including one with a colorful pattern and another with a black and white photograph. The background shows a slum environment with makeshift buildings made of corrugated metal and wood, and some laundry hanging in the distance.

**WHY DATA SHARING DOESN'T HAPPEN
MORE OFTEN DURING OUTBREAKS?
ONE REASON IS CONCERN ABOUT
PATIENT PRIVACY**

Image source: Andres Colubri, Sabeti Lab



Boston Area Research Initiative

Area Research Initiative

[EM](#)[About ▾](#)[News & Events ▾](#)[BARI's Spring Conference ▾](#)[Projects ▾](#)[Boston Data Portal](#)[Grants & Fellowships ▾](#)[BARI Comm](#)

Community

[Read more »](#)

WorldMap Bicycle Collisions in Boston (2009-2012)

Sign in | Create Map | View Map | Help

Add Layers Save Identify Link Print Gazetteer About Notes Google Earth Street View Measure

Overlays

- Bicycle Collisions (2009-2012)**
 - ☐ Bicycle Collisions in Boston (2009-2012)
 - ☒ Locations of Bicycle Collisions in Boston
 - 3 Collisions
 - 4-5 Collisions
 - >=6 Collisions**
 - ☐ Road Segments with Bicycle Accidents
- Place Locations
- Transportation
- Health & Human Ecology
- Society & Demographics
- Census 2010
- Census 2000
- American Community Survey (2005-2009)
- Historic Maps
- Local Projects
- Parcels
- Political Boundaries and Areas
- General

Base Maps

- ☐ OpenStreetMap
- ☐ Google Hybrid
- ☐ ESRI World Imagery
- ☒ **Google Terrain**
- ☐ Google Satellite

Map showing Bicycle Collisions in Boston (2009-2012) using Google Terrain. The map displays various neighborhoods including Cambridge, Allston, Brighton, Brookline, and South Boston. Collision locations are marked with green dots of varying sizes, indicating the frequency of collisions (3, 4-5, or 6 or more). The interface includes a top navigation bar with 'Sign in | Create Map | View Map | Help' and a bottom toolbar with 'Add Layers', 'Save', 'Identify', 'Link', 'Print', 'Gazetteer', 'About', 'Notes', 'Google Earth', 'Street View', and 'Measure'. A left sidebar shows a tree view of map layers, with 'Bicycle Collisions (2009-2012)' selected. The map itself shows major roads, parks, and landmarks like MIT and the Boston Common.

Data published at Harvard Dataverse

Cyclist Safety Report



2013



City of Boston
Thomas M. Menino, Mayor

Bicycle data released by BARI was the centerpiece of Boston Mayor's Bike Safety Report



Ethical. Trustworthy. Real.

With crash data, Boston tries to make bicycling safer

Ridership up; 50% drop in injury rate is the goal



23

By Martine Powers | GLOBE STAFF MAY 15, 2013

Bike crashes in Boston are most likely to occur on Massachusetts Avenue, between 4 and 7 p.m., to men between 19 and 31 years old. And they most often occur when a driver fails to see a cyclists or opens a door on a bike.

