

Data Sharing and Open Science in Neuroimaging

Adam Thomas

Data Science and Sharing Team, FMRI, NIMH



Credits

Material borrowed, adapted, and/or stolen from:

- Russ Poldrack



- Chris Gorgolewski



- Brian Nosek



- Tal Yorkoni



- Niko Kriegeskorte



- Tom Nichols



- Phil Bourne



Outline

- Why do we need Open Science?
- What is Open Science?
- How do I do Open Science?

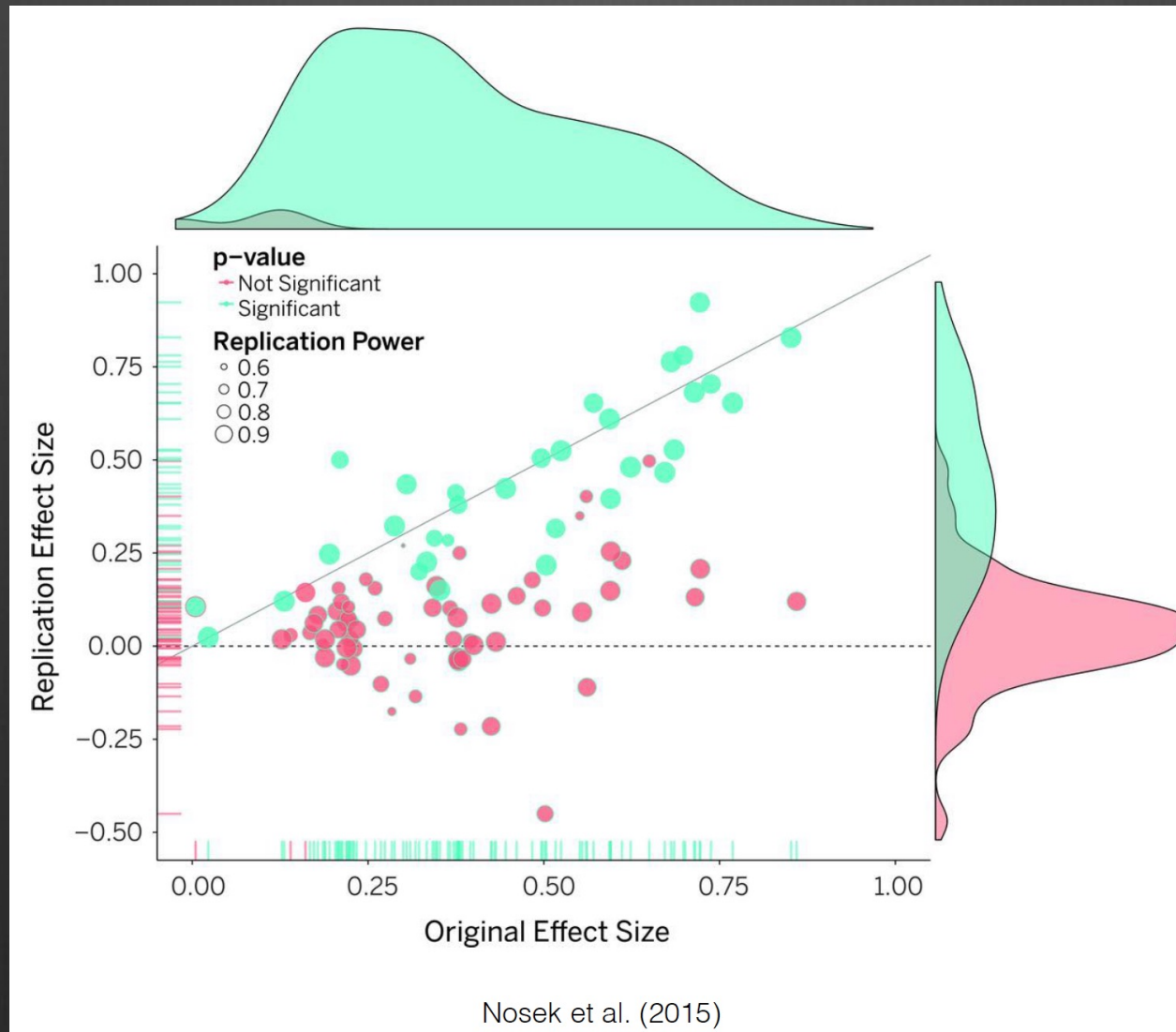
Outline

- Why do we need Open Science?
- What is Open Science?
- How do I do Open Science?

Outline

- The Problem
- What is Open Science?
- How do I do Open Science?

Problem: Reproducibility



The Problem: Reproducibility



Essay

Why Most Published Research Findings Are False

John P.A. Ioannidis

Summary

There is increasing concern that most current published research findings are false. The probability that a research claim is true may depend on study power and bias, the number of other studies on the same question, and, importantly, the ratio of true to no relationships among the relationships probed in each scientific field. In this framework, a research finding is less likely to be true when the studies conducted in a field are smaller; when effect sizes are smaller; when there is a greater number and lesser preselection of tested relationships; where there is greater flexibility in designs, definitions, outcomes, and analytical modes; when there is greater financial and other interest and prejudice; and when more teams are involved in a scientific field in chase of statistical significance. Simulations show that for most study designs and settings, it is more likely for

factors that influence this problem and some corollaries thereof.

Modeling the Framework for False Positive Findings

Several methodologists have pointed out [9–11] that the high rate of nonreplication (lack of confirmation) of research discoveries is a consequence of the convenient, yet ill-founded strategy of claiming conclusive research findings solely on the basis of a single study assessed by formal statistical significance, typically for a p -value less than 0.05. Research is not most appropriately represented and summarized by p -values, but, unfortunately, there is a widespread notion that medical research articles

It can be proven that most claimed research findings are false.

is characteristic of the field and can vary a lot depending on whether the field targets highly likely relationships or searches for only one or a few true relationships among thousands and millions of hypotheses that may be postulated. Let us also consider, for computational simplicity, circumscribed fields where either there is only one true relationship (among many that can be hypothesized) or the power is similar to find any of the several existing true relationships. The pre-study probability of a relationship being true is $R/(R + 1)$. The probability of a study finding a true relationship reflects the power $1 - \beta$ (one minus the Type II error rate). The probability of claiming a relationship when none truly exists reflects the Type I error rate, α . Assuming that c relationships are being probed in the field, the expected values of the 2×2 table are given in Table 1. After a research finding has been claimed based on



Problem: Wasted time & resources



“How much time do you spend handling, reorganizing, and managing your data as opposed to actually *doing* science?”

- Median answer is 80%

Problem: Wasted Time & resources

Unpublished Data

- File drawer problem
- Lost staff & lost metadata
- Underutilized data



The Problem

Lack of transparency and reproducibility hinders integration



NIMH

National Institute of Mental Health
Intramural Research Program (IRP)

Blue Ribbon Panel

June 6, 2008

Final Report

“The Blue Ribbon Panel proposes that basic and clinical groups in NIMH IRP be linked more closely than is generally the case in universities. Linking basic and clinical teams of investigators may facilitate the translational goals of understanding disease mechanisms and developing novel therapies.”

The Problem... is not new

Research in the Service of Mental Health

Summary Report of the Research
Task Force of the National
Institute of Mental Health

A comprehensive and detailed report of the NIMH Research Task Force, totaling over 400 pages, is for sale by the Superintendent of Documents, Government Printing Office, Washington, D.C. 20402. Order DHEW Publication No. (ADM) 75-236 Printed 1975

3. *The Need to Broaden the Use of Research Findings*

The greatest single need in this area is an explicit policy on which to base an Institute-wide effort to disseminate research findings, and, whenever appropriate, to foster their use.

4. *The Need for Synthesis and Integration*


There has been a natural tendency to use research funds mainly for the development of new knowledge. Relatively neglected has been the need to bring together and evaluate findings in a given area, consider them in relation to findings from other mental health research areas, and determine the implications for further research and for application. NIMH should recognize that the synthesis and integration of research results may often be as important as the research itself.

“Relatively neglected has been the need to bring together and and evaluate findings in a given area and consider them in relation to findings from other mental health research areas [...] NIMH should recognize that the synthesis and integration of research results may often be as important as the research itself”

- Research Task Force of the NIMH, 1975



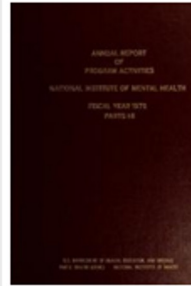

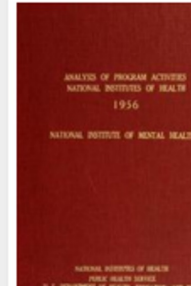
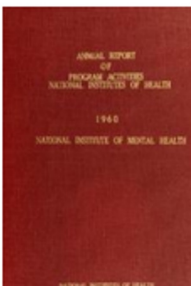
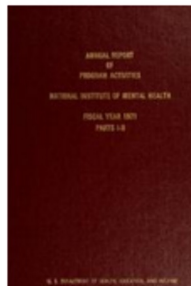
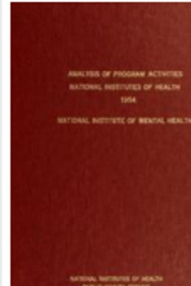
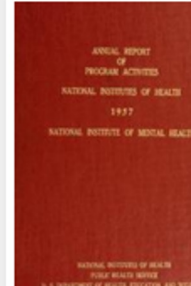
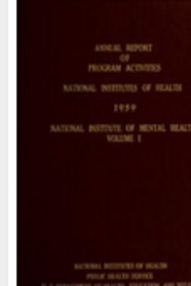
The Problem... is not new

<https://archive.org/details/nimh-nihlibrary>

 **National Institute of Mental Health (NIMH) Publications**
Annual reports and publications created by the National Institute of Mental Health (NIMH) have been digitized and made freely available on the Web as a service of the [NIH Library](#).

