# Data and Software:
# Important Building Blocks of Research

## 12 October 2021

**Shelley Stall**
**Senior Director, Data Leadership**
**American Geophysical Union**
**0000-0003-2926-8353**
**@ShelleyStall | sstall@agu.org**

AGU
ADVANCING EARTH
AND SPACE SCIENCE

# AGU's position statement on data affirms that

**"Earth and space science data are a world heritage, and an essential part of the science ecosystem. All players in the science ecosystem—researchers, repositories, publishers, funders, institutions, etc.— should work to ensure that relevant scientific evidence is processed, shared, and used ethically, and is available, preserved, documented, and fairly credited."**

https://www.agu.org/Share-and-Advocate/Share/Policymakers/Position-Statements/Position_Data

# AGU's position statement on data affirms that

"**Earth and space science data are a world heritage, and an essential part of the science ecosystem. All players in the science ecosystem—researchers, repositories, publishers, funders, institutions, etc.—should work to ensure that relevant scientific evidence is processed, shared, and used ethically, and is available, preserved, documented, and fairly credited.**"

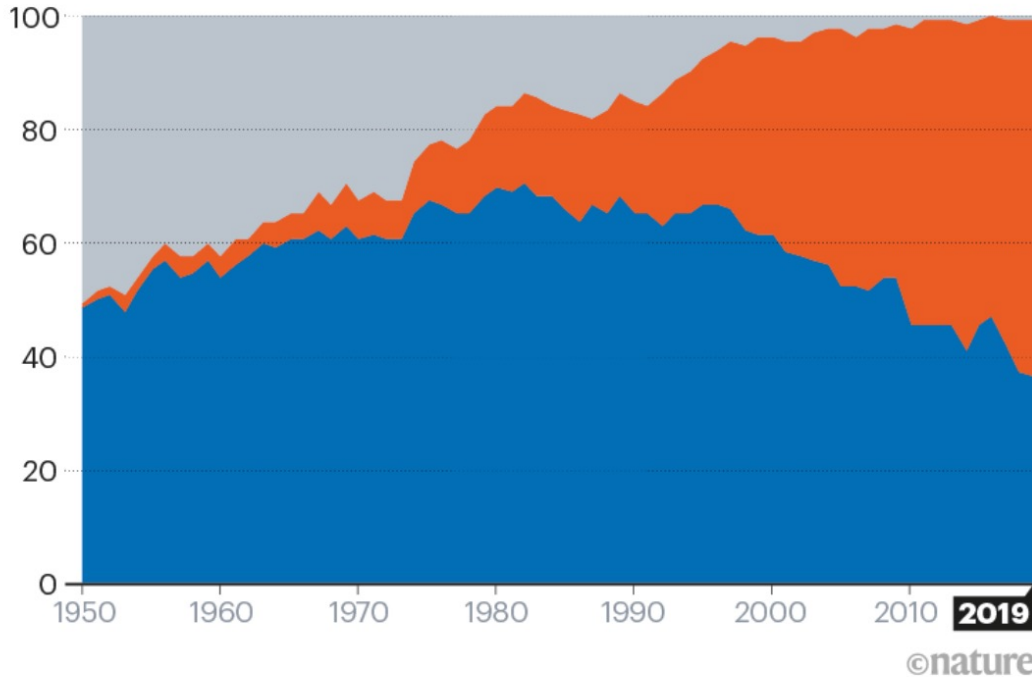# 150 years of Nature

A century and a half of research and discovery.

nature

# INTERNATIONAL COLLABORATIONS

Author lists on research publications show a shift towards multinational teams; fewer teams are composed entirely of researchers from one country.

**Proportion of papers**

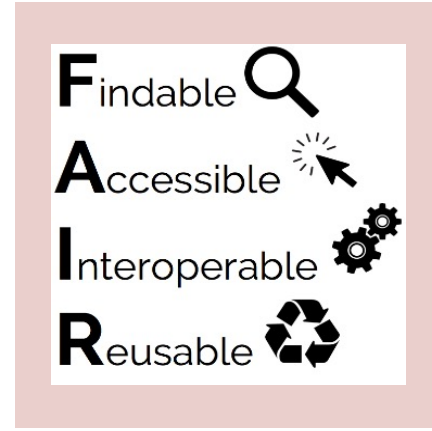■ Multinational  ■ Domestic  ■ Single author



Design: Jasiek Krzysztofiak; Data analysis: Digital Science Consultancy

Monastersky, R., & Van Noorden, R. (2019). 150 years of Nature: a data graphic charts our evolution. Nature, 575(7781), 22–23. https://doi.org/10.1038/d41586-019-03305-w

# The Future of your Research

- Research Teams (not individuals)
- International Collaborations

- Robust **tools to discover** relevant research worldwide
- Good **documentation** to understand that research, data, and/or software

- Data that is **interoperable**, no matter which research team created it
- Software that is **accessible** and developed in current tools (e.g., Jupyter Notebooks)

- Licenses that support **reuse**.