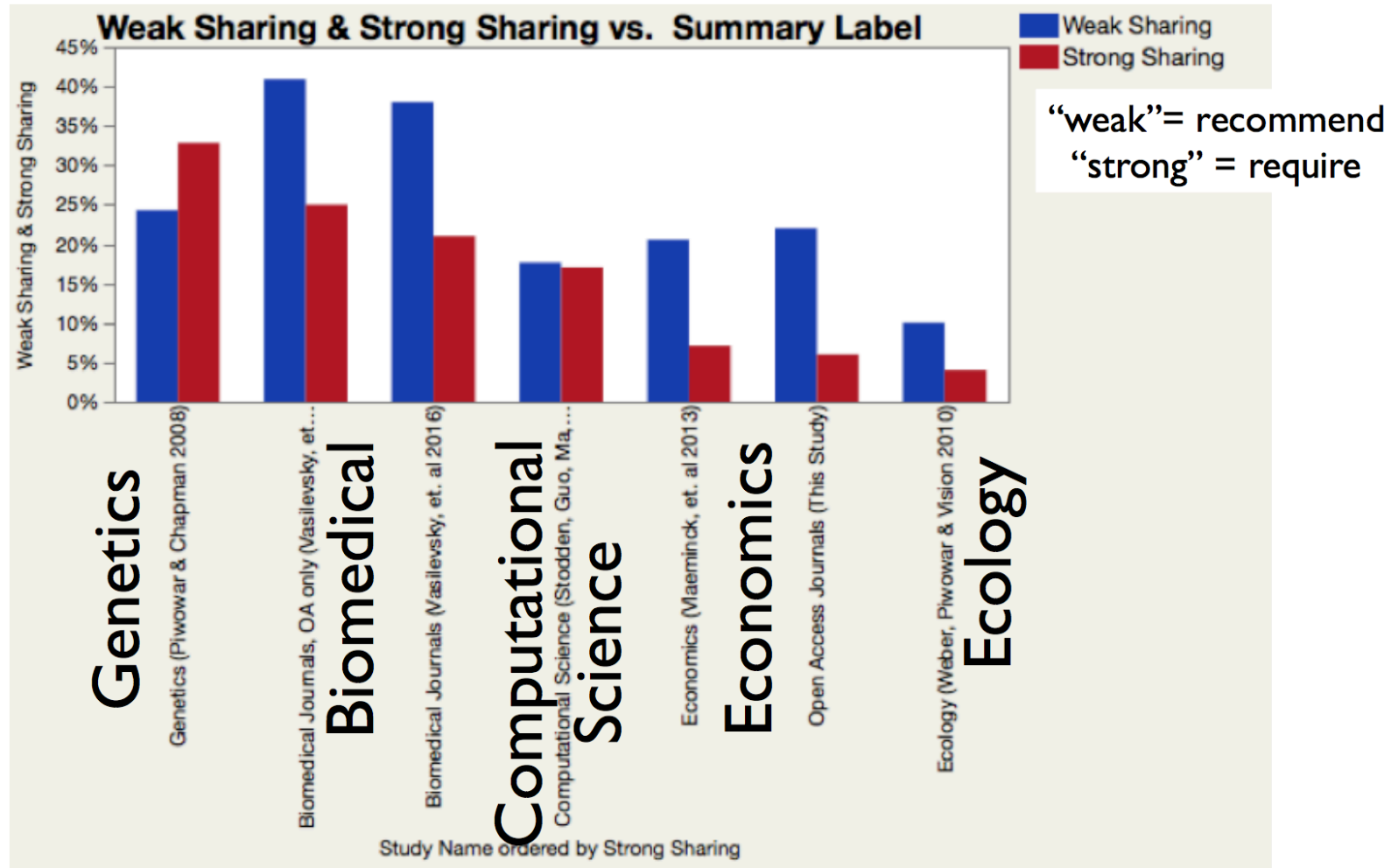
Those insights are contained in a report of unprecedented scale that compiles years of data on bike collisions in the city. With a slew of



HOW CAN WE INCREASE DATA SHARING?