[About](#) [Collection](#)

SORT BY RELEVANCE · VIEWS · TITLE · **DATE PUBLISHED** · CREATOR

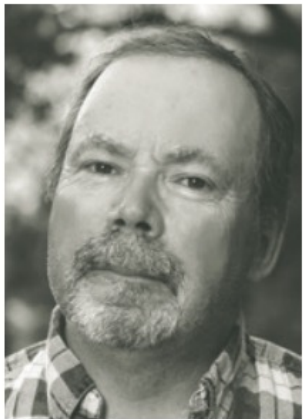
 <p>Contagious and infectious diseases among the Indians 1913</p> <p>194 0 0</p>	 <p>Vol 1967-68 v.1-2: Report of program activities : 1954</p> <p>903 0 0</p>	 <p>Vol 1972 pt.1-2: Report of program activities : 1954</p> <p>990 0 0</p>	 <p>Vol 1955: Report of program activities : 1954</p> <p>1,094 0 0</p>	 <p>Vol 1956: Report of program activities : 1954</p> <p>899 0 0</p>
 <p>Vol 1960: Report of program activities : 1954</p> <p>646 0 0</p>	 <p>Vol 1971 pt.1-2: Report of program activities : 1954</p> <p>1,122 0 0</p>	 <p>Vol 1954: Report of program activities : 1954</p> <p>643 0 0</p>	 <p>Vol 1957: Report of program activities : 1954</p> <p>1,139 0 0</p>	 <p>Vol 1959 v.1: Report of program activities : 1954</p> <p>846 0 0</p>

Problems: The big-data revolution

PERSPECTIVE

Sustaining the big-data ecosystem

Organizing and accessing biomedical big data will require quite different business models, say Philip E. Bourne, Jon R. Lorsch and Eric D. Green.



Biomedical big data offer tremendous potential for making discoveries, but the cost of sustaining these digital assets and the resources needed to make them useful have received relatively little attention. Research budgets are flat or declining in inflation-

recorded. All of this means that absolute numbers are hard to interpret.

These caveats notwithstanding, more details of data usage are needed to inform funding decisions. Over time, such usage patterns could tell us how best to target annotation and curation efforts, establish which data should receive the most attention and therefore incur the largest cost, and determine which data should be kept in the longer term. The cost of data regeneration can also influence decisions about keeping data.

Funders should encourage the development of new metrics to ascertain the usage and value of data, and persuade data resources to provide such statistics for all of the data they maintain. We can learn here from the private sector: understanding detailed data usage patterns through data analytics forms the basis of highly successful companies such as Amazon and Netflix.

FAIR AND EFFICIENT

OPEN SCIENCE:

WHY



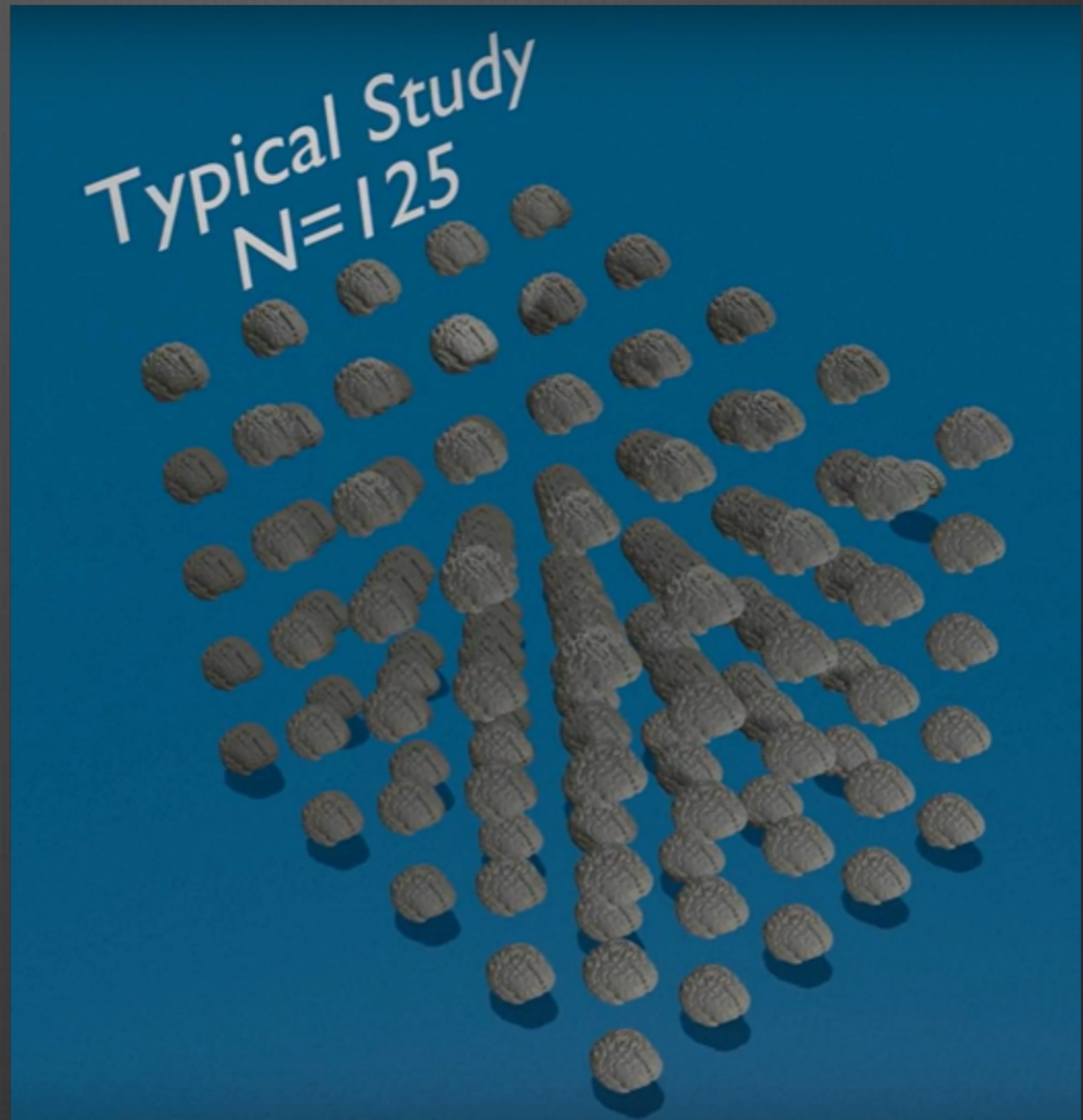
WHAT



HOW

Problems: The big-data revolution

UK Biobank
Imaging
Initiative



OPEN SCIENCE:

