

# Creating Incentives for data sharing and transparency

15 May 2024

Adam Thomas

Intramural Research Program  
National Institute of Mental Health  
Bethesda, Maryland, USA  
<http://cmn.nimh.nih.gov/dsst>



# Thank You!

## Tanenbaum Open Science Institute (TOSI)



Guy Rouleau

Director, The Neuro  
Co-Founder of the Tanenbaum Open Science Institute



Annabel Seyller

Chief of Staff, The Neuro  
CEO, Tanenbaum Open Science Institute



**Jean-Baptiste Poline**



Gabriel Pelletier

Open Science Data Manager  
Interim Open Science Alliance Officer



Luisa Pimentel

Open Science Community Officer

# Outline

- Background (Why?)
- Proposal (What?)
- Questions/feedback (How?)

# Recent Progress in Data Sharing

NEWS | 16 February 2022 | Correction [16 February 2022](#)

nature

## NIH issues a seismic mandate: share data publicly

The data-sharing policy could set a global standard for biomedical research, scientists say, but they have questions about logistics and equity.

Nature, Feb 2022 [DOI:10.1038/d41586-022-00402-1](#)



Science, Jan 2023 [DOI:10.1126/science.adg8142](#)

# Recent Progress in Data Sharing

NEWS | 16 February 2022 | Correction [16 February 2022](#)

nature

## NIH issues a seismic mandate: share data publicly

The data-sharing policy could set a global standard for biomedical research, scientists say, but they have questions about logistics and equity.

Nature, Feb 2022 [DOI:10.1038/d41586-022-00402-1](#)

## Data Sharing is the new norm

FEATURES



## READY, SET, SHARE!

As funders roll out new requirements for making data freely available, researchers weigh costs and benefits *By Jocelyn Kaiser and Jeffrey Brainard*

Science, Jan 2023 [DOI:10.1126/science.adg8142](#)

# Recent Progress in Data Sharing

NEWS | 16 February 2022 | Correction [16 February 2022](#)

nature

## NIH issues a seismic mandate: share data publicly

The data-sharing policy could set a global standard for biomedical research, scientists say, but they have questions about logistics and equity.

Nature, Feb 2022 [DOI:10.1038/d41586-022-00402-1](#)

## Data Sharing is the new norm, or is it?

**FEATURES**

**READY, SET, SHARE!**

As funders roll out new requirements for making data freely available, researchers weigh costs and benefits *By Jocelyn Kaiser and Jeffrey Brainard*

Science, Jan 2023 [DOI:10.1126/science.adg8142](#)

# History of Data Sharing in the US

FINAL NIH STATEMENT ON SHARING RESEARCH DATA

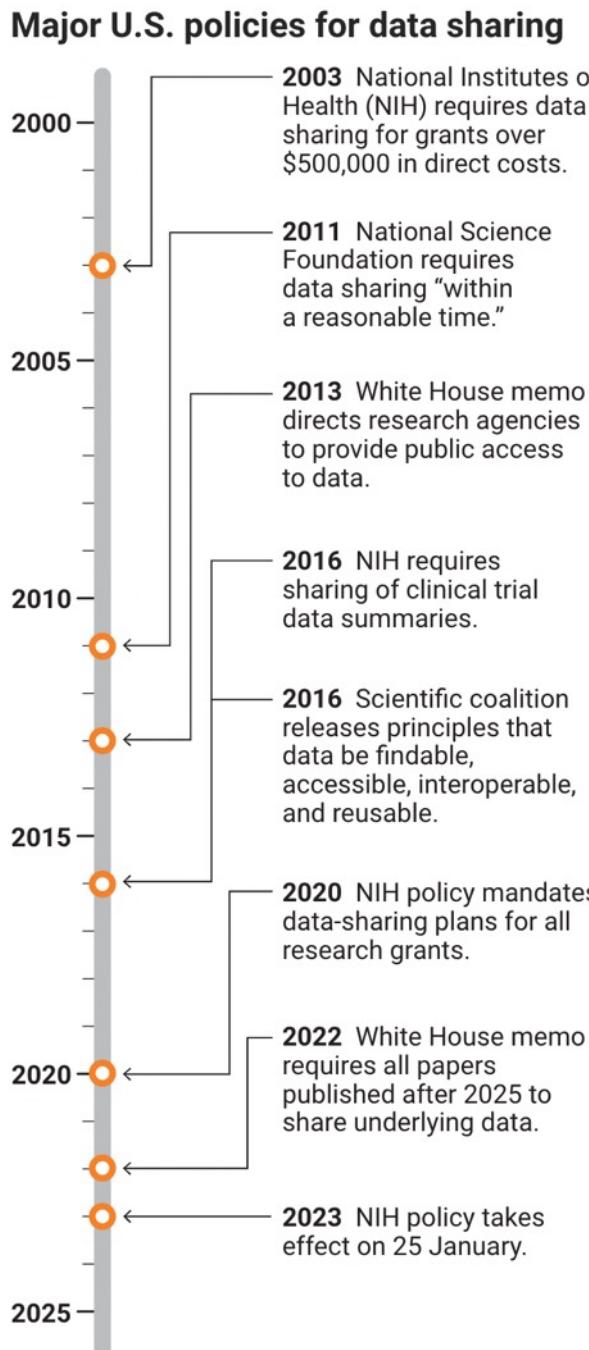
RELEASE DATE: February 26, 2003

NOTICE: NOT-OD-03-032

National Institutes of Health ((NIH))

As part of NIH's long-standing policy to share and make available to the public the results and accomplishments of the activities that it funds, NIH announced and invited comments on a draft statement about the sharing of final research data on March 1, 2002. Since that time,