OPEN

FAIR

# Closed vs Open – A continuum

# FAIR Guiding Principles

## FAIR is…
### Findable
### Accessible
### Interoperable
### Reusable

Article in Nature journal *Scientific Data*: Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and stewardship. *Sci. Data* 3:160018 doi: 10.1038/sdata.2016.18 (2016).

# FAIR Data Principles (applies to software and all digital objects)

- **Findable**
  - Assign persistent IDs (PIDs), provide rich metadata, register in a searchable resource, …
- **Accessible**
  - Retrievable by their ID using a standard protocol, metadata remain accessible even when data are no longer available…
- **Interoperable**
  - Use formal, broadly applicable languages, use standard vocabularies, qualified references…
- **Reusable**
  - Rich, accurate metadata, clear licenses, provenance, use of community standards…

Repository

Repository

Community

Community / Repository

# Is FAIR Open?   In short, "It depends."



Data can be FAIR or Open, both or neither.

The greatest potential for **reuse** comes when data are **both FAIR and Open.**

Higman, Rosie, Daniel Bangert, and Sarah Jones. 2019. "Three Camps, One Destination: The Intersections of Research Data Management, FAIR and Open". *Insights* 32 (1): 18. DOI: http://doi.org/10.1629/uksg.468

**Data should be as open as possible, as closed as necessary.**

# What is the Continuum for FAIR?   Two examples...

**Persistent Identifier**    No unique ID    Globally resolvable persistent identifier

**Robust Metadata**    No metadata    Metadata that aligns to community recommended standards/vocabulary/guidelines; Machine Readable

# Why do we Care about FAIR?

# The Future of your Research

- Research Teams (not individuals)
- International Collaborations

- Robust **tools to discover** relevant research worldwide
- Good **documentation** to understand that research, data, and/or software

- Data that is **interoperable**, no matter which research team created it
- Software that is **accessible** and developed in current tools (e.g., Jupyter Notebooks)
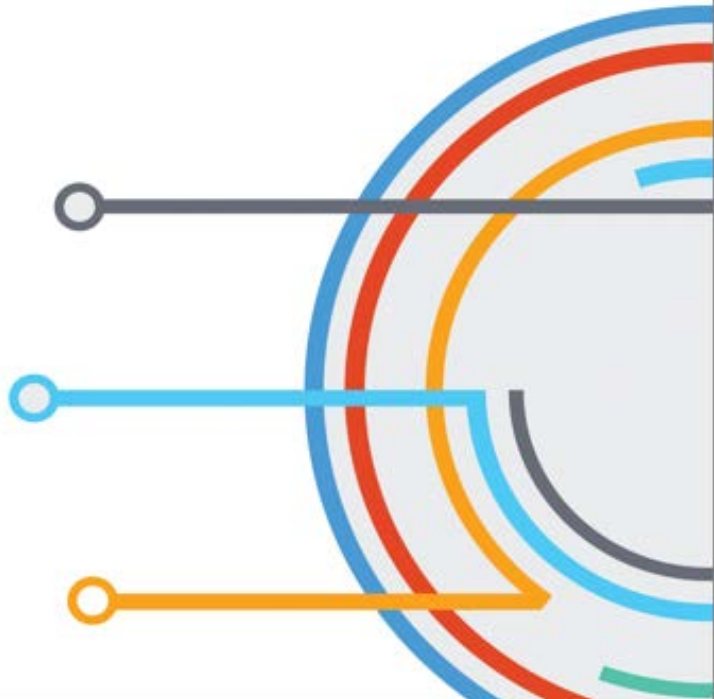
- Licenses that support **reuse**.

F-A

F-A-R

F-A-I-R

R

*Open Science by Design* is aimed at overcoming barriers and **moving toward open science as the default approach across the research enterprise**. This report explores specific examples of open science and discusses a range of challenges, focusing on stakeholder perspectives. It is meant to provide guidance to the research enterprise and its stakeholders as they build strategies for achieving open science and take the next steps.

Report released July 2018

# NASEM Definition of Open Science

Open science aims to ensure the open availability and usability:

- scholarly publications
- the data that result from scholarly research
- and the methodologies, including code or algorithms, that were used to generate those data.