WHY



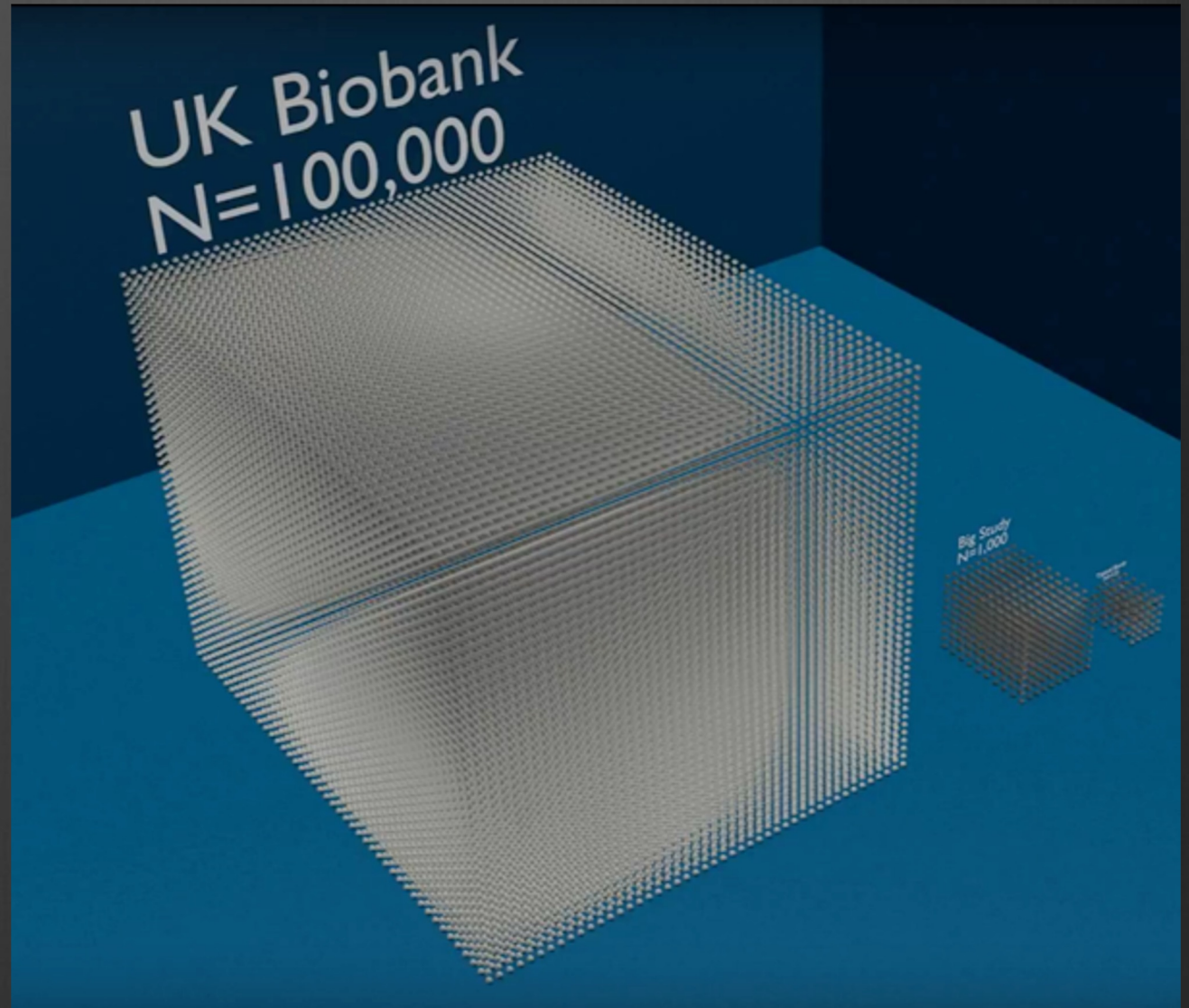
WHAT



HOW

Problems: The big-data revolution

UK Biobank
Imaging
Initiative



OPEN SCIENCE:

WHY



WHAT



HOW

Problems: The big-data revolution



Obama's precision medicine initiative will aim to enroll a large number of people in a genetic database representing the U.S. population.

Amy West/Flickr (CC BY 2.0)

President Obama's 1-million-person health study kicks off with five recruitment centers

By [Jocelyn Kaiser](#) | Jul. 7, 2016, 5:00 PM

OPEN SCIENCE:

WHY



WHAT



HOW

Problems: The big-data revolution



OPEN SCIENCE: WHY → WHAT → HOW

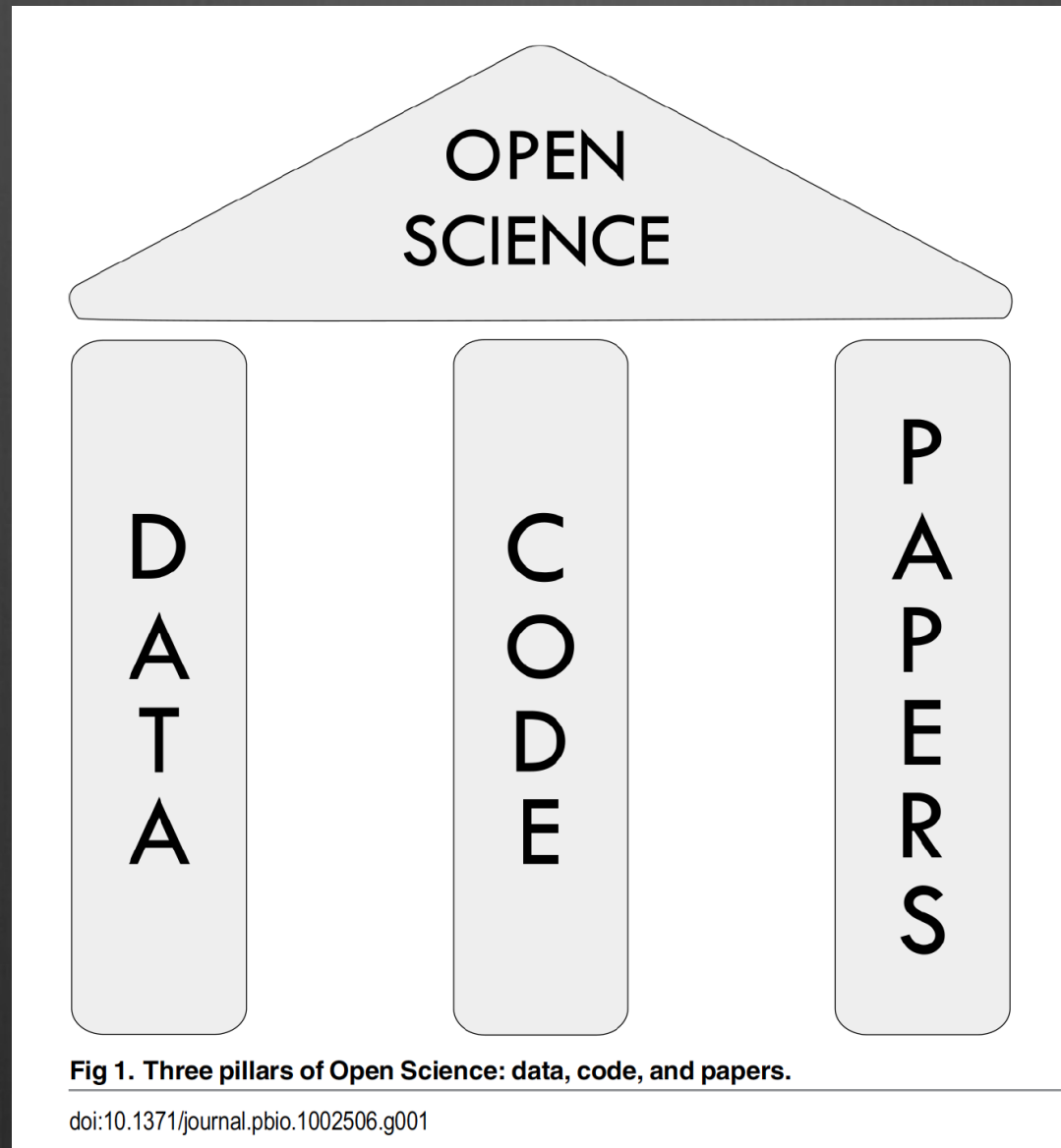
Outline

- The Problem
 - Reproducibility
 - Wasted resources
 - Lack of integration
 - Ill-prepared to work with big datasets
- What is Open Science?
- How do I do Open Science?

Outline

- Why do we need Open Science?
- What is Open Science?
- How do I do Open Science?

What is Open Science?



What is Open Data?

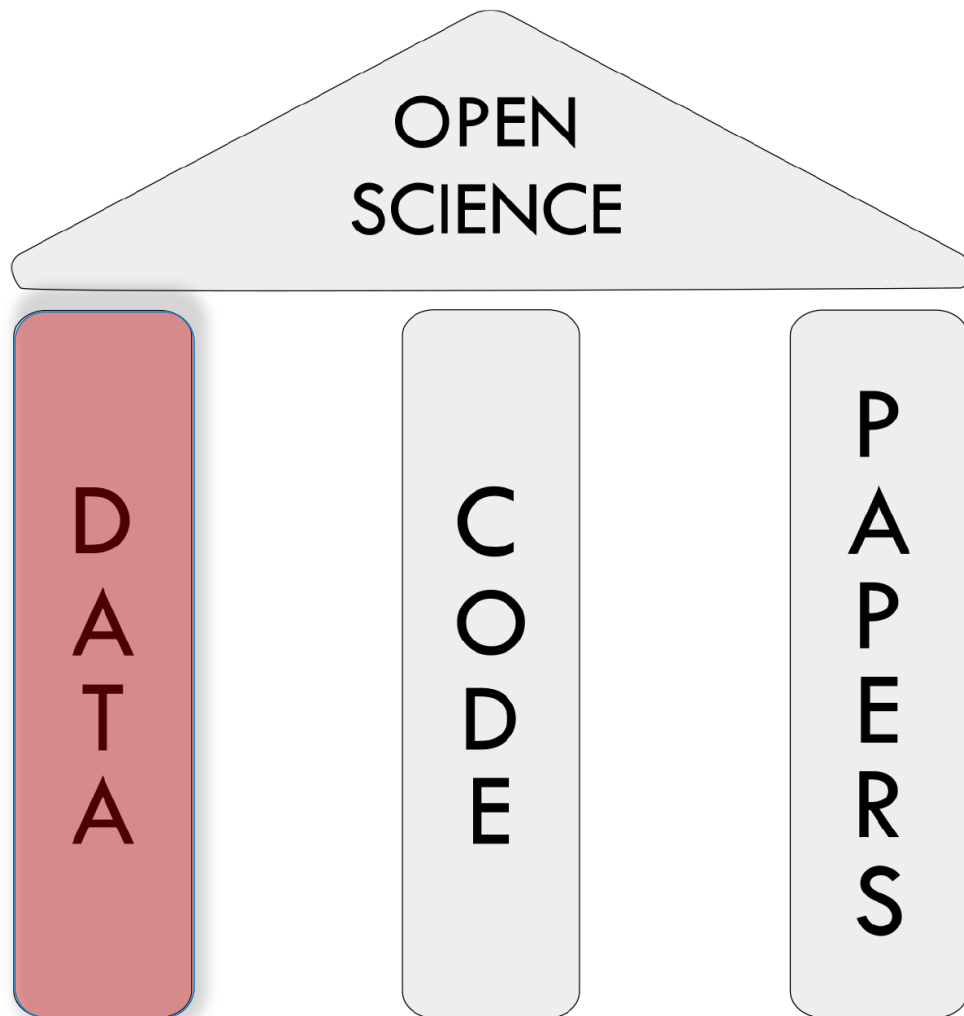


Fig 1. Three pillars of Open Science: data, code, and papers.