<https://grants.nih.gov/grants/guide/notice-files/NOT-OD-03-032.html>

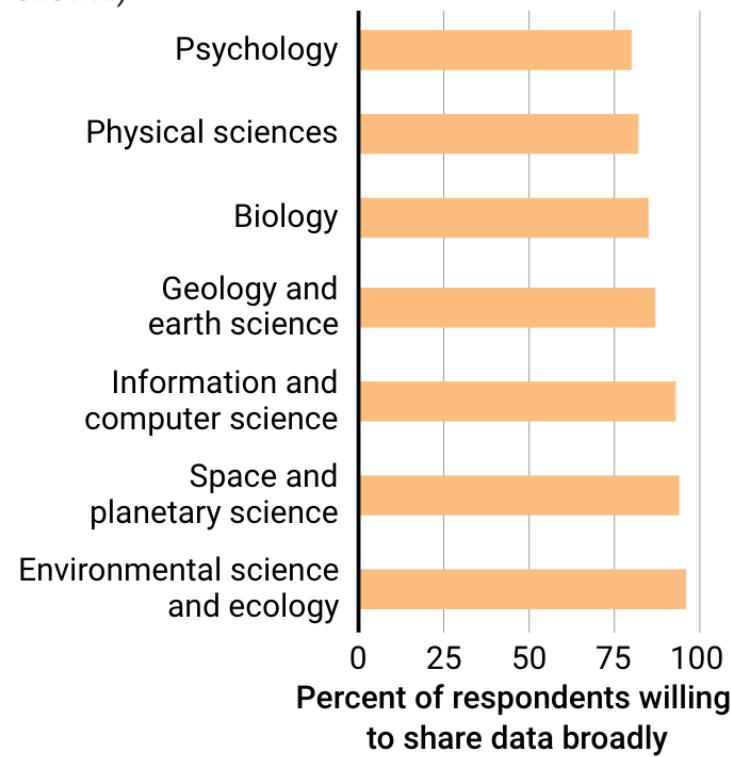


# Data Sharing Attitudes

## Data sharing, by the numbers

### Scientists express interest in sharing their data ...

Across fields, most scientists voice interest in the idea of sharing data, according to a 2017–18 survey of more than 2000 respondents from multiple countries. (Selected categories are shown.)

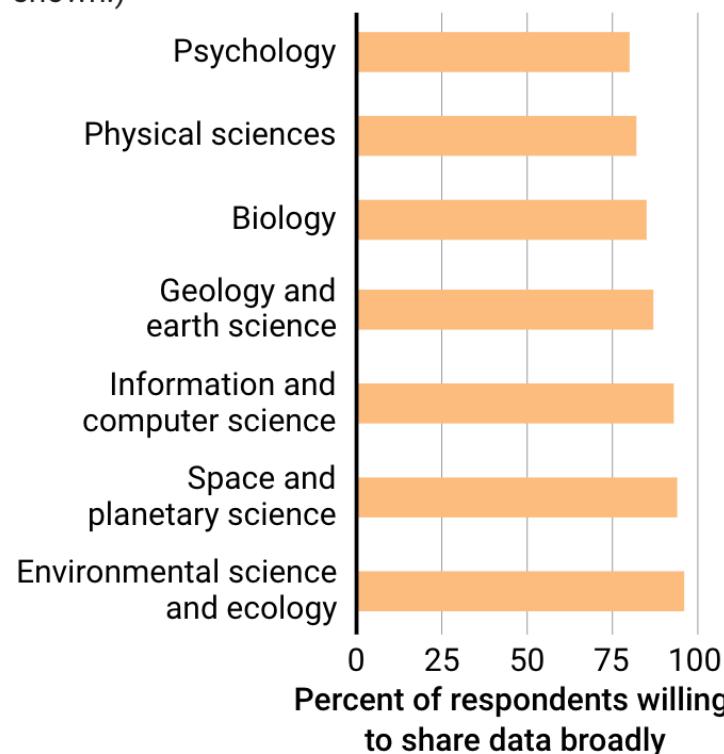


# Data Sharing Attitudes vs Practice

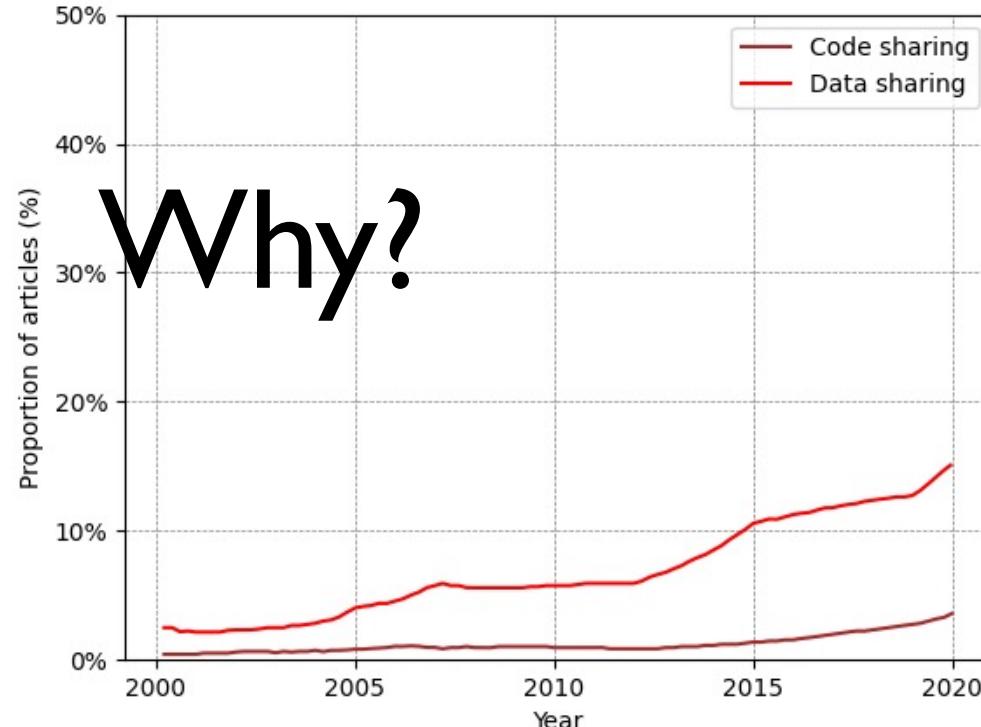
## Data sharing, by the numbers

### Scientists express interest in sharing their data ...

Across fields, most scientists voice interest in the idea of sharing data, according to a 2017–18 survey of more than 2000 respondents from multiple countries. (Selected categories are shown.)