# 5 Findings and Recommendations

1. Building a Supportive Culture - internationally

2. Training for Open Science by Design

3. Ensuring long term preservation and stewardship - Developing and sustaining the infrastructure

4. Facilitating data discovery, and reuse – FAIR Data Principles

5. Developing new approaches to fostering Open Science by Design – with public and private funder support
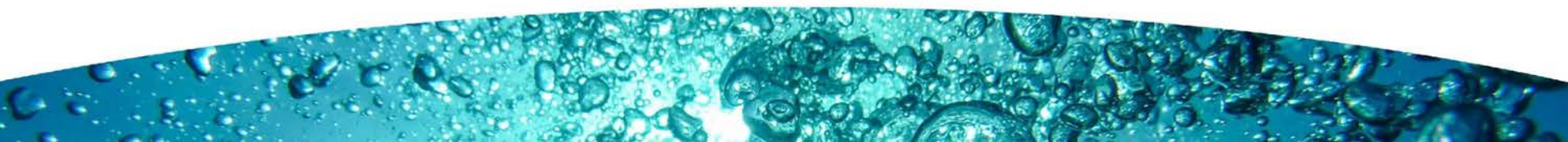
# Developing a Toolkit for Fostering Open Science Practices

## Released 30 Sept 2021 from The National Academies Roundtable on Aligning Incentives for Open Science

National Academies of Sciences, Engineering, and Medicine. 2021. *Developing a Toolkit for Fostering Open Science Practices: Proceedings of a Workshop*. Washington, DC: The National Academies Press. https://doi.org/10.17226/26308.

PROCEEDINGS OF A WORKSHOP

Developing a Toolkit for Fostering
OPEN SCIENCE PRACTICES

The National Academies of
SCIENCES · ENGINEERING · MEDICINE

# Open Science – The Key Pillars

- Open scientific knowledge
  - **Scientific publications**
  - **Open research data**
  - **Open source software and source code**
  - **Open hardware**
- Open science infrastructures
- Science communication
- Open engagement of societal actors
- Open dialogue with other knowledge systems.

Source: Final report on the draft text of the UNESCO Recommendation on Open Science; final adoption expected November 2021

https://en.unesco.org/science-sustainable-future/open-science

From the very beginning of the research process,

the researcher **both contributes** to open science and

**takes advantage** of the open science practices of other members of the research community.

# NSF Public Access Repository (PAR) 1.0 Current System Functions and Characteristics

**PAR 1.0 focused on Peer-Reviewed Articles**

- Enables researchers to enter metadata for peer-reviewed articles or auto-populate by means of Digital Object Identifiers (DOI)
- Metadata recorded in PAR is also transmitted and synchronized with Award Search database
- The metadata for articles can be searched and displayed
- Researchers may deposit (or retrieve) public access articles as PDF/A files or point to download locations through DOI
- Note that PAR is not a single system, but rather an infrastructure of various interacting software systems in several different parts of the NSF infrastructure, also making use of modified modules from the DOE Office of Scientific and Technical Information (OSTI) system

**URL: http://par.nsf.gov**

Slide Credit: Martin Halbert, NSF Senior Advisor for Public Access, 19 November 2020

# NSF PAR 2.0 Development Planning

- Primary goal for the second version of PAR is to accommodate submissions of research data sets

- Will not require that data sets be deposited in PAR 2.0 but instead will expect data sets to be deposited in repositories which demonstrate FAIR best practices such as maintaining sustainable long-term access through the assignment of a persistent identifier (PID) with quality descriptive metadata

- Now planning the upgrades to the NSF PAR infrastructure to address the additional needed features

- Aim to accomplish major development efforts in calendar 2021

Slide Credit: Martin Halbert, NSF Senior Advisor for Public Access, 19 November 2020

# Final NIH Policy for Data Management and Sharing (NOT-OD-21-013)

Release Date: **October 29, 2020 |** Effective Date: **January 25, 2023**

NIH requires researchers to prospectively plan for how scientific data will be preserved and shared through submission of a Data Management and Sharing Plan

Submission of a Data Management and Sharing Plan outlining how scientific data and any accompanying metadata will be managed and shared, taking into account any potential restrictions or limitations.

Plan is part of the budget Justification section of the application for extramural awards and as part of the technical evaluation for contracts

The DMS Policy applies to all research, funded or conducted in whole or in part by NIH, that results in the generation of scientific data. This includes research funded or conducted by extramural grants, contracts, Intramural Research Projects, or other funding agreements regardless of NIH funding level or funding mechanism. **The DMS Policy does not apply to research and other activities that do not generate scientific data, including training, infrastructure development, and non-research activities.**

# Final NIH Policy for Data Management and Sharing

**NOT-OD-21-015: Allowable Costs for Data Management and Sharing** include curation, data management and sharing data through repositories

- Shared scientific data should be made accessible as soon as possible, and no later than the time of an associated publication, or the end of performance period, whichever comes first.