doi:10.1371/journal.pbio.1002506.g001

Data deposited in a public, community-recognized repository with a stable DOI

Follows FAIR Principle

- Findable
- Accessible
- Intra-operable
- Reusable

Should be deposited *before* publication

Open Data: Community recognized Repositories

MRI Specific Repos

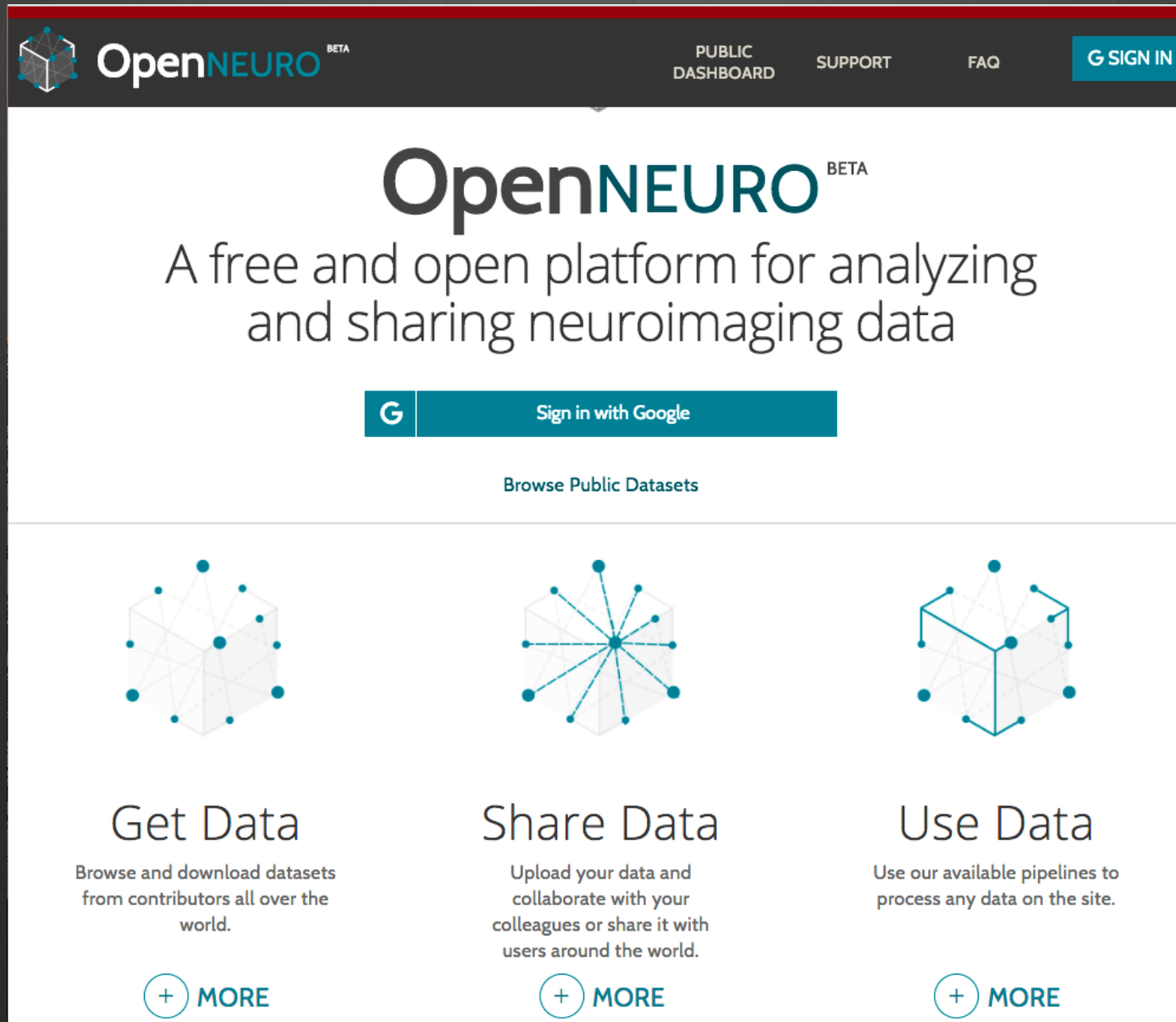
- OpenfMRI / OpenNeuro
- COINS
- FCP/INDI
- LONI
- LORIS
- NITRC
- XNAT Central
- ANIMA*
- BALSA*
- Neuovault*

Data Agnostic Repos

- FigShare
- Dryad
- DataVerse
- Open Science Framework
- NIMH Data Archive

* Statistical & derived data only

Open Data: Community recognized Repositories



The screenshot shows the OpenNEURO website homepage. At the top left is the OpenNEURO logo with a cube icon. To its right are navigation links: PUBLIC DASHBOARD, SUPPORT, and FAQ. A blue button labeled 'SIGN IN' is in the top right. The main heading is 'OpenNEURO BETA' followed by the tagline 'A free and open platform for analyzing and sharing neuroimaging data'. Below this is a 'Sign in with Google' button and a 'Browse Public Datasets' link. The page features three columns: 'Get Data' (browse and download datasets), 'Share Data' (upload and collaborate), and 'Use Data' (use available pipelines). Each column has a 'MORE' button.

OpenNEURO^{BETA}

PUBLIC DASHBOARD SUPPORT FAQ [SIGN IN](#)


OpenNEURO

^{BETA}

A free and open platform for analyzing and sharing neuroimaging data

[G Sign in with Google](#)


[Browse Public Datasets](#)



Get Data

Browse and download datasets from contributors all over the world.


[+ MORE](#)



Share Data

Upload your data and collaborate with your colleagues or share it with users around the world.

[+ MORE](#)



Use Data

Use our available pipelines to process any data on the site.

[+ MORE](#)

OPEN SCIENCE:

WHY



WHAT



HOW

Open Data: NIMH IRP's Repository



OPEN SCIENCE:

WHY



WHAT



HOW

Open Data: NIMH IRP's Repository

The screenshot displays the NIDO (NIMH Open Data) repository interface. The browser address bar shows the URL: <https://nido.nih.gov/datasets/5941a372063a1f000ae15eba>. The page header includes the NIDO logo and navigation links: MY DASHBOARD, PUBLIC DASHBOARD, CONTACT, and an UPLOAD DATASET button. The main content area is titled '100 checkerboard runs JGC 2012' and includes the following information:

- Files:** 1371, **Size:** 8.52GB, **Subjects:** 3, **Sessions:** 11
- Available Tasks:** checkerboard, rest
- Available Modalities:** bold, T1w
- AUTHORS:** Javier Gonzalez-Castilloa, Ziad S. Saad, Daniel A. Handwerker, Souheil J. Inati, Noah Brenowitz, Peter A. Bandettini
- README:** Abstract: The brain is the body's largest energy consumer, even in the absence of demanding tasks. Electrophysiologists report on-going neuronal firing during stimulation or task in regions beyond those of primary relationship to the perturbation. Although the biological origin of consciousness remains elusive, it is argued that it emerges from complex, continuous whole-brain neuronal collaboration. Despite converging evidence

On the right side of the page, there is a 'BIDS Validation' section showing a green checkmark and 'Valid' status, along with '3 WARNINGS'. Below this is a 'Dataset File Tree' showing the following structure:

- 100 checkerboard runs JGC 2012
 - dataset_description.json
 - README
 - task-checkerboard_events.tsv
 - sub-001
 - sub-002
 - sub-003

At the bottom right, there is a figure labeled 'A' showing six axial brain slices with colored overlays representing functional regions. The slices are labeled with coordinates: 15I, 11I, 7I, 9S, 13S, and 17S. A vertical blue bar on the left of the figure contains the text 'PNAS'.

OPEN SCIENCE:

WHY



WHAT



HOW

What is Open Code?