## 2.75M Publications in PubMed Central (PMC)



Adapted from Serghiou et al. 2021 [DOI:10.1371/journal.pbio.3001107](https://doi.org/10.1371/journal.pbio.3001107)

# The Challenge of Cultural Change

- What's need is not just a change in policy but a cultural change
- The existing culture in academia is self-reinforcing
- Changing culture (any culture) is notoriously difficult, especially in the presence of intense competition
- Cultural change can only succeed through a sustained, multi-tiered approach

# The Pyramid of Social Change

Make it REQUIRED - Policy



Make it REWARDING - Incentives



Make it NORMATIVE - Examples



Make it EASY - Support



Make it POSSIBLE - Infrastructure



image src: Center for Open Science

# The Pyramid of Social Change

Killam Seminar Series: Data Sharing and Transparency in Policy and Practice at NIH

Adam Thomas

Data Science Team Lead, NIMH Intramural Research Program

Host: Jean-Baptiste Poline

Tuesday, May 28, 2024

4 p.m. EDT

Jeanne Timmins Amphitheatre, The Neuro, 3801 University

Registration Required

Make it EAST

- Support

Make it POSSIBLE

- Infrastructure



image src: Center for Open Science

# The Pyramid of Social Change

Killam Seminar Series: Data Sharing and Transparency in Policy and Practice at NIH

Adam Thomas

Data Science Team Lead, NIMH Intramural Research Program

Host: Jean-Baptiste Poline

Tuesday, May 28, 2024

4 p.m. EDT

Jeanne Timmins Amphitheatre, The Neuro, 3801 University

Registration Required

Make it EAST

- Support

Make it POSSIBLE

- Infrastructure



Data Science and Sharing Team  
Board of Scientific Counselors Review

[https://bit.ly/dsst\\_bsc2024](https://bit.ly/dsst_bsc2024)

December 2023

## Contents

<b>1 Making Data Sharing Possible</b>	3
1.1 Creating and Expanding Data Standards	3
1.2 Creating and Expanding Data Repositories	4
<b>2 Making Data Sharing Easy</b>	4
2.1 Providing Training	4
2.2 Creating Tools for Managing Data	5
2.3 Using Large Shared Datasets	5
2.4 Sharing IRP Datasets	5
<b>3 Making Data Sharing Normative</b>	6
3.1 fMRI Correlates of Positive and Negative Phenotypes	7
3.2 Factors Influencing Mental Well-being in Response to COVID Pandemic	8
3.3 Leveraging Large Open Data for Machine Learning	8

# Alternative Incentives



Make it REWARDING

Alternative incentives

The large blue triangle contains the text "Make it REWARDING" on the left and "Alternative incentives" on the right. In the center is a white icon of a medal with a letter "I" on it.



# Measuring Data Sharing



- To reward data sharing there must be a mechanism to measure it
- Text mining tools exist to identify data sharing reported in publications
- However, aggregate information on reported data sharing is not easily accessible to funders and policy makers



~25,000 Unique  
Pubs 2019-2023



# Measuring Data Sharing



- NIH is composed of 25 distinct institutes and centers (ICs)



# Measuring Data Sharing



- To reward data sharing there must be a mechanism to measure it
- Text mining tools exist to identify data sharing reported in publications
- However, aggregate information on reported data sharing is not easily accessible to funders and policy makers