- **New Norms**
- **New Incentives**
- **New Technology**

JOURNAL DATA POLICIES APPLIED ACROSS DISCIPLINES



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

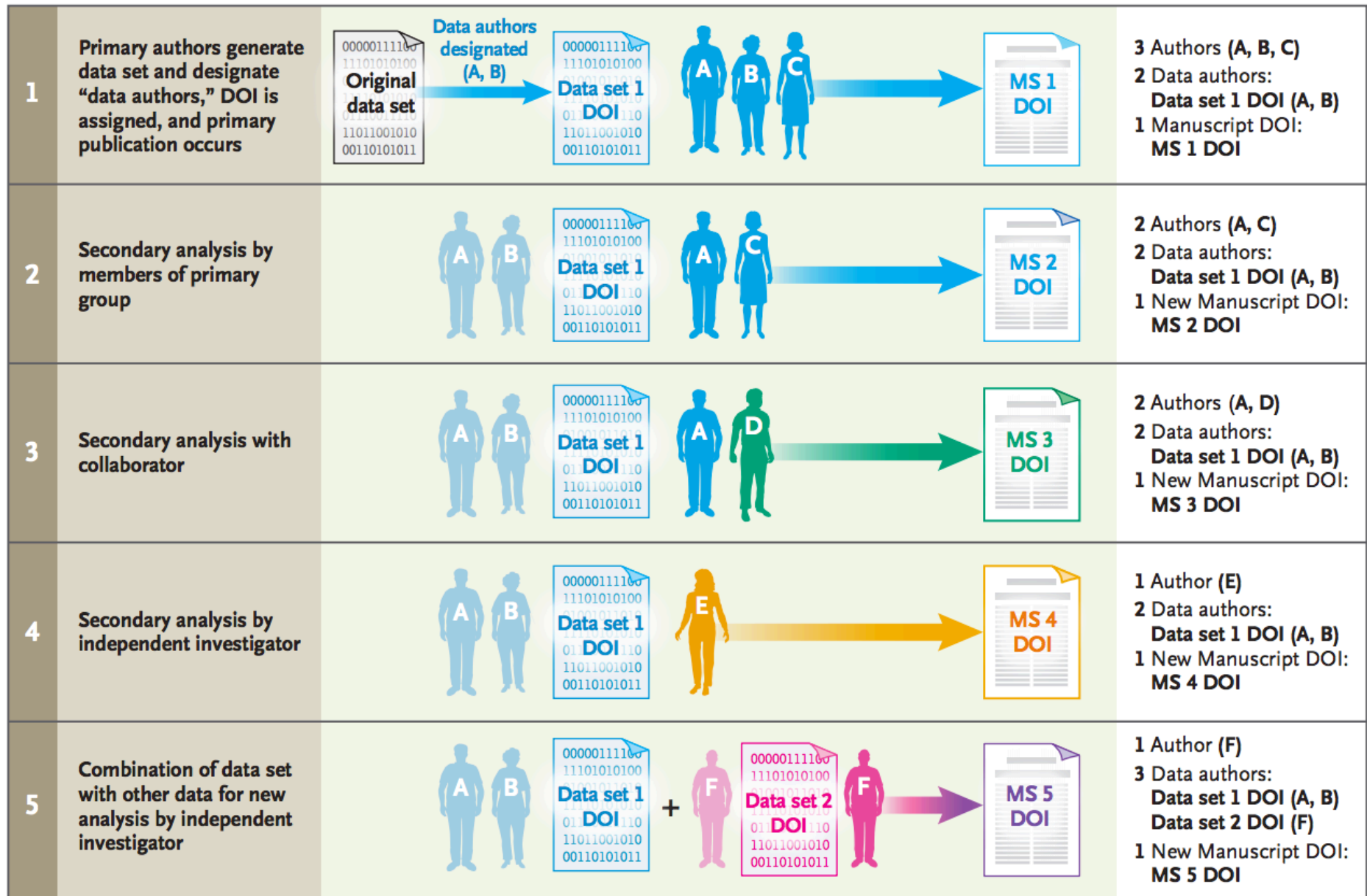
Barbara E. Bierer, M.D., Mercè Crosas, Ph.D., and Heather H. Pierce, J.D., M.P.H.

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."