Open code enables greater reproducibility (includes non-code methods)

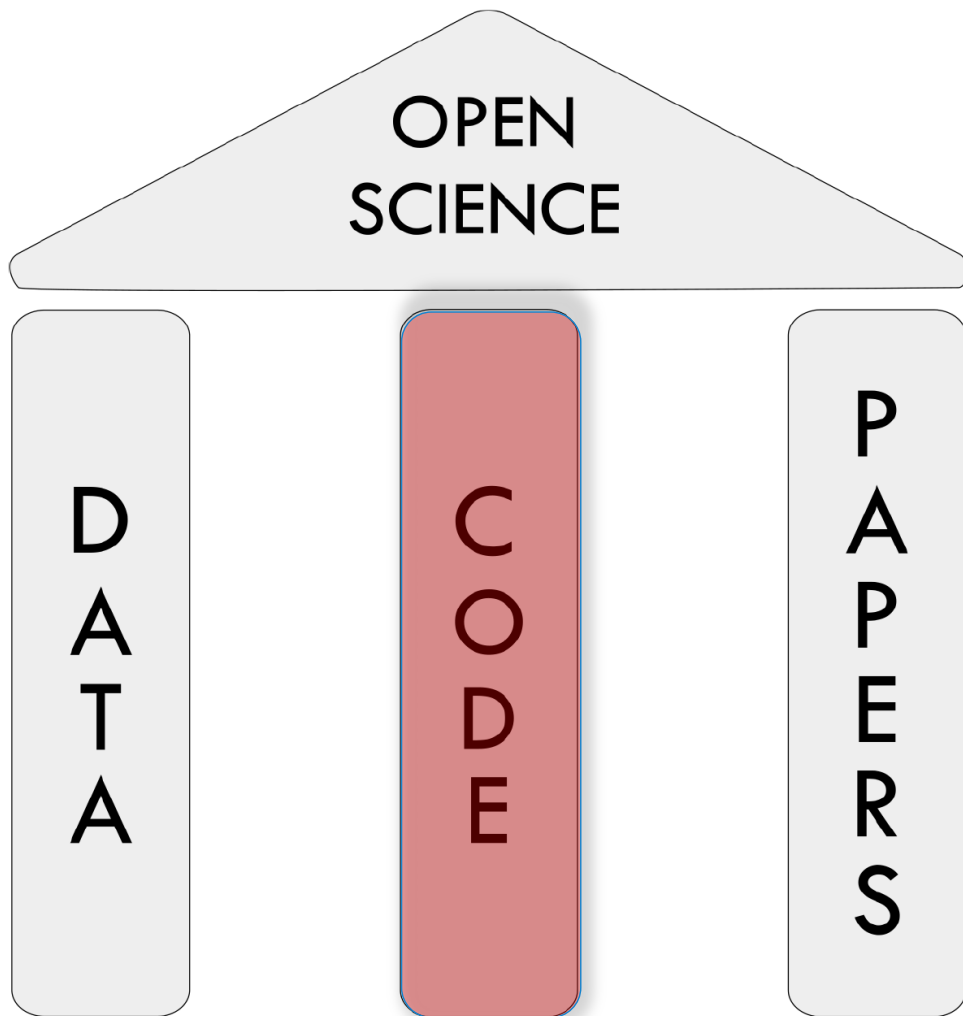
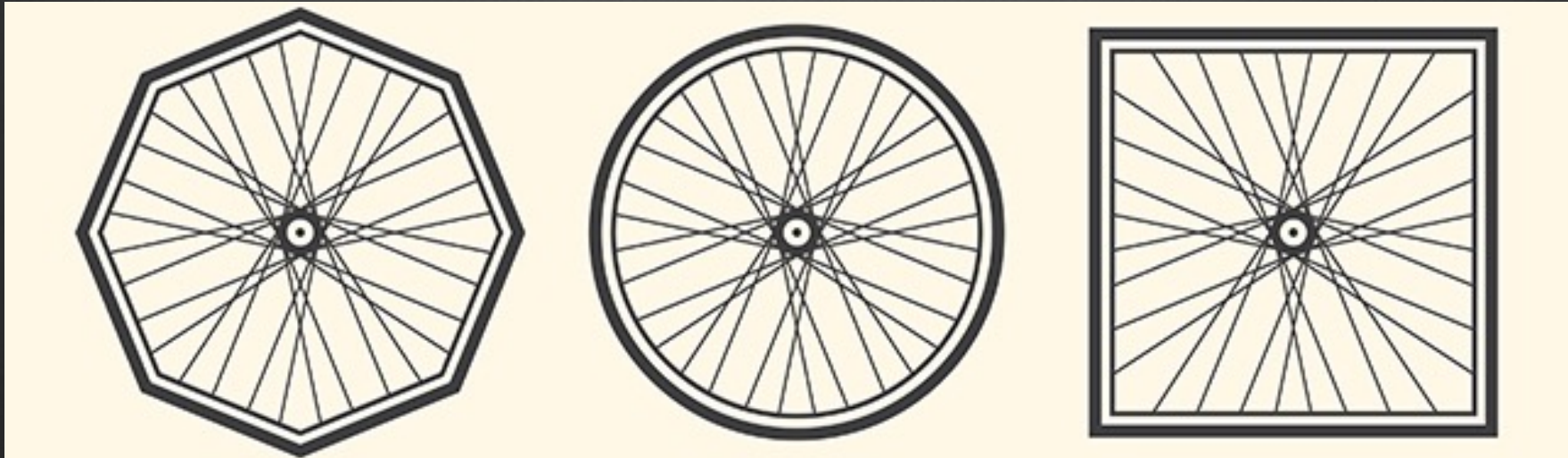


Fig 1. Three pillars of Open Science: data, code, and papers.

doi:10.1371/journal.pbio.1002506.g001

Open Code – Don't Reinvent



Reuse and improve



OPEN SCIENCE:

WHY



WHAT



HOW

Open Code - Version Control

Version control systems allows you to:

- Store all of your analysis in a central repository
- Keep a history of “snapshots” of your evolving analysis
- Quickly switch between different versions of your analysis
- Adopt and modify code from other scientists
- Collaborate



What are Open Papers?

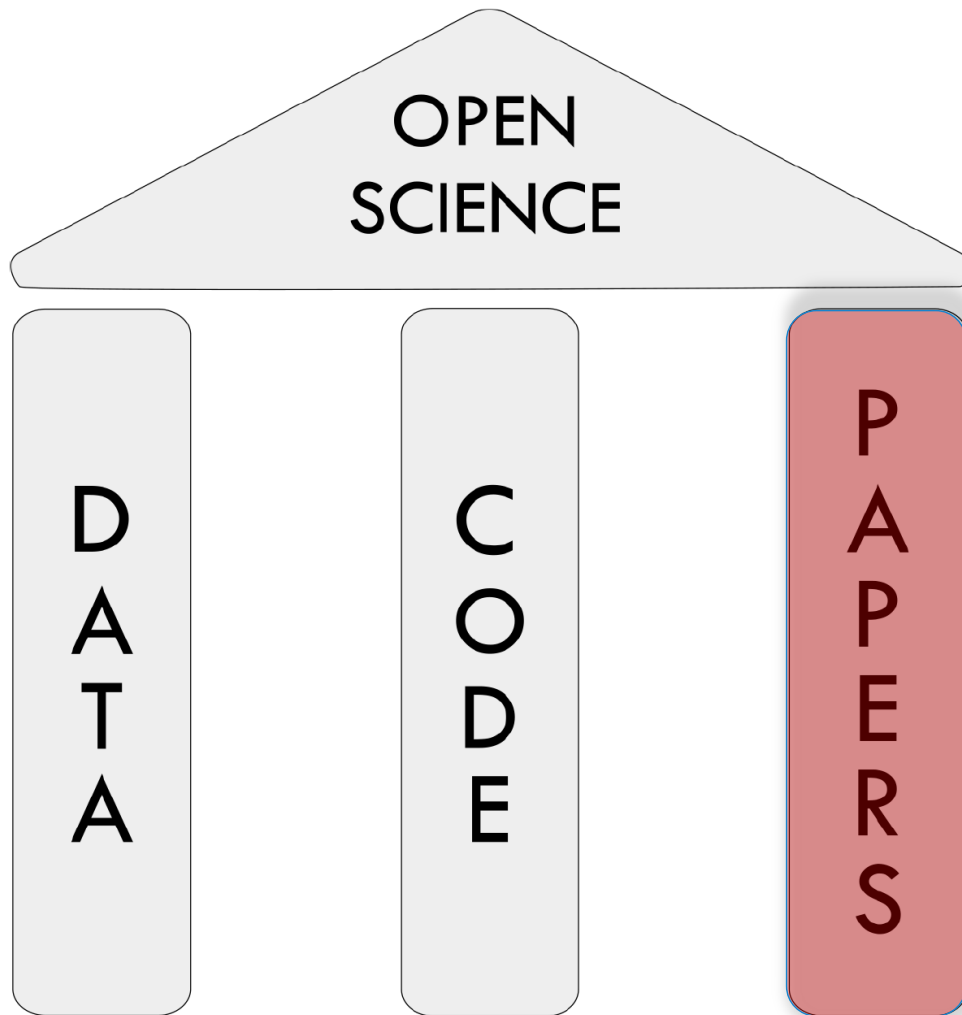


Fig 1. Three pillars of Open Science: data, code, and papers.