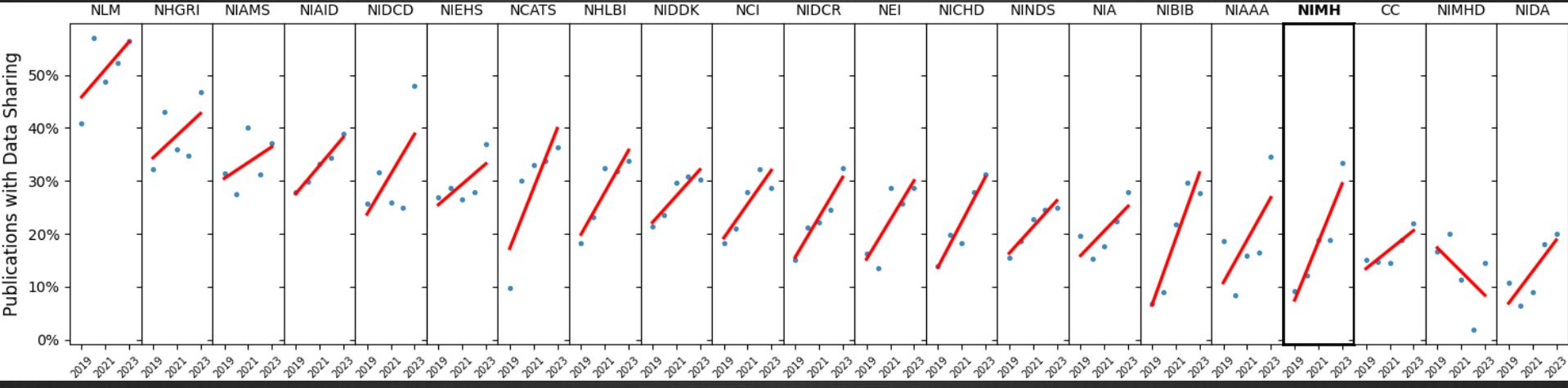


~25,000 Unique Pubs  
2019-2023



## Open Data Detection in Publications (ODDPub)

build error codecov 90% License MIT DOI 10.5281/zenodo.4071699



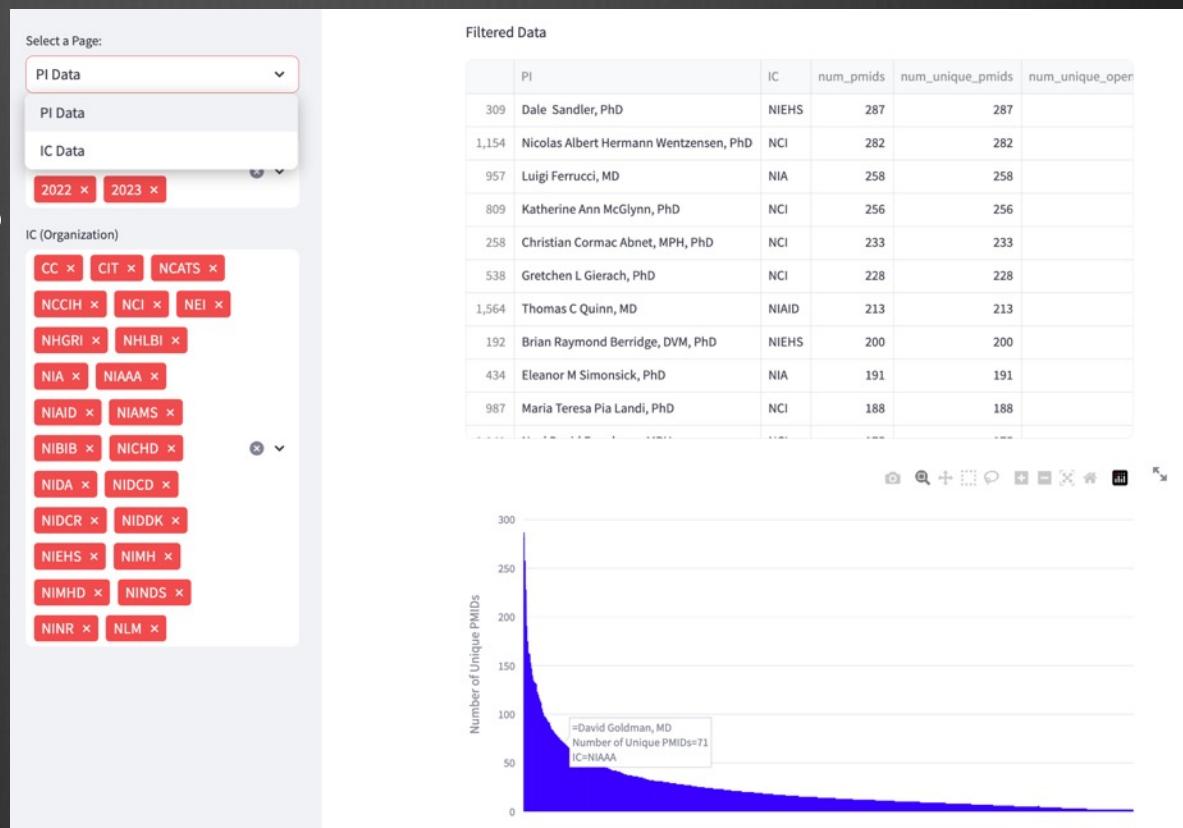
# Measuring Data Sharing



- To reward data sharing there must be a mechanism to measure it
- Text mining tools exist to identify data sharing reported in publications
- However, aggregate information on reported data sharing is not easily accessible to funders and policy makers

## The IRP ShareStats Dashboard

- Allows NIH intramural leadership to quickly identify and reward investigators who have strong records of data sharing



# Similar Efforts

- Riedel et al., 2020 (OddPub)
  - ~ 11,000 pubs from the Charité – U. Berlin
- Serghiou et al., 2021 (rTransparent)
  - ~ 2.75M pubs from PubMed Central Open Access
- Menke et al., 2022 (SciScore Rigor Transparency Index)
  - ~2.15M pubs from PubMed Central Open Access
- Piękniewska et al., 2023 (SciCrunch, RRIDs)
  - ~1.3M pubs PubMed Central Open Access

# What can we add?

- LLM-based metrics
  - Improved accuracy
    - In our analysis, OddPub had ~10% error rate, equally balanced between false positives and false negatives
  - Extracting more specific details from text

# Measuring Data Sharing



## Large Language Models to parse data sharing statements

EM densities and protein models have been deposited in the Electron Microscopy Data Bank and Protein Data Bank for the NaCT-citrate (EMD-22457, 7JSK) and NaCT-PF2 (EMD-22456, 7JSJ) complexes.



Paper DOI	Data Types	Shared Data	Data Repository	Unique ID	Repository URL	Metadata Accessible
10.1038/s41586-021-03230-x	Cryo-EM image processing and model building	TRUE	Electoscopic Data Bank (EMDB)	EMD-22456	<a href="https://www.ebi.ac.uk/emdb/">https://www.ebi.ac.uk/emdb/</a>	TRUE
10.1038/s41586-021-03230-x	Cryo-EM image processing and model building	TRUE	Electoscopic Data Bank (EMDB)	EMD-22457	<a href="https://www.ebi.ac.uk/emdb/">https://www.ebi.ac.uk/emdb/</a>	TRUE
10.1038/s41586-021-03230-x	Protein structure determination, Biochemical assay results, Thermostability assay results	TRUE	Protein Data Bank (PDB)	7JSJ	<a href="https://www.rcsb.org/">https://www.rcsb.org/</a>	TRUE
10.1038/s41586-021-03230-x	Protein structure determination, Biochemical assay results, Thermostability assay results	TRUE	Protein Data Bank (PDB)	7JSK	<a href="https://www.rcsb.org/">https://www.rcsb.org/</a>	TRUE

...