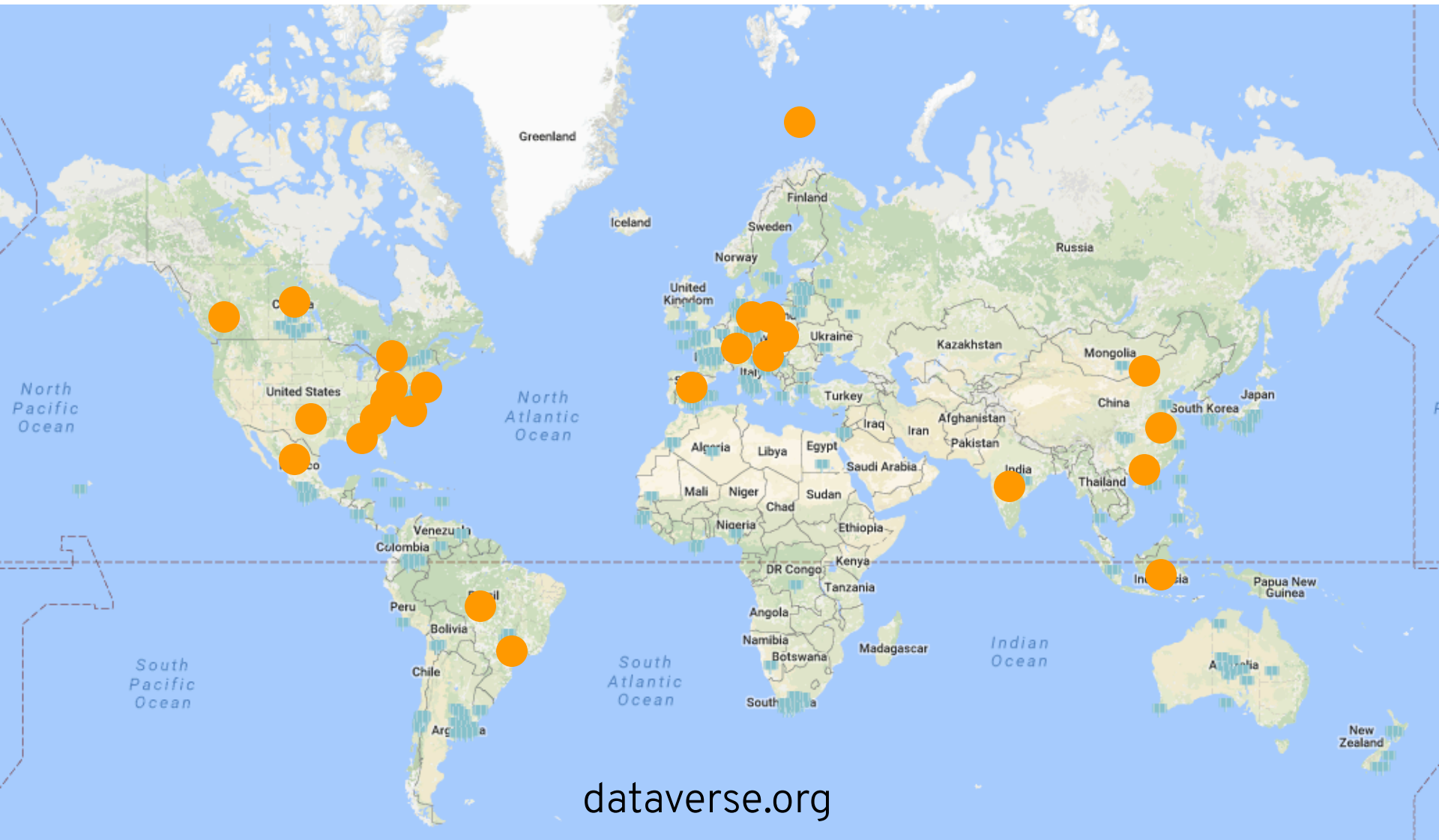
**OUR INSTITUTE PROVIDES A
TECHNOLOGY SOLUTION TO
DATA SHARING**



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

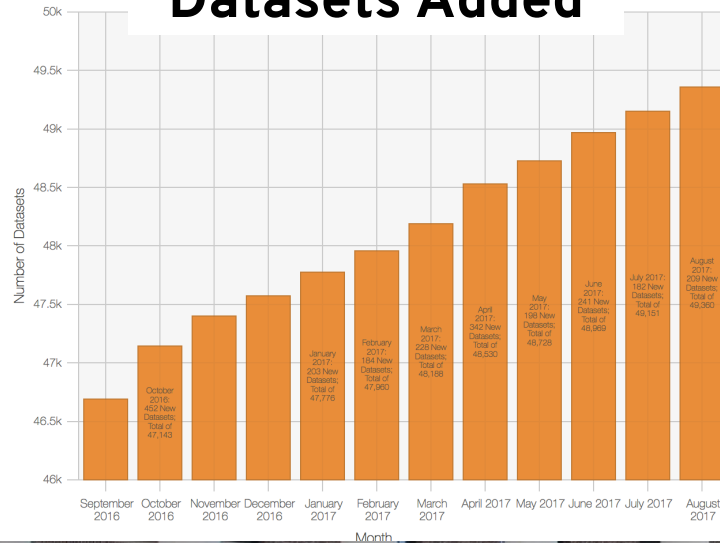
2006 (we started)