[doi:10.1371/journal.pbio.1002506.g001](https://doi.org/10.1371/journal.pbio.1002506.g001)

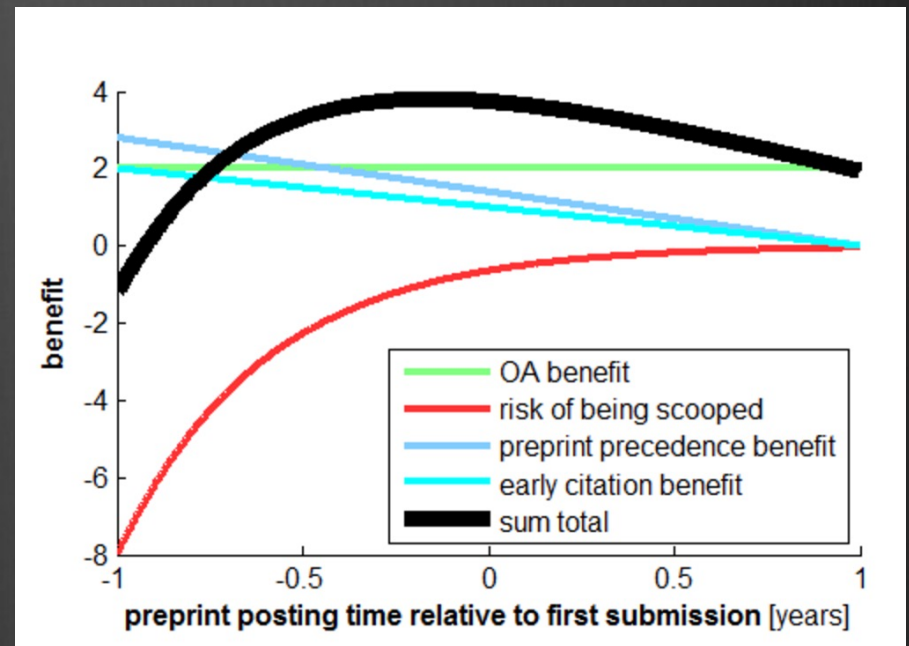
- Preprint posting
- Open access
- Open review

Open Papers: Preprint posting

arXiv.org

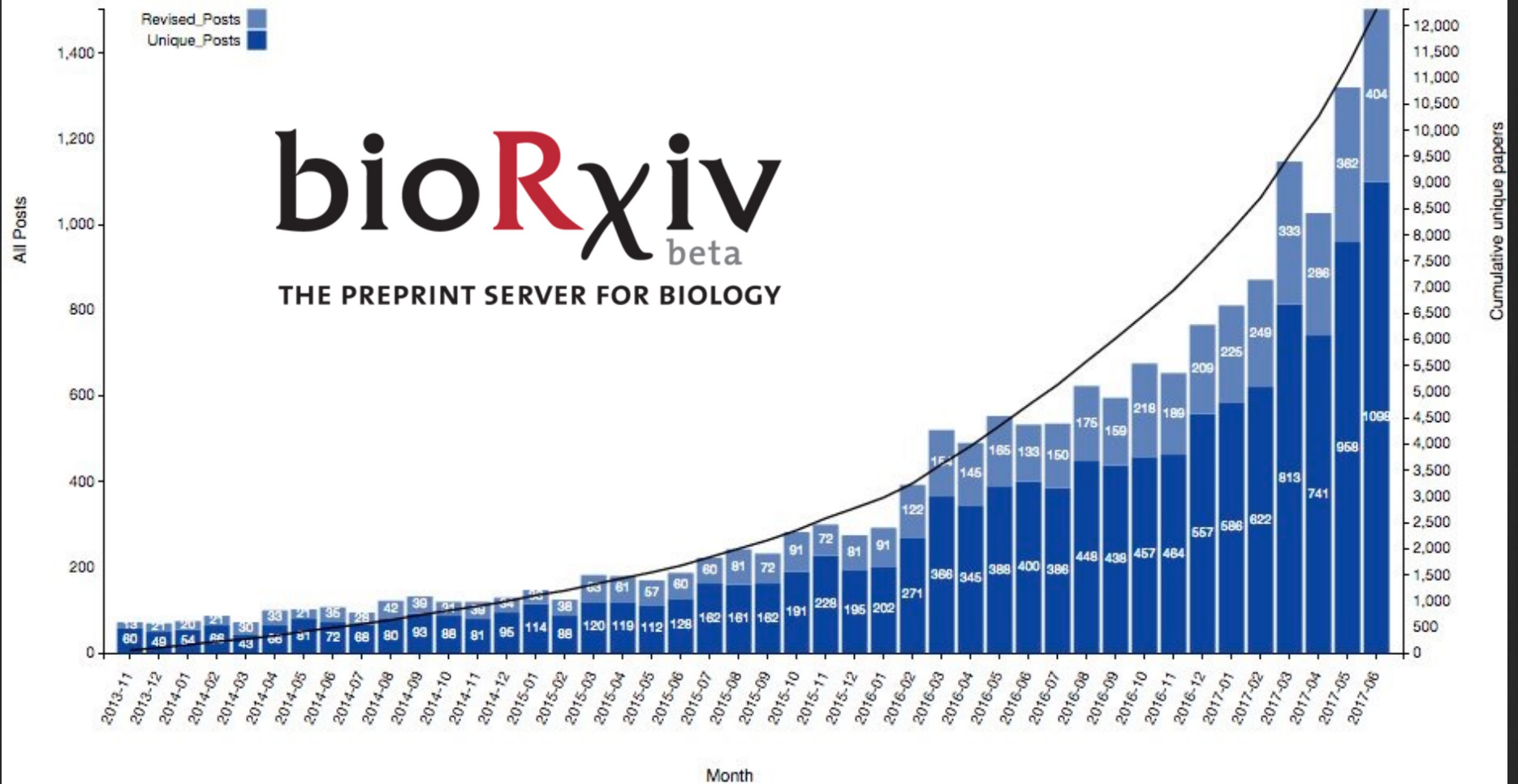
bioRxiv
beta
THE PREPRINT SERVER FOR BIOLOGY

- Benefits:
 - Open access
 - Catch errors
 - Earlier citation
 - Earlier precedence, prevent scooping
 - Speed and improve final submission



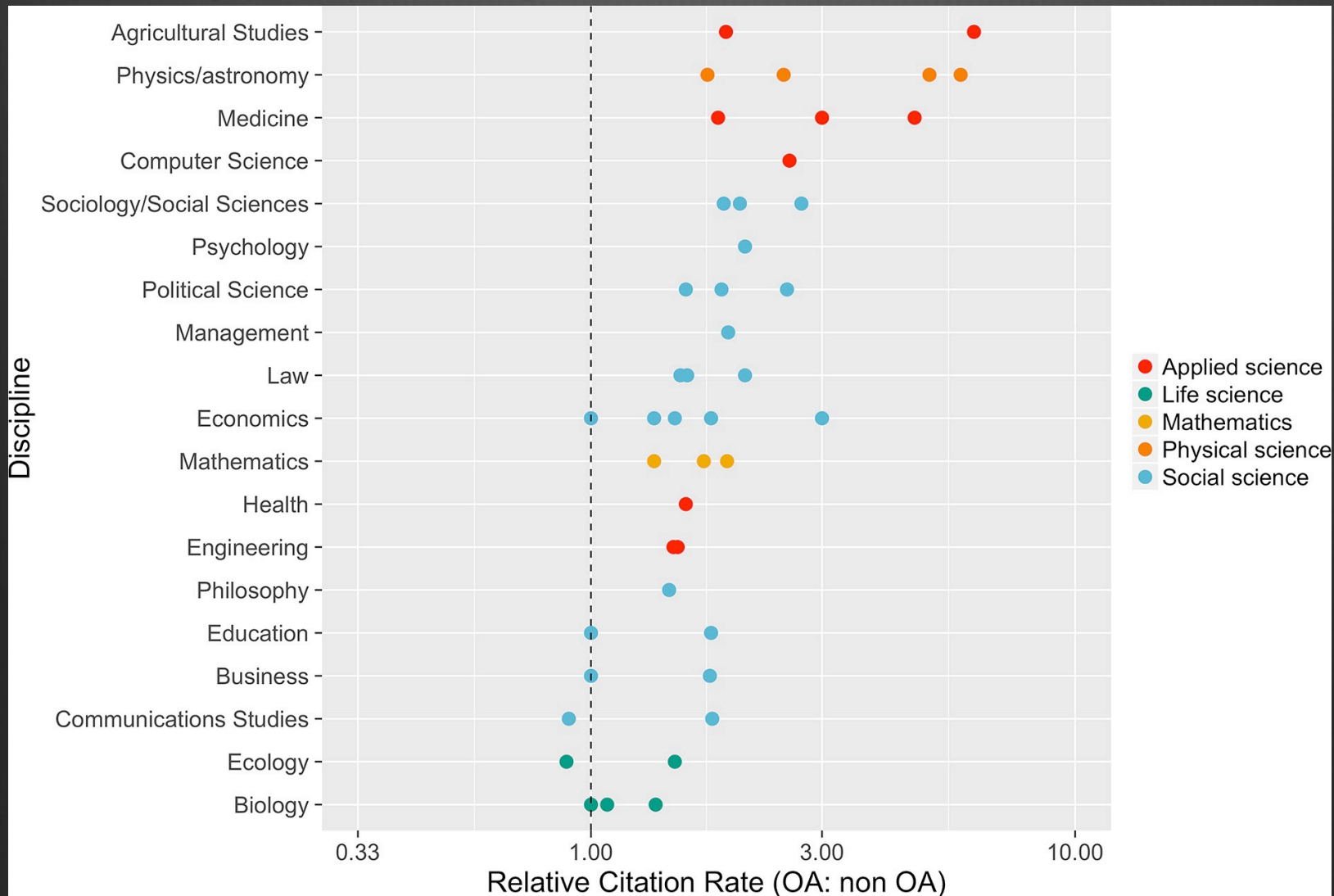
Open Papers: Preprint posting

bioRxiv Content by Month



Open Access

Open access publication are cited more



<https://elifesciences.org/content/5/e16800%20>

mean citation rate of OA articles divided by mean citation rate of non-OA articles

OPEN SCIENCE:

WHY



WHAT



HOW

Open Review

PubPeer
The online journal club

🔍 Search by DOI, PMID, arXiv ID, keyword, author, etc.