# Using LLMs for complex extraction



## Data Sharing Statement Extraction

Extract data from papers

+ 3 columns added Sort: Most relevant CSV

Paper	Shared Data	links	Data Availability
<input type="checkbox"/> A highly replicable decline in mood during rest and simple tasks David C Jangraw +10 <i>Nature Human Behaviour</i> 2023 6 citations	The shared data for this paper can be downloaded from the Open Science Framework for Online Participants and from Dryad for Mobile App Participants. The data was collected during specific time periods mentioned in the paper.	- Online Participants' data: <a href="https://osf.io/km69z">https://osf.io/km69z</a> - Mobile App Participants' data: <a href="https://doi.org/10.5061/dryad.prr4xgxkk">https://doi.org/10.5061/dryad.prr4xgxkk</a>	All data used in the manuscript have been made publicly available. Online Participants' data can be found on the Open Science Framework at <a href="https://osf.io/km69z">https://osf.io/km69z</a> . Mobile App Participants' data can be found on Dryad at <a href="https://doi.org/10.5061/dryad.prr4xgxkk">https://doi.org/10.5061/dryad.prr4xgxkk</a> (ref. 101).  The code for the task and survey is available on GitLab at <a href="https://gitlab.pavlovia.org/mooddrift">https://gitlab.pavlovia.org/mooddrift</a> . Our data analysis software, as well as the means to create a Python environment that automatically installs it on a user's machine, has been made available online at <a href="https://github.com/djangraw/MoodDrift">https://github.com/djangraw/MoodDrift</a> .

# What can we add?

- LLM-based metrics
  - Improved accuracy
    - In our analysis, OddPub had ~10% error rate, equally balanced between false positives and false negatives
  - Extracting more specific details from pubs
  - More nuanced evaluations

# Many Open Science Rubrics are vague

- FAIR (Wilkinson et al, 2016)

ID	Questions from <a href="https://www.go-fair.org/fair-principles/">https://www.go-fair.org/fair-principles/</a>
F1	(Meta)data are assigned a globally unique and persistent identifier
F2	Data are described with rich metadata (defined by R1 below)
F3	Metadata clearly and explicitly include the identifier of the data they describe
F4	(Meta)data are registered or indexed in a searchable resource
A1	(Meta)data are retrievable by their identifier using a standardised communications protocol
A1-1	The protocol is open, free, and universally implementable
A1-2	The protocol allows for an authentication and authorisation procedure, where necessary
A2	Metadata are accessible, even when the data are no longer available
I1	(Meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation
I2	(Meta)data use vocabularies that follow FAIR principles
I3	(Meta)data include qualified references to other (meta)data
R1	(Meta)data are richly described with a plurality of accurate and relevant attributes
R1-1	(Meta)data are released with a clear and accessible data usage license
R1-2	(Meta)data are associated with detailed provenance
R1-3	(Meta)data meet domain-relevant community standards

# Many Open Science Rubrics are vague

- FAIR (Wilkinson et al., 2016)
- Delphi Consensus (Cobey et al., 2023)

Rank	Practice	Unit of analysis
1	whether clinical trials were registered before they started recruitment	Paper, protocol
2	whether study data were shared openly at the time of publication	paper
3	published open access (time delay?)	paper
4	whether study code was shared openly at the time of publication	paper
5	whether systematic reviews have been registered.	paper
6	were registered clinical trials were reported in the registry within 1 year of study completion	paper
7	whether there was a statement about study materials sharing with publications	paper
8	whether study reporting guideline checklists were used	paper
9	citations to data.	dataset
10	trial results in a manuscript-style publication (peer reviewed or preprint).	paper
11	the number of preprints.	Institution
12	systematic review results in a manuscript-style publication (peer reviewed or preprint).	Systematic review

# Many Open Science Rubrics are vague

- FAIR (Wilkinson et al., 2016)
- Delphi Consensus (Cobey et al., 2023)
- TOSI Open  
Science Principles

## **1. Public release of data and other scientific resources**

The Neuro and its researchers will publish open access articles and render all positive and negative numerical data, models used, data sources, materials, reagents, algorithms, software and other scientific resources publicly available no later than the publication date of the first article that relies on this data or resource.

## **2. External research partnerships**

All publications, data and scientific resources generated through research partnerships –whether with commercial, philanthropic, or public sector actors – are to be released publicly no later than the publication date of the first article that relies on this data or resource.

## **3. Research materials and tools**

The Neuro supports knowledge creation and innovation by maximizing the long-term value of the physical contributions made by research participants and the physical scientific resources created by Neuro researchers and their collaborators. The Neuro will manage these resources in such a way as to remain financially self-sustaining, while continuing to enrich and strengthen its informational content and the knowledge it provides.

In handling materials originating from patients, The Neuro recognizes the primacy of safeguarding the dignity and privacy of patient-participants, and respecting the rights and duties owed them through the informed consent process.

## **4. Intellectual property**

Subject to patient confidentiality and informed consent given, neither The Neuro nor its researchers in their capacity as employees or consultants of McGill or The Neuro will obtain restrictive intellectual property protection in respect of any of their research outputs, whether done internally or with collaborators.

## **5. Researcher and patient autonomy**

The Neuro supports the autonomy of its stakeholders, including but not limited to researchers, staff, trainees and patients, through recognizing their right to decline to participate in research and associated activities under an Open Science framework. However, The Neuro will not support activities that compromise the Open Science principles outlined above.