2017

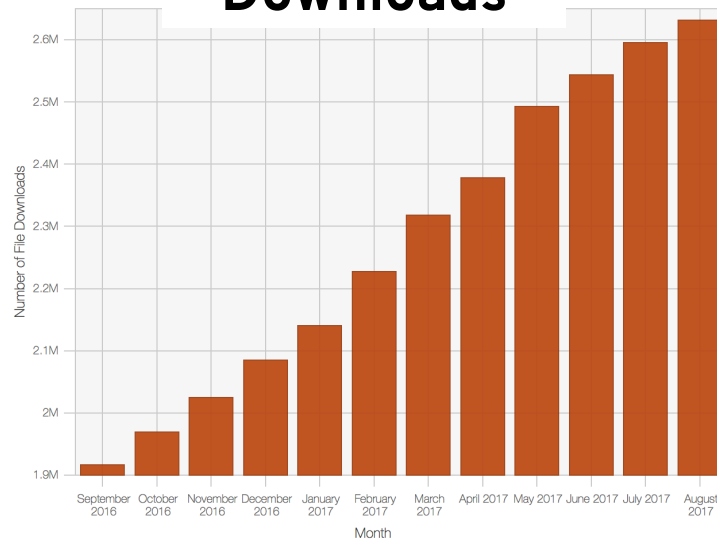


HOW RESEARCHERS SHARE & USE DATA WITH DATAVERSE

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

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 →

DataVerse Q About User Guide Support Sign Up Log In

Clinical Illness and Outcomes in Patients with Ebola in Sierra Leone Version 1.0

Schieffelin, John; Shaffer, Jeffrey; Goba, Augustine; Gbakie, Michael; Gire, Stephen; Colubri, Andres; Sealfon, Rachel; Kanneh, Lansana; Moigboi, Alex; Momoh, Mambu; Fullah, Mohammed; Moses, Lina; Brown, Bethany; Andersen, Kristian; Winnicki, Sarah; Schaffner, Stephen; Park, Daniel; Yozwiak, Nathan; Jiang, Pan-Pan; Kargbo, David; Jalloh, Simbirie; Fonnies, Mbalu; Sinnah, Vandi; French, Issa; Kovoma, Alice; Kamara, Fatima; Tucker, Veronica; Konuwa, Edwin; Sellu, Josephine; Mustapha, Ibrahim; Foday, Momoh; Yillah, Mohamed; Kanneh, Franklyn; Saffa, Sidiki; Massally, James; Boisen, Matt; Branco, Luis; Vandi, Mohamed; Grant, Donald; Happi, Christian; Gevao, Sahr; Fletcher, Thomas; Fowler, Robert; Bausch, Daniel; Sabeti, Parris; Khan, Humarr; Garry, Robert, 2015, "Clinical Illness and Outcomes in Patients with Ebola in Sierra Leone", doi:10.7910/DVN/29296, Harvard DataVerse, V1, UNF:5:wNv/DjKH9a6ELNFuIFm72w==

[Cite Dataset](#) Learn about [Data Citation Standards](#).

Description This data comprises a total of 213 cases evaluated for Ebola virus infection at the Kenema Government Hospital in Sierra Leone between May 25 and June 18, 2014. Outcome data was available for 87 of 106 EBOV positive cases. Metabolic panels were performed on 98 Ebola virus disease and non-Ebola virus disease illness patients with adequate samples volumes. Ebola virus load was determined in 63 cases with adequate samples volumes by quantitative polymerase chain reaction (qPCR) at Harvard University. Sign and symptom data was obtained on 44 patients with a clinical chart that were admitted to Kenema Hospital. The metabolic panels were obtained from serum samples analyzed with a Piccolo Blood Chemistry Analyzer and Comprehensive Metabolic Reagent Discs (Abaxis).

Keyword Ebola, clinical data, laboratory data, outbreak, fatality rate

Related Publication Clinical Illness and Outcomes in Patients with Ebola in Sierra Leone. John S. Schieffelin, et al. N Engl J Med 2014; 371:2092-2100 November 27, 2014 DOI: 10.1056/NEJMoa1411680 doi: 10.1056/NEJMoa1411680

[Files](#) [Metadata](#) [Terms](#) [Versions](#)

Search this dataset... Q Find

3 Files

[Download](#)

ebola-data.tab
Tabular Data - 148.2 KB - Mar 1, 2015 - 94 Downloads
215 Variables, 214 Observations - UNF:5:wNv/DjKH9a6ELNFuIFm72w==
Mirador dataset converted into SPSS format
[SPSS](#)

[Explore](#) [Download](#)

ebola-mirador.zip
ZIP Archive - 15.3 KB - Mar 1, 2015 - 27 Downloads
MD5: 673f468a59d716b9c31b31d22d1b7d04
This dataset is formatted as a project to load into Mirador. It is essentially a collection of CSV files, but with some extra information to make the navigation of the data easier and to display values such as dates and categories more conveniently.
[Mirador](#)

[Download](#)

WHAT ARE WE WORKING ON NOW?

