

# Introduction to Research Data Management and Open Science (aka Research)

S. Venkataraman, PhD  
[sv1uk@netscape.net](mailto:sv1uk@netscape.net)

*12<sup>th</sup> & 13<sup>th</sup> August 2024, Trieste*



# Agenda

<b>Day 1 (12<sup>th</sup> August)</b>	
14:00	Introduction to research data management (RDM)
15:00	Exercise: Practical session on RDM
15:30	Introduction to Open Science (Research)
16:00	<b>Break</b>
16:30	Introduction to Open Science (Research) (cont'd)
17:00	Exercise: Open Science
18:00	<b>End Day 1</b>
<b>Day 2 (13<sup>th</sup> August)</b>	
08:30	Introduction to DMPs
09:30	<b>End</b>



# Learning outcomes

---

- Be familiar with the curation lifecycle.
  - Understand the standardisation methods and principles available to add value to your data.
  - Learn about resources to aid your workflows.
  - Increase/encourage your level of openness.
  - Learn about data management plans and the value in implementing them.
-

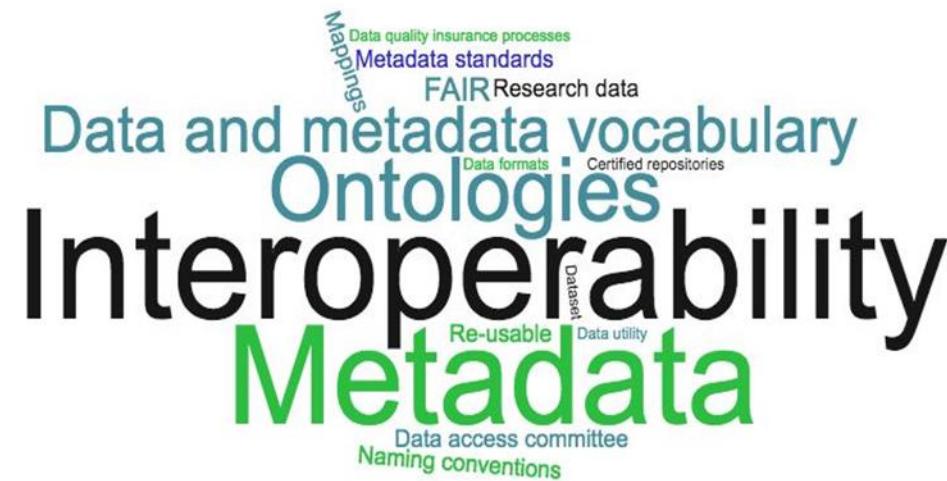


## Language is a barrier...

Respondents mentioned 40 terms which were unclear to them in European Commission DMP:

*“Researchers are not familiar with the following terms/phrases : Metadata, standards for metadata/data, ontologies, mapping with ontologies, interoperability, . . . All the ICT jargon”*

*“With the help from Swedish National Data Service we could clarify many questions. Without this help we would not be able to finish the DMP.”*