The PubPeer database contains all articles. Search results return articles with comments.
To leave a new comment on a specific article, paste a unique identifier such as a DOI, PubMed ID, or arXiv ID into the search bar.

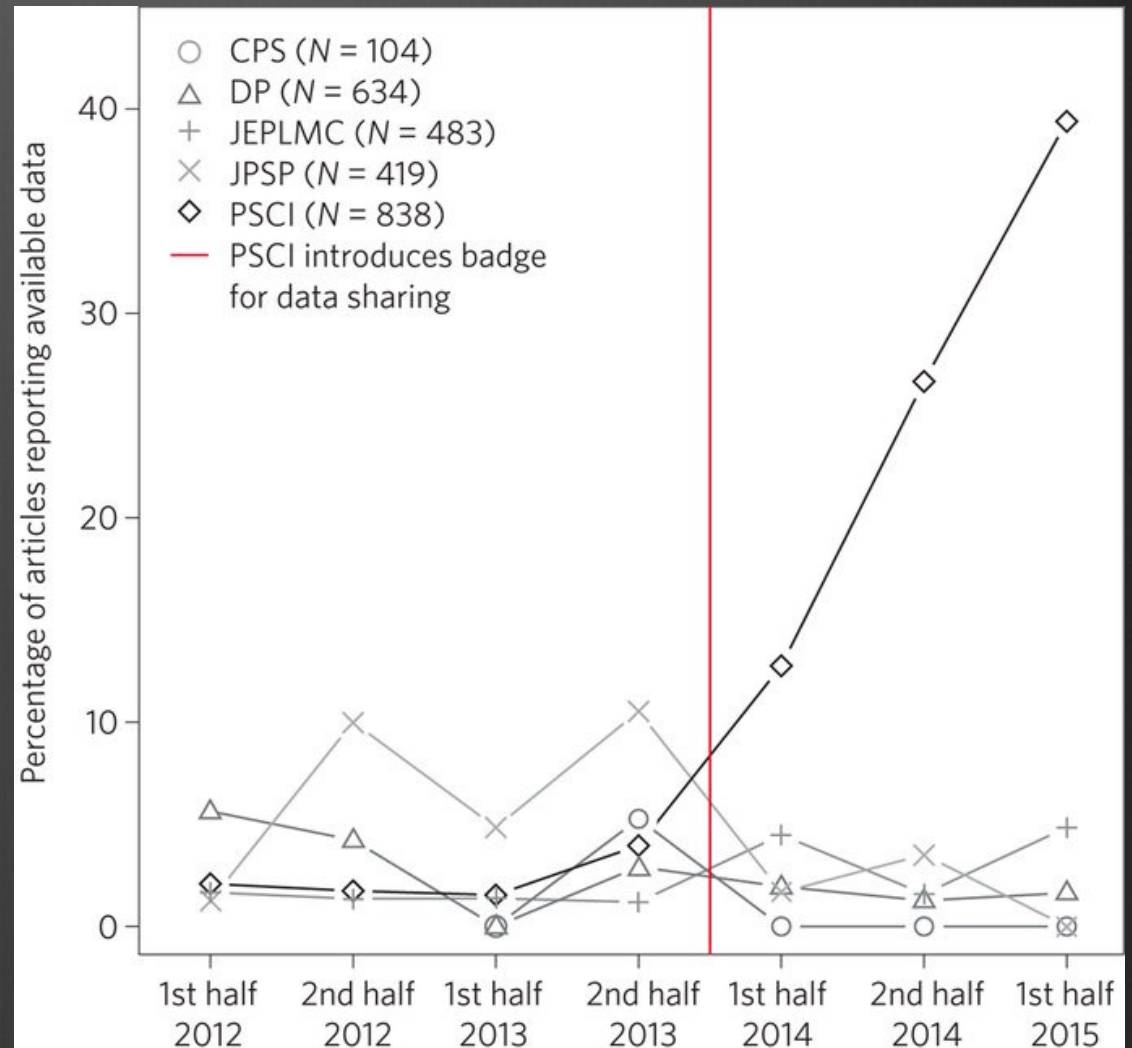
Search Publications

the
WINNOWER

The Winnower is founded on the principle that all ideas should be openly discussed, debated, and archived.

- Public discussion of pros and cons of submission
- Optional anonymity
- Prevent low-quality and or biased review

Incentives: Badges

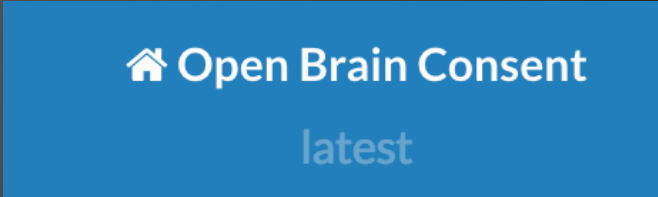


Outline

- Why do we need Open Science?
- What is Open Science?
- How do I do Open Science?

How – Plan Ahead

- Get data sharing in your protocol:
 - NIMH Data Sharing Committee
 - <https://open-brain-consent.readthedocs.io>
- When designing, collecting, and analyzing consult with standards documents:
 - Enhancing Quality and Transparency of Health Research (EQUATOR) <http://www.equator-network.org>
 - Best Practices in Data Analysis and Sharing in Neuroimaging using MRI (COBIDAS) <http://dx.doi.org/10.1101/054262>



Standards – EQUATOR & COBIDAS

- EQUATOR: Different standards for different designs
 - RCT, crossover, observational, etc.
- COBIDAS Sections
 1. Experimental Design
 2. Image Acquisition
 3. Preprocessing
 4. Statistical Modeling
 5. Results
 6. Data Sharing
 7. Reproducibility
- Both EQUATOR and COBIDAS focus on reporting,
- Reviewing them in advance will help you plan and design your study
- Also useful reference when reviewing papers

Standards – EQUATOR & COBIDAS

Checklists



CONSORT 2010 checklist of information to include when reporting a randomised trial*

Section/Topic	Item No	Checklist item	Reported on page No
Title and abstract			
	1a	Identification as a randomised trial in the title	_____
	1b	Structured summary of trial design, methods, results, and conclusions (for specific guidance see CONSORT for abstracts)	_____
Introduction			
Background and objectives	2a	Scientific background and explanation of rationale	
	2b	Specific objectives or hypotheses	
Methods			
Trial design	3a	Description of trial design (such as parallel, crossover, or cluster)	
	3b	Important changes to methods after trial commencement (such as amendments, protocol deviations, and changes in personnel, funding sources, or contracts)	
Participants	4a	Eligibility criteria for participants	
	4b	Settings and locations where the data were collected	
Interventions	5	The interventions for each group with sufficient precision to describe the methods (e.g., drug, dose, duration, and route of administration) and procedures (e.g., randomisation, blinding, and allocation concealment)	
Outcomes	6a	Completely defined pre-specified primary and secondary outcome measures, including definitions and measurement methods	
	6b	Any changes to trial outcomes after commencement	
Sample size	7a	How sample size was determined	
	7b	When applicable, explanation of any interim analyses and stopping rules	
Randomisation:			

Table D.1. Experimental Design Reporting

Aspect	Notes	Mandatory
Number of subjects	<i>Elaborate each by group if have more than one group.</i>	
Subjects approached		N
Subjects consented		N
Subjects refused to participate	Provide reasons.	N
Subjects excluded	Subjects excluded after consenting but before data acquisition; provide reasons.	N
Subjects participated and analyzed	Provide the number of subjects scanned, number excluded after acquisition, and the number included in the data analysis. If they differ, note the number of subjects in each particular analysis.	Y
Inclusion criteria and descriptive statistics	<i>Elaborate each by group if have more than one group.</i>	
Age	Mean, standard deviation and range.	Y
Sex	Absolute counts or relative frequencies.	Y
Race & ethnicity	Per guidelines of NIH or other relevant agency.	N