# Using LLMs for nuanced open science questions

## Can large language models provide useful feedback on research papers? A large-scale empirical analysis.

Weixin Liang<sup>1\*</sup>, Yuhui Zhang<sup>1\*</sup>, Hancheng Cao<sup>1\*</sup>, Binglu Wang<sup>2</sup>, Daisy Yi Ding<sup>3</sup>, Xinyu Yang<sup>4</sup>, Kailas Vodrahalli<sup>5</sup>, Siyu He<sup>3</sup>, Daniel Scott Smith<sup>6</sup>, Yian Yin<sup>4</sup>, Daniel A. McFarland<sup>6</sup>, and James Zou<sup>1,3,5+</sup>

<sup>1</sup>Department of Computer Science, Stanford University, Stanford, CA 94305, USA

<sup>2</sup>Kellogg School of Management, Northwestern University, Evanston, IL 60208, USA

<sup>3</sup>Department of Biomedical Data Science, Stanford University, Stanford, CA 94305, USA

<sup>4</sup>Department of Information Science, Cornell University, Ithaca, NY 14850, USA

<sup>5</sup>Department of Electrical Engineering, Stanford University, Stanford, CA 94305, USA

<sup>6</sup>Graduate School of Education, Stanford University, Stanford, CA 94305, USA

\*Correspondence should be addressed to: jamesz@stanford.edu

+these authors contributed equally to this work

### ABSTRACT

<https://doi.org/10.48550/arXiv.2310.01783>

“... more than half (57.4%) of the users found GPT-4 generated feedback helpful/very helpful and 82.4% found it more beneficial than feedback from at least some human reviewers.”

# Using LLMs for nuanced open science questions

## Data Sharing Statement Extraction

Extract data from papers

+ 3 columns added Sort: Most relevant CSV

Paper	Shared Data	links	Data Availability
<input type="checkbox"/> A highly replicable decline in mood during rest and simple tasks David C Jangraw +10 <i>Nature Human Behaviour</i> 2023 6 citations	The shared data for this paper can be downloaded from the Open Science Framework for Online Participants and from Dryad for Mobile App Participants. The data was collected during specific time periods mentioned in the paper.	- Online Participants' data: <a href="https://osf.io/km69z">https://osf.io/km69z</a> - Mobile App Participants' data: <a href="https://doi.org/10.5061/dryad.prr4xgxkk">https://doi.org/10.5061/dryad.prr4xgxkk</a>	All data used in the manuscript have been made publicly available. Online Participants' data can be found on the Open Science Framework at <a href="https://osf.io/km69z">https://osf.io/km69z</a> . Mobile App Participants' data can be found on Dryad at <a href="https://doi.org/10.5061/dryad.prr4xgxkk">https://doi.org/10.5061/dryad.prr4xgxkk</a> (ref. 101).  The code for the task and survey is available on GitLab at <a href="https://gitlab.pavlovia.org/mooddrift">https://gitlab.pavlovia.org/mooddrift</a> . Our data analysis software, as well as the means to create a Python environment that automatically installs it on a user's machine, has been made available online at <a href="https://github.com/djangraw/MoodDrift">https://github.com/djangraw/MoodDrift</a> .

# What can we add?

- LLM-based metrics
  - Improved accuracy
    - In our analysis, OddPub had ~10% error rate, equally balanced between false positives and false negatives
  - Extracting more specific details from pubs
  - More nuanced evaluations
- An open, method-agnostic, repository of open science indicators

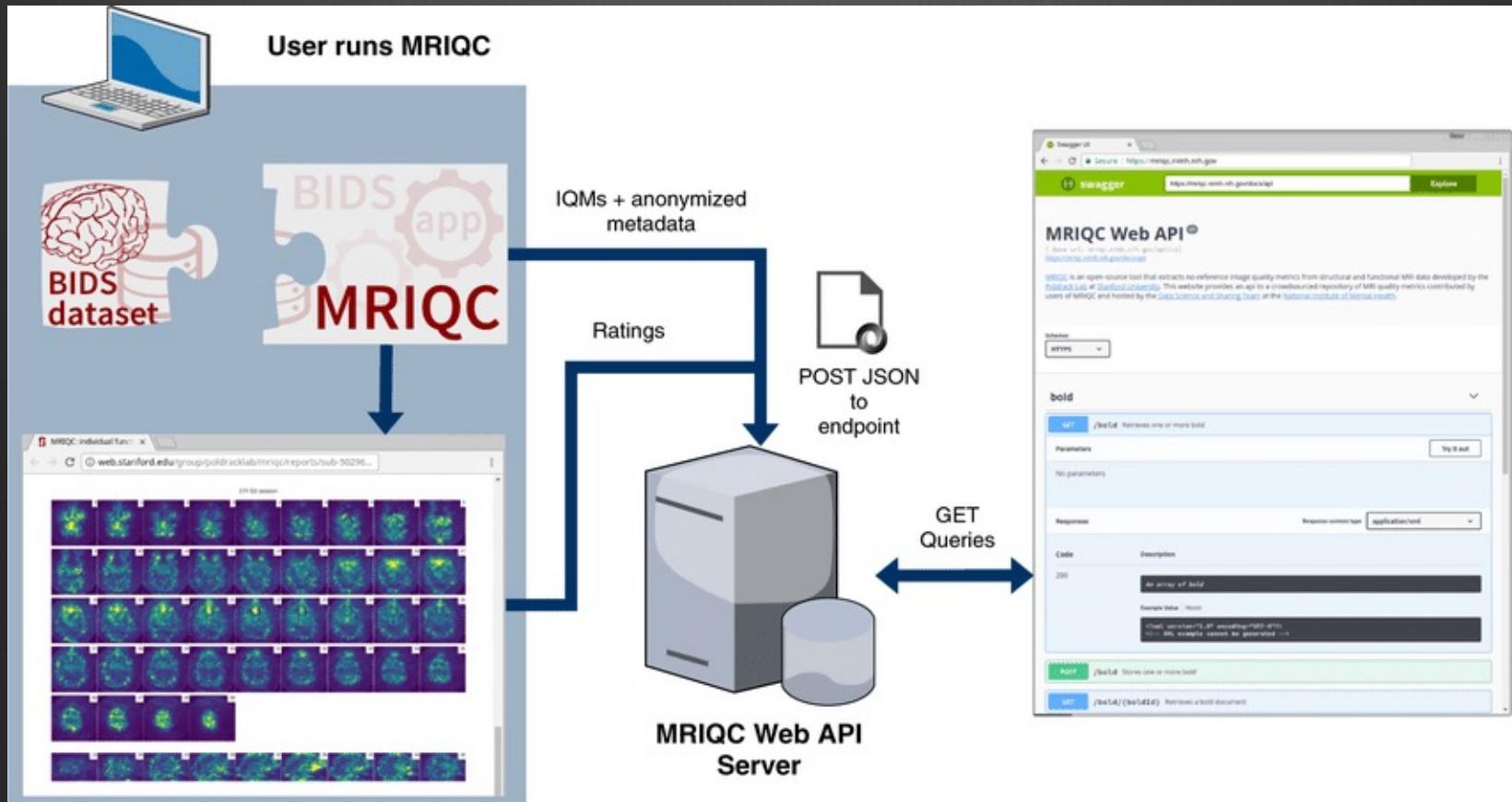
# Lots of automatically and manually label data available