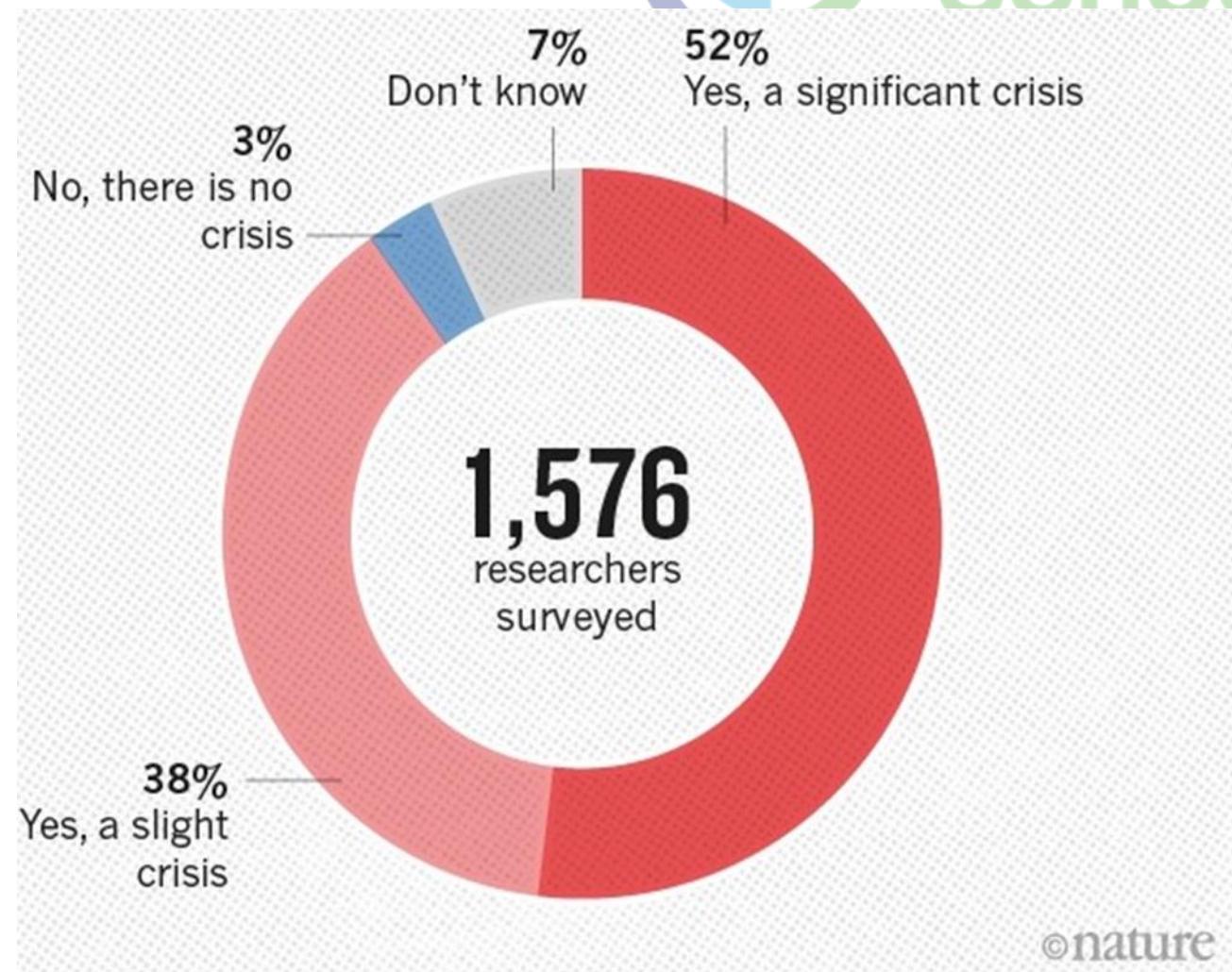


# Is there a reproducibility crisis?

Baker, M. "1,500 scientists lift the lid on reproducibility" *Nature* 533: 452-454 (2016).

<http://www.nature.com/news/1-500-scientists-lift-the-lid-on-reproducibility-1.19970>

doi:10.1038/533452a



# The “data deluge”

- The volume of data is growing exponentially with >90% of all data in the world having been generated in just the last few years.
  - **How to safeguard for the future?**
    - Good RDM is essential!
  - **And what about the environmental impact??**



Image by Pete from [Pixabay](#)



# The wider context

Set of goals outlined by the United Nations

## SUSTAINABLE DEVELOPMENT GOALS



Developed in collaboration with TROLLBÄCK+COMPANY | [TheGlobalGoals@trollback.com](mailto:TheGlobalGoals@trollback.com) | +1.212.529.1010  
For queries on usage, contact: [dpicampaigns@un.org](mailto:dpicampaigns@un.org)