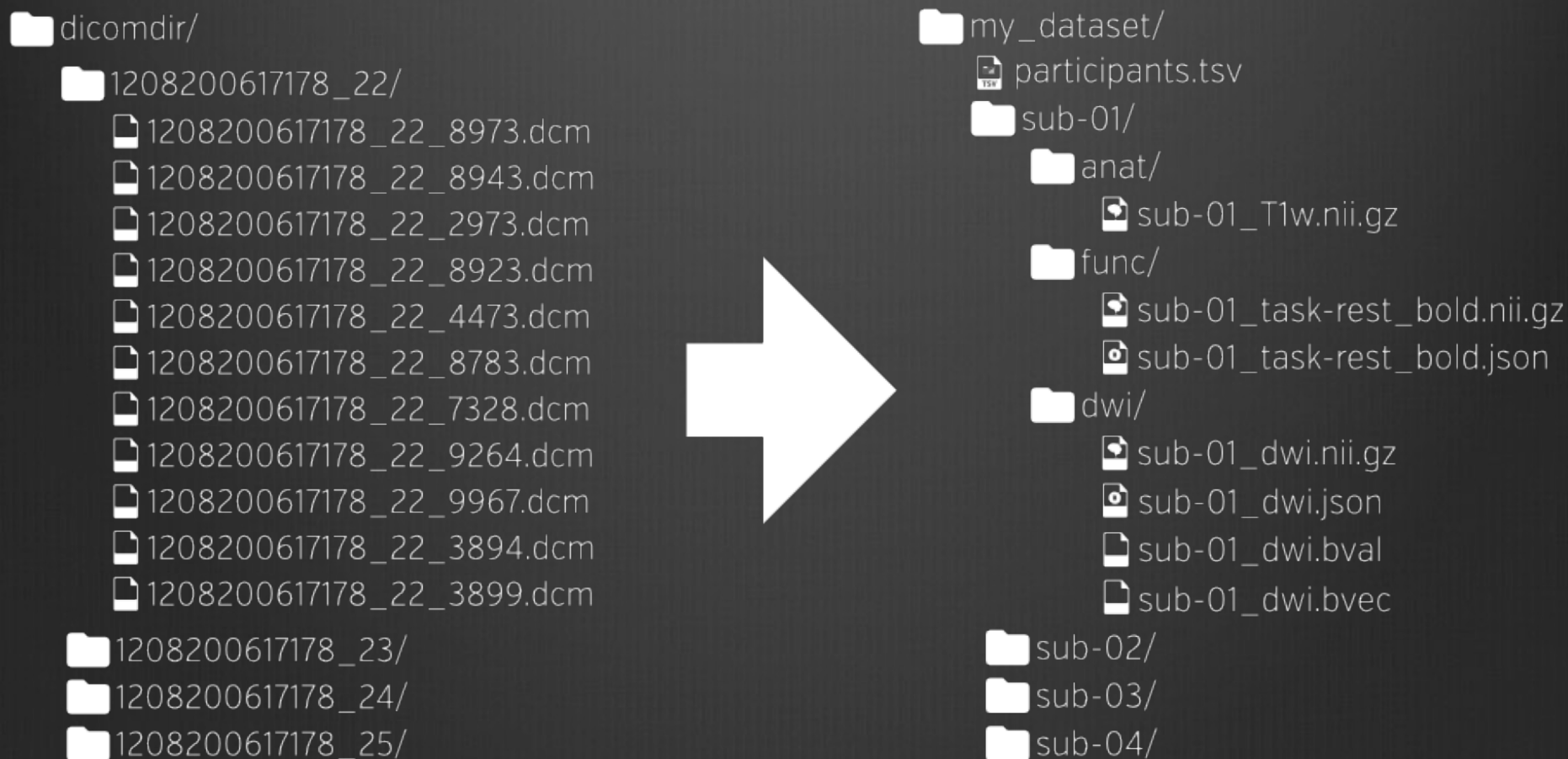
COBIDAS – Highlights

- Report scan parameters by exporting exam cards
- Preprocessing include *all* steps applied to the data before and must be reported
- For maximal transparency, report all regions of interest (ROIs) and/or experimental conditions examined as part of the research, so that the reader can gauge the degree of any HARKing
 - Hypothesizing After The Results are Known
 - It's OK to explore your data, just be clear that that is what you're doing

Organizing your data - BIDS

A simple and intuitive way to organize and describe your neuroimaging and behavioral data.

<http://bids.neuroimaging.io>



How to be Open – Choose your battles

Be open when you can, as you can

Summary of the eight standards and three levels of the TOP guidelines

Levels 1 to 3 are increasingly stringent for each standard. Level 0 offers a comparison that does not meet the standard.

	LEVEL 0	LEVEL 1	LEVEL 2	LEVEL 3
Citation standards	Journal encourages citation of data, code, and materials—or says nothing.	Journal describes citation of data in guidelines to authors with clear rules and examples.	Article provides appropriate citation for data and materials used, consistent with journal's author guidelines.	Article is not published until appropriate citation for data and materials is provided that follows journal's author guidelines.
Data transparency	Journal encourages data sharing—or says nothing.	Article states whether data are available and, if so, where to access them.	Data must be posted to a trusted repository. Exceptions must be identified at article submission.	Data must be posted to a trusted repository, and reported analyses will be reproduced independently before publication.
Analytic methods (code) transparency	Journal encourages code sharing—or says nothing.	Article states whether code is available and, if so, where to access them.	Code must be posted to a trusted repository. Exceptions must be identified at article submission.	Code must be posted to a trusted repository, and reported analyses will be reproduced independently before publication.
Research materials transparency	Journal encourages materials sharing—or says nothing.	Article states whether materials are available and, if so, where to access them.	Materials must be posted to a trusted repository. Exceptions must be identified at article submission.	Materials must be posted to a trusted repository, and reported analyses will be reproduced independently before publication.
Design and analysis transparency	Journal encourages design and analysis transparency or says nothing.	Journal articulates design transparency standards.	Journal requires adherence to design transparency standards for review and publication.	Journal requires and enforces adherence to design transparency standards for review and publication.
Preregistration of studies	Journal says nothing.	Journal encourages preregistration of studies and provides link in article to preregistration if it exists.	Journal encourages preregistration of studies and provides link in article and certification of meeting preregistration badge requirements.	Journal requires preregistration of studies and provides link and badge in article to meeting requirements.
Preregistration of analysis plans	Journal says nothing.	Journal encourages preanalysis plans and provides link in article to registered analysis plan if it exists.	Journal encourages preanalysis plans and provides link in article and certification of meeting registered analysis plan badge requirements.	Journal requires preregistration of studies with analysis plans and provides link and badge in article to meeting requirements.
Replication	Journal discourages submission of replication studies—or says nothing.	Journal encourages submission of replication studies.	Journal encourages submission of replication studies and conducts blind review of results.	Journal uses Registered Reports as a submission option for replication studies with peer review before observing the study outcomes.

How to Open – You don't have to do it alone

- Training

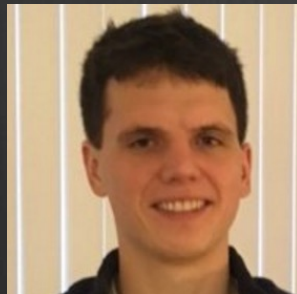


- Asking for help

- Data Science and Sharing Team



Adam Thomas



John Lee

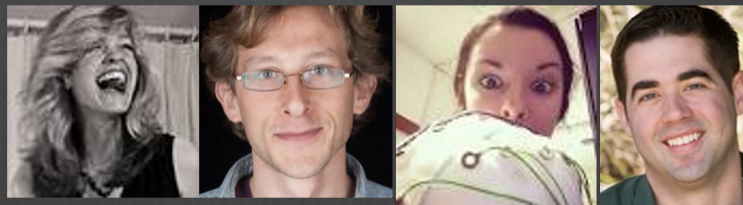


Dylan Nielson



Data Science and Sharing Team's Workshop on Open and Reproducible Neuroscience Mar 13-17th, 2017

- 45 applications, 25 students attended
- 16 hours of instruction on Python, Git, Data Repositories, Biowulf integration, Pre-registration, and statistical rigor
- Instructors from Gallaudet, King's College London, AFNI and Biowulf Teams



- All course material available online:
https://github.com/nih-fmrif/NIMH_repro_2017
- Next course Nov 2017



Data Science and Sharing Team's 2nd Workshop on Open and Reproducible Neuroscience Aug 3-4th, 2017

- I Cover Python, Git, Data Repositories, Biowulf integration, Pre-registration, and statistical rigor
- Instructors from Gallaudet, MIT, & Princeton
 - Regina Nuzzo (Statistics)
 - Satra Ghosh (NiPy)
 - Yarik Halchenko (NeuroDebian)
 - Anisha Keshavan (MindControl)



Summary and Take Homes

- Science is changing (for the better) in both scope (big) and culture (open) to address future challenges
- Open science strives to maximize reproducibility and transparency of data, code, and papers
- Adopting Open Science practices yields benefits in productivity, impact, and reach
- You don't have to do it all at once, and you don't have to do it alone

Thanks!

See online slides for more URLs and references:

<https://github.com/agt24>

Questions?

A manifesto for reproducible science

Marcus R. Munafò^{1,2*}, Brian A. Nosek^{3,4}, Dorothy V. M. Bishop⁵, Katherine S. Button⁶,
Christopher D. Chambers⁷, Nathalie Percie du Sert⁸, Uri Simonsohn⁹, Eric-Jan Wagenmakers¹⁰,
Jennifer J. Ware¹¹ and John P. A. Ioannidis^{12,13,14}

- The Problem
- Methods
 - Cognitive Bias
 - Methodological Training
 - Independent Method support (and oversight)
 - Encouraging Team Science
- Reporting and dissemination
 - Pre-registration
 - Quality of Reporting (checklist & guidelines)
- Reproducibility
 - Transparency
 - Data Sharing
- Evaluation
 - Peer Review
- Incentives
 - Changing cultural norms
 - Badges