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

DataTags Project: datatags.org

CLOUD DATAVERSE

**COMBINE DATA REPOSITORIES WITH CLOUD
COMPUTING**

Users, External Tools, Services



Deposit

Access

Compute

Software: Services & Tools



Giji

Data Storage



Swift



Cloud Computing



Sahara



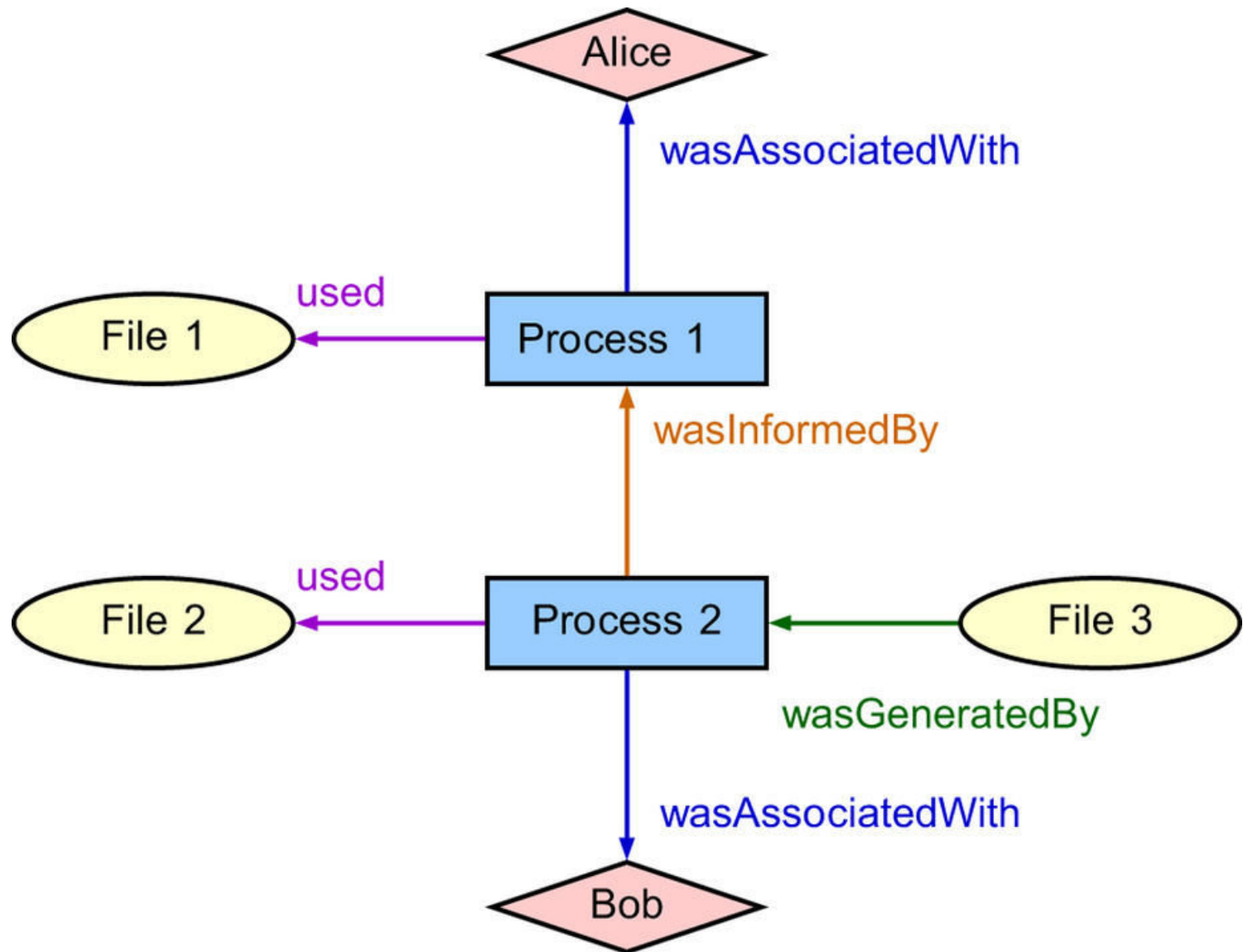
+



FAIR Cloud Dataaverse

DATA PROVENANCE

TRACK THE ORIGINAL SOURCE OF A DATASET



Pasquier, Lau, Trisovic, Boose, Coutierer, Crosas, Ellison, Glibson, Jones, Seltzer, 2017, *If These Data Could Talk*, Nature Scientific Data

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[®]
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

dataverse.org