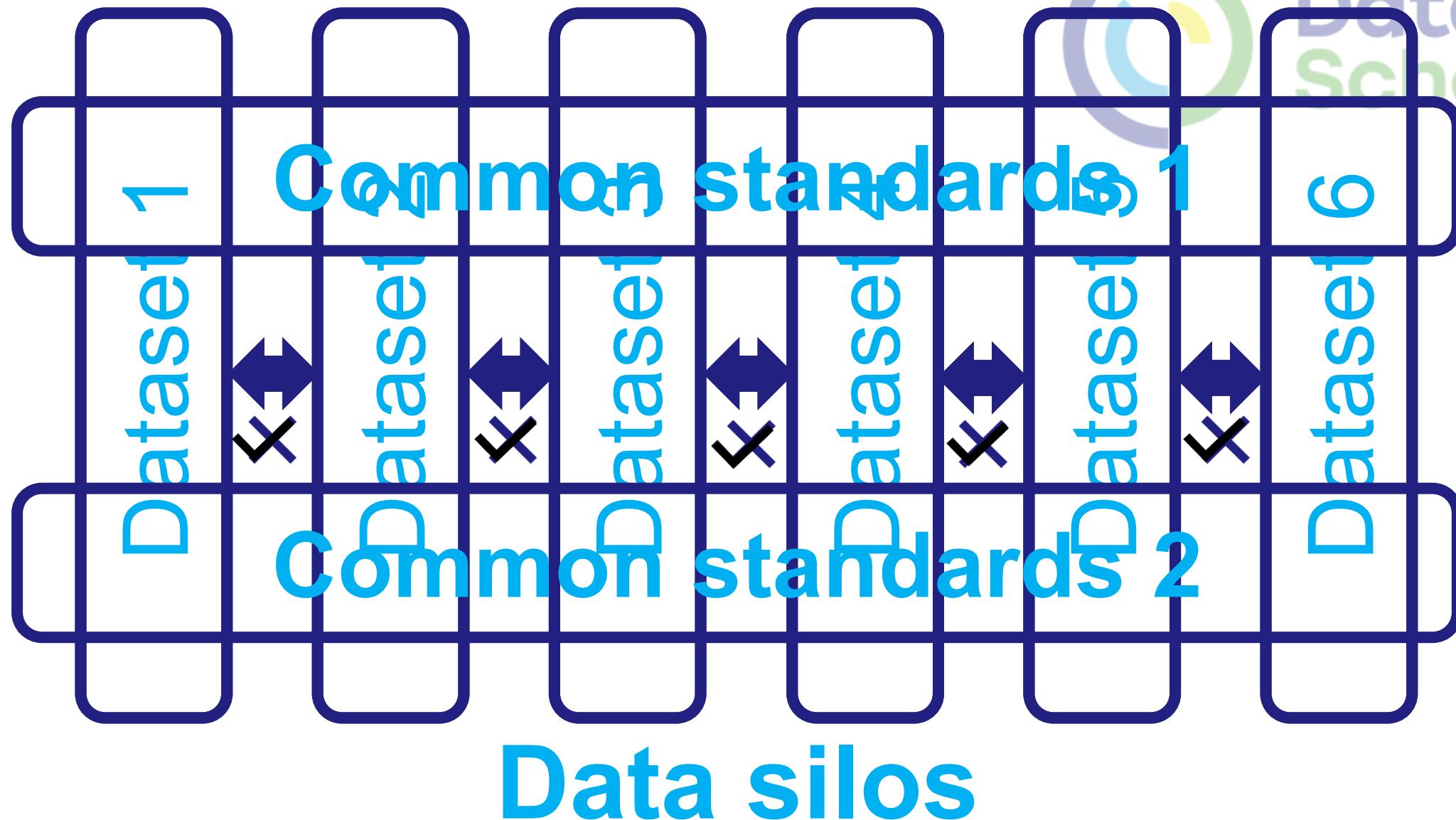


Data  
Schools

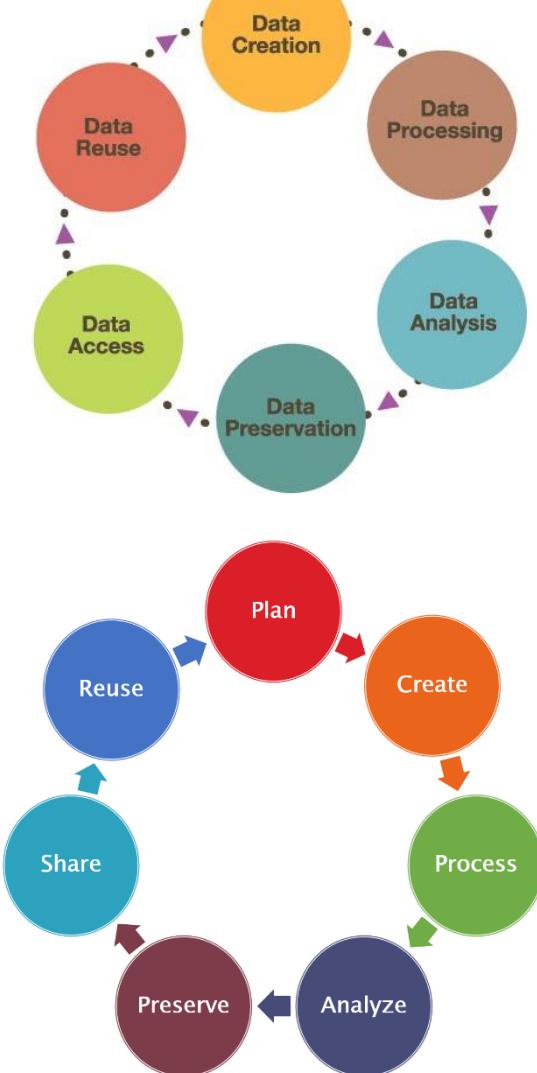