- Riedel et al.
  - ~ 11,000 images
- Serghiou et al.
  - ~ 2.75M images
- Menke et al.
  - ~2.15M images
- Piękniewski et al.
  - ~1.3M images



Lots of disparate,  
unharmonized CSV files

# Modelled on previous success: MRIQC



# Modelled on previous success: MRIQC DB

Stanford University

## MRIQC Web API

- Crowdsourced database of MR QC metrics
- QC metrics from ~375K unique BOLD scans and ~280K T1w scans
- Publicly available:
  - <https://mriqc.nimh.nih.gov/>

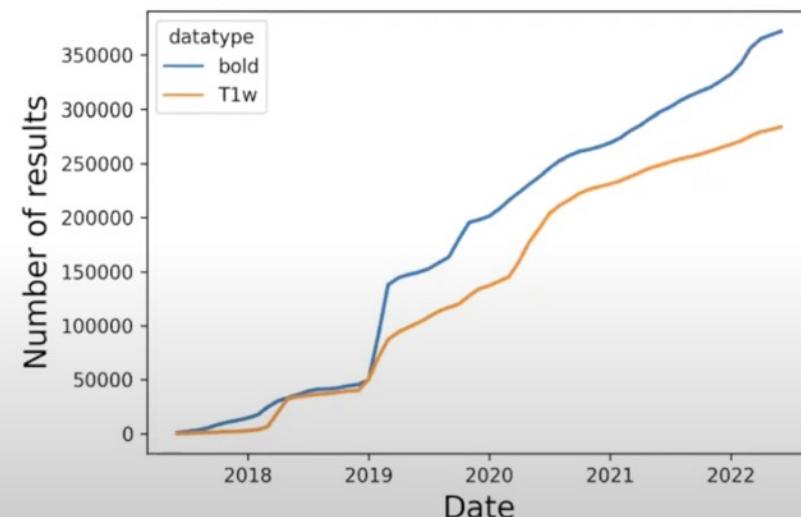
## SCIENTIFIC DATA

OPEN  
DATA DESCRIPTOR

Received: 19 September 2018  
Accepted: 12 March 2019

Crowdsourced MRI quality metrics and expert quality annotations for training of humans and machines

Oscar Esteban<sup>1</sup>, Ross W. Blair<sup>1</sup>, Dylan M. Nielson<sup>2</sup>, Jan C. Varada<sup>3</sup>, Sean Marrett<sup>3</sup>, Adam G. Thomas<sup>2</sup>, Russell A. Poldrack<sup>1</sup> & Krzysztof J. Gorgolewski<sup>1</sup>



# Modelled on previous success: MRIQC DB

Stanford University

## MRIQC Web API

- Crowdsourced database of MR QC metrics
- QC metrics from ~375K unique BOLD scans and ~280K T1w scans
- Publicly available:
  - <https://mriqc.nimh.nih.gov/>

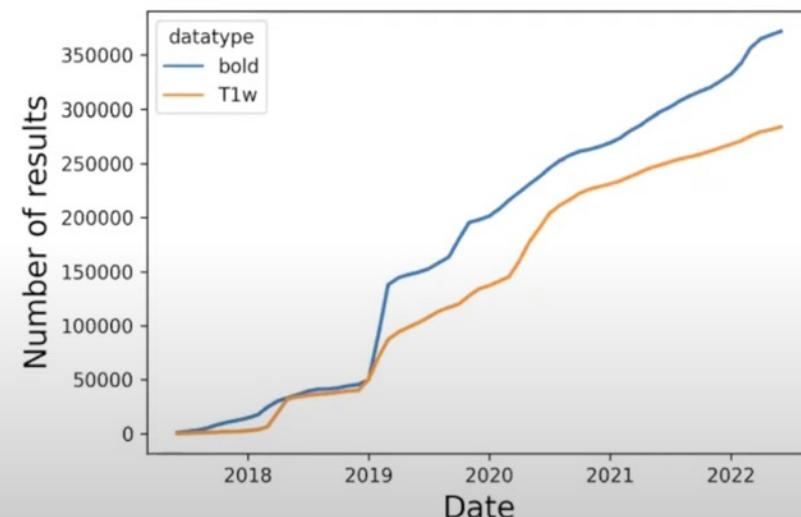
## SCIENTIFIC DATA

OPEN  
DATA DESCRIPTOR

Received: 19 September 2018  
Accepted: 12 March 2019

Crowdsourced MRI quality metrics and expert quality annotations for training of humans and machines

Oscar Esteban<sup>1</sup>, Ross W. Blair<sup>1</sup>, Dylan M. Nielson<sup>2</sup>, Jan C. Varada<sup>3</sup>, Sean Marrett<sup>3</sup>, Adam G. Thomas<sup>2</sup>, Russell A. Poldrack<sup>1</sup> & Krzysztof J. Gorgolewski<sup>1</sup>



# Work In Progress

The screenshot shows a GitHub repository page for 'osm\_cli'. The repository is owned by 'nimh-dsst'. The main tab is 'Code'. The repository is public. It has 1 branch and 0 tags. There are 6 commits from 'agt24' made last week. The commits are:

- add xml outputs from sciencebeam-parser container (a7e4549 · last week)
- example\_pdf\_inputs (add link to spec and example input PDF · last week)
- xmls\_sciencebeam (add xml outputs from sciencebeam-parser container · last week)
- LICENSE (Initial commit · last week)
- README.md (Initial commit · last week)
- draft\_of\_specification\_2024-05-06.md (Update draft\_of\_specification\_2024-05-06.md · last week)

GitStart



# Code as a Service

Assign tickets, get ***high-quality production code*** powered by AI agents and our developer community.