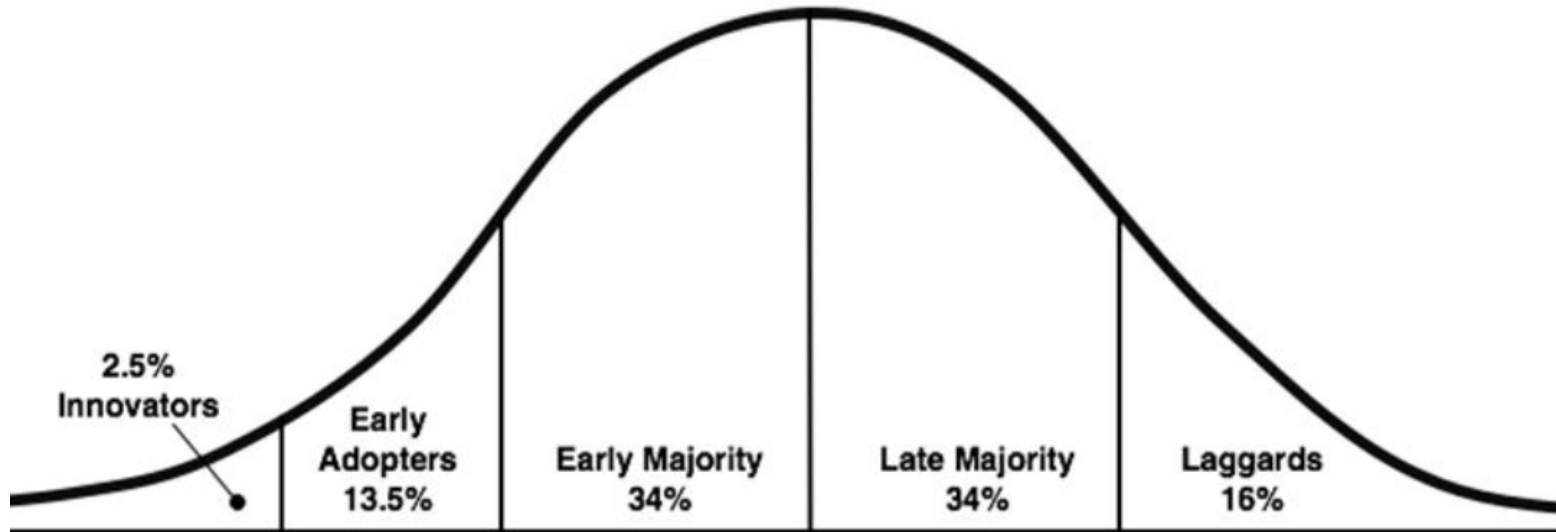**NOT-OD-21-016: Selecting a Repository** include desirable characteristics for a data sharing repository

- **NIH strongly encourages the use of established repositories** to the extent possible for preserving and sharing scientific data.

# What is the Funder Focus?

1. Data Management Plans
   – PI required to consider **how** data will managed and preserved (i.e., repository selection)
2. Linking Data (and software) to a grant, publication, researcher, etc
   – Persistent Identifiers (**PIDs**) and linked infrastructure
3. Culture change
   – Evolving from recommendations to requirements

# Changing a Research Culture



2.5% Innovators

Early Adopters 13.5%

Early Majority 34%

Late Majority 34%

Laggards 16%

Source: Everet Rogers *Diffusion of Innovations model*

Diffusion of Innovations; Rogers, 1963

Laggards
16%

Late Majority
34%

Early Majority
34%

Early Adopters
13.5%

2.5% Innovators

Adapted from Brian Nosek (COS) by Marcia McNutt, President of the National Academy of Science

Make it required – with policies

Make it rewarding – with incentives

Make it accepted – by building communities

Community of Practice:  Promote data sharing, metadata standards, criteria for interoperability

Make it easy – with user interface

Make it possible – with infrastructure

# AGU Data & Software Sharing Guidance

What is covered:

- What data needs to be available?
- Repository Selection
- Availability Statement
- Data & Software Citation
- Physical Samples

# Thank you

**Shelley Stall**
**AGU Sr. Director, Data Leadership**
**sstall@agu.org**
**@ShelleyStall**
**https://orcid.org/0000-0003-2926-8353**