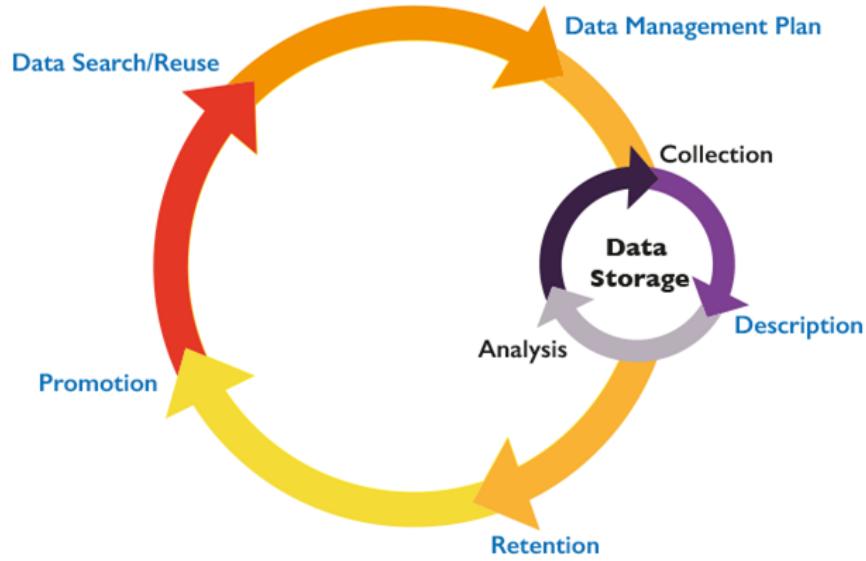
# RDM & the Data Lifecycle