## OpenSciMetrics

OpenSciMetrics (OSM) applies NLP and LLM-based metrics and indicators related to transparency, data sharing, rigor, and open science on biomedical publications.

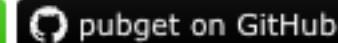
# Synergies on Extraction & Manual Labelling

## Mining the neuroimaging literature.

Jérôme Dockès<sup>1\*</sup>, Kendra Oudyk<sup>2\*</sup>, Mohammad Torabi<sup>2</sup>, Alejandro I de la Vega<sup>3</sup>, and Jean-Baptiste Poline<sup>2</sup>

**pubget**



 build passing  codecov 100%  pubget on GitHub

**labelbuddy**



# What can we add?

- LLM-based metrics
  - Improved accuracy
    - In our analysis, OddPub had ~10% error rate, equally balanced between false positives and false negatives
  - Extracting more specific details from pubs
  - More nuanced evaluations
- An open, method-agnostic, repository of open science indicators
- A new “more bang-for-your-buck” metric

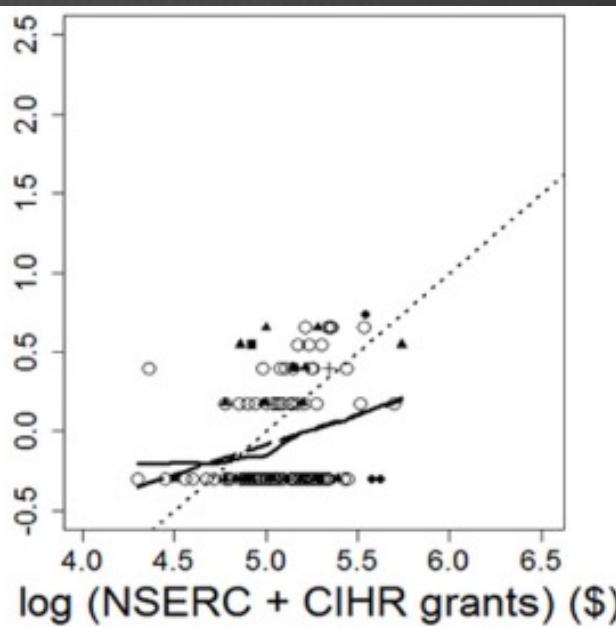
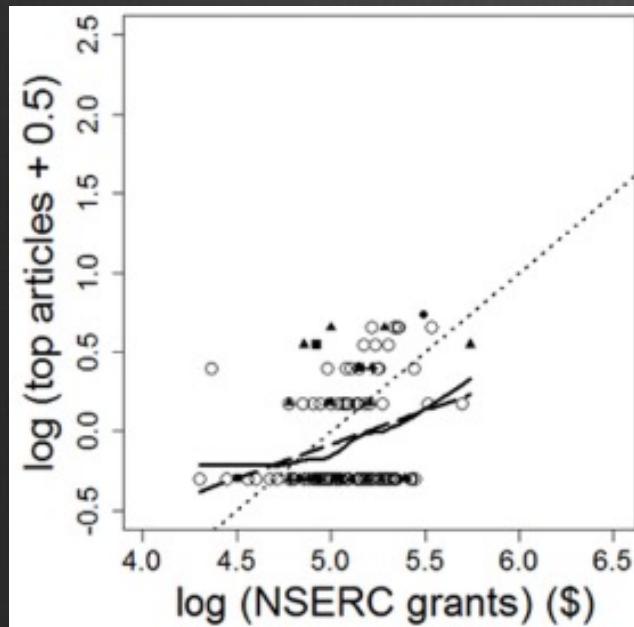
# Correlating Funding Levels and Open Science Practices

## Big Science vs. Little Science: How Scientific Impact Scales with Funding

Jean-Michel Fortin, David J. Currie\*

Ottawa-Carleton Institute of Biology, University of Ottawa, Ottawa, Ontario, Canada

June 2013 | Volume 8 | Issue 6 | e65263



# Correlating Funding Levels and Open Science Practices



- Identify investigators with relatively low levels of funding, but relatively strong open science practices
-

# Thank You!

## Tanenbaum Open Science Institute (TOSI)



Guy Rouleau

Director, The Neuro  
Co-Founder of the Tanenbaum Open Science Institute

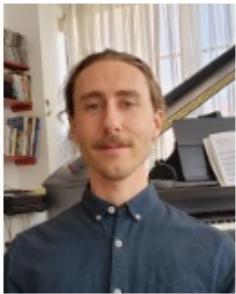


Annabel Seyller

Chief of Staff, The Neuro  
CEO, Tanenbaum Open Science Institute



**Jean-Baptiste Poline**



Gabriel Pelletier

Open Science Data Manager  
Interim Open Science Alliance Officer



Luisa Pimentel

Open Science Community Officer

# 2023 Data Science and Sharing Team



Adam Thomas  
**Team Lead**



Dustin Moraczewski  
**Data Scientist**



Eric Earl  
**Data Scientist**



Mia Zwally  
**Post-bac IRTA**



Arshitha Basavaraj  
**Data Engineer**



Jessica Dafflon  
**Data Scientist**

<http://cmn.nimh.nih.gov/dsst>

# Feedback? Where to focus?

- LLM-based metrics
  - Improved accuracy
- In our analysis, OddPub had ~10% error rate, equally balanced between false positives and false negatives
  - Extracting more specific details from pubs
  - More nuanced evaluations
- An open, method-agnostic, repository of open science indicators
- A new “more bang-for-your-buck” metric

END