Data  
Schools



# RDM lifecycles



# The curation lifecycle

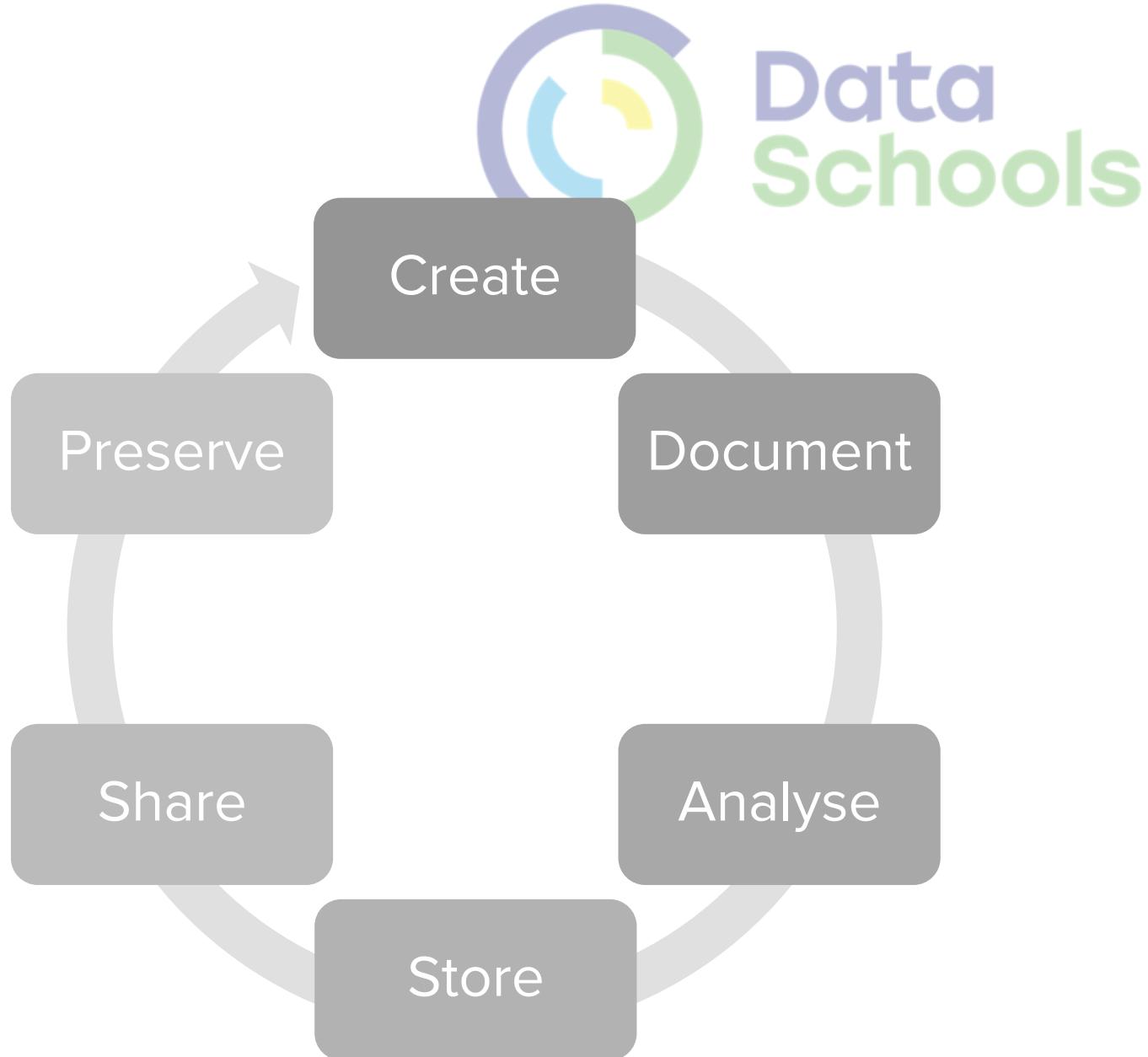
---



# What is Research Data Management?

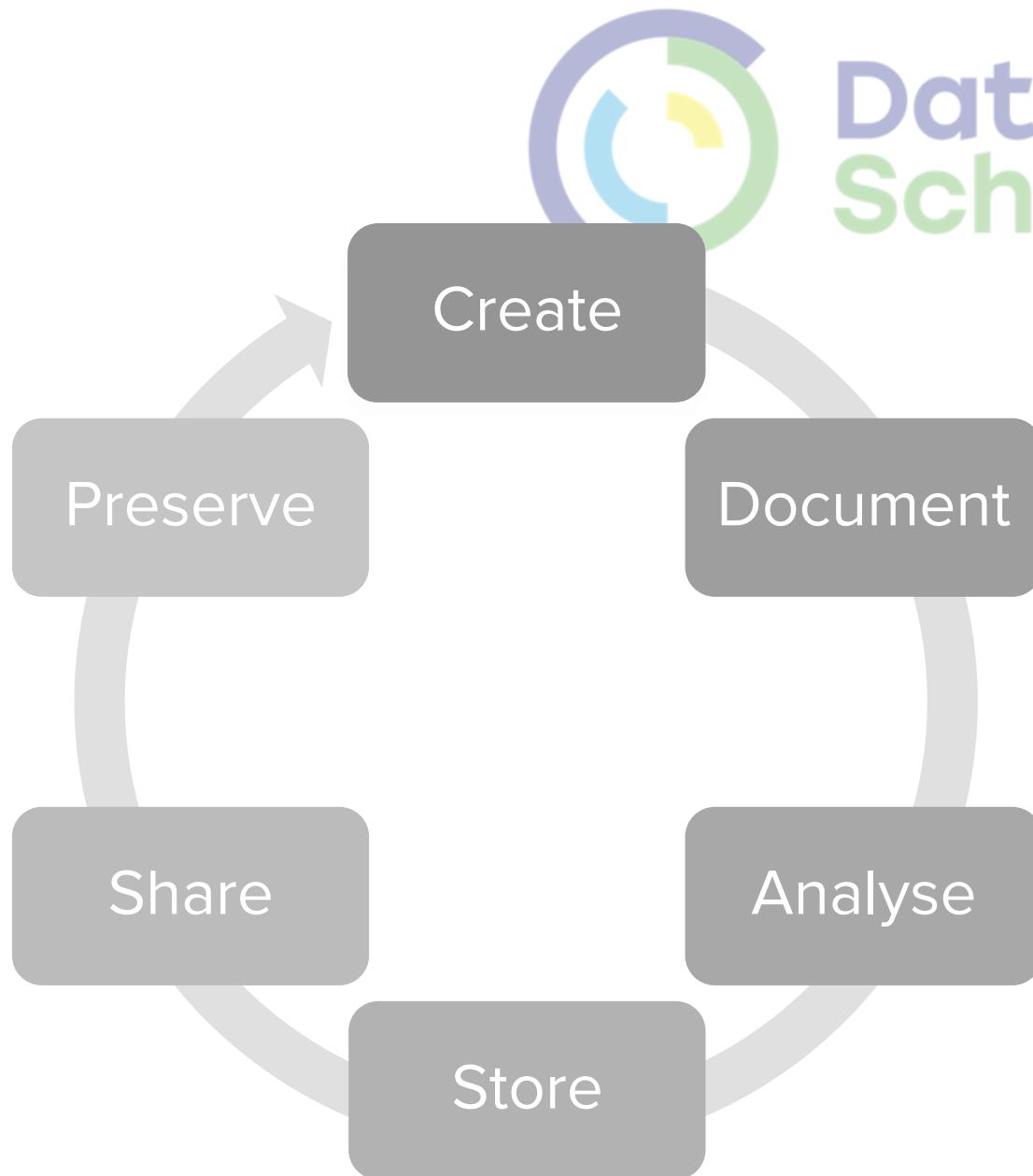
“the active management and appraisal of data over the lifecycle of scholarly and scientific interest”

Data management is part of good research practice.



## Data creation tips

- Ensure consent forms, licences and agreements don't restrict opportunities to share data.
- Choose appropriate formats.
- Adopt a file naming convention.
- Create metadata and documentation as you go.





## Ask for consent for data sharing

If not, data centres won't be able to accept the data – regardless of any conditions on the original grant.

### SAMPLE CONSENT STATEMENT FOR QUANTITATIVE SURVEYS

Thank you very much for agreeing to participate in this survey.

The information provided by you in this questionnaire will be used for research purposes. It will not be used in any manner which would allow identification of your individual responses.

Anonymised research data will be archived at ..... in order to make them available to other researchers in line with current data sharing practices.

# Choose appropriate file formats

---

- Different formats are good for different things.
  - *open, lossless* formats are more sustainable e.g. rtf, xml, tif, wav.
  - proprietary and/or compressed formats are less preservable but are often in widespread use e.g. doc, jpg, mp3.
  - One format for analysis then convert to a standard format.
  - Data centres may suggest preferred formats for deposit.
-

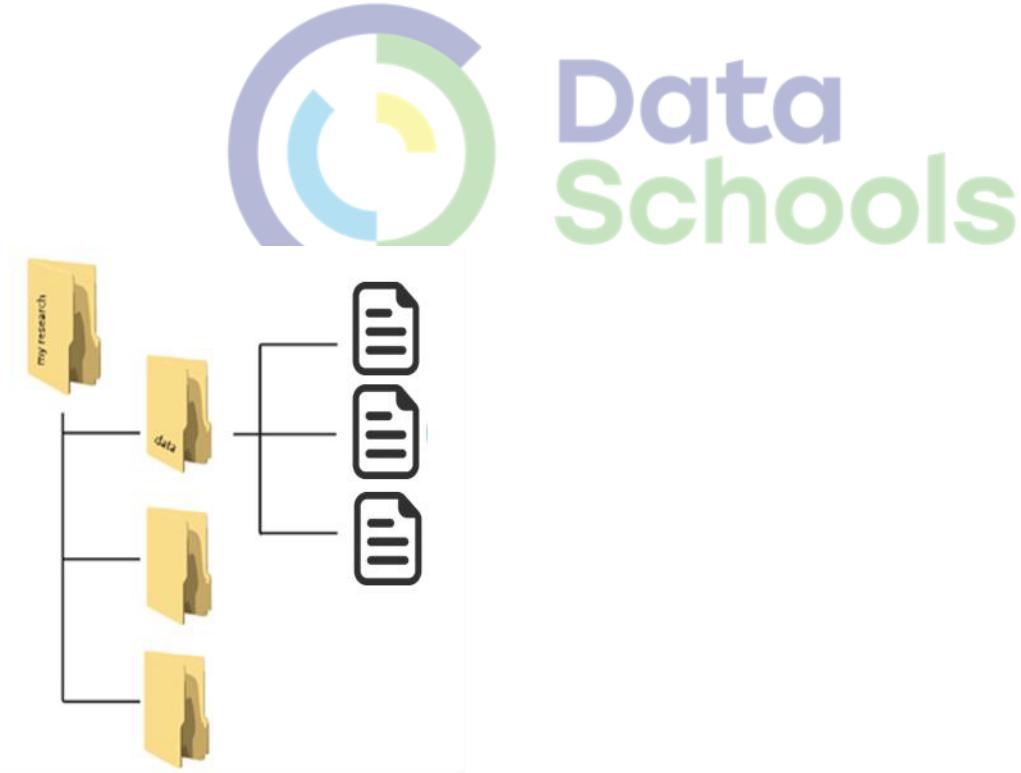


S

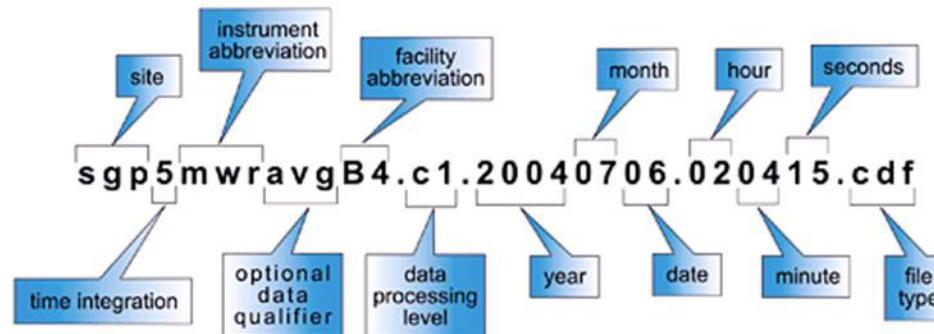
Type of data	Recommended formats	Acceptable formats
Tabular data with extensive metadata variable labels, code labels, and defined missing values	SPSS portable format (.por) delimited text and command ('setup') file (SPSS, Stata, SAS, etc.) structured text or mark-up file of metadata information, e.g. DDI XML file	proprietary formats of statistical packages: SPSS (.sav), Stata (.dta), MS Access (.mdb/.accdb)
Tabular data with minimal metadata column headings, variable names	comma-separated values (.csv) tab-delimited file (.tab) delimited text with SQL data definition statements	delimited text (.txt) with characters not present in data used as delimiters widely-used formats: MS Excel (.xls/.xlsx), MS Access (.mdb/.accdb), dBase (.dbf), OpenDocument Spreadsheet (.ods)
Geospatial data vector and raster data	ESRI Shapefile (.shp, .shx, .dbf, .prj, .sbx, .sbn optional) geo-referenced TIFF (.tif, .tfw) CAD data (.dwg) tabular GIS attribute data Geography Markup Language (.gml)	ESRI Geodatabase format (.mdb) MapInfo Interchange Format (.mif) for vector data Keyhole Mark-up Language (.kml) Adobe Illustrator (.ai), CAD data (.dxf or .svg) binary formats of GIS and CAD packages
Textual data	Rich Text Format (.rtf) plain text, ASCII (.txt) eXtensible Mark-up Language (.xml) text according to an appropriate Document Type Definition (DTD) or schema	Hypertext Mark-up Language (.html) widely-used formats: MS Word (.doc/.docx) some software-specific formats: NUD*IST, NVivo and ATLAS.ti
Image data	TIFF 6.0 uncompressed (.tif)	JPEG (.jpeg, .jpg, .jp2) if original created in this format GIF (.gif) TIFF other versions (.tif, .tiff) RAW image format (.raw) Photoshop files (.psd) BMP (.bmp) PNG (.png) Adobe Portable Document Format (PDF/A, PDF) (.pdf)
Audio data	Free Lossless Audio Codec (FLAC) (.flac)	MPEG-1 Audio Layer 3 (.mp3) if original created in this format Audio Interchange File Format (.aif) Waveform Audio Format (.wav)
Video data	MPEG-4 (.mp4) OGG video (.ogv, .ogg) motion JPEG 2000 (.mj2)	AVCHD video (.avchd)
Documentation and scripts	Rich Text Format (.rtf) PDF/UA, PDF/A or PDF (.pdf) XHTML or HTML (.xhtml, .htm) OpenDocument Text (.odt)	plain text (.txt) widely-used formats: MS Word (.doc/.docx), MS Excel (.xls/.xlsx) XML marked-up text (.xml) according to an appropriate DTD or schema, e.g. XHMTL 1.0

# How will you organise your data?

- Keep file and folder names short, but meaningful.
- Agree a method for versioning.
- Include dates in a set format e.g. YYYYMMDD.
- Avoid using non-alphanumeric characters in file names.
- Use hyphens or underscores not spaces e.g. day-sheet, day sheet.
- Order the elements in the most appropriate way to retrieve the record.



An example netCDF data file name is depicted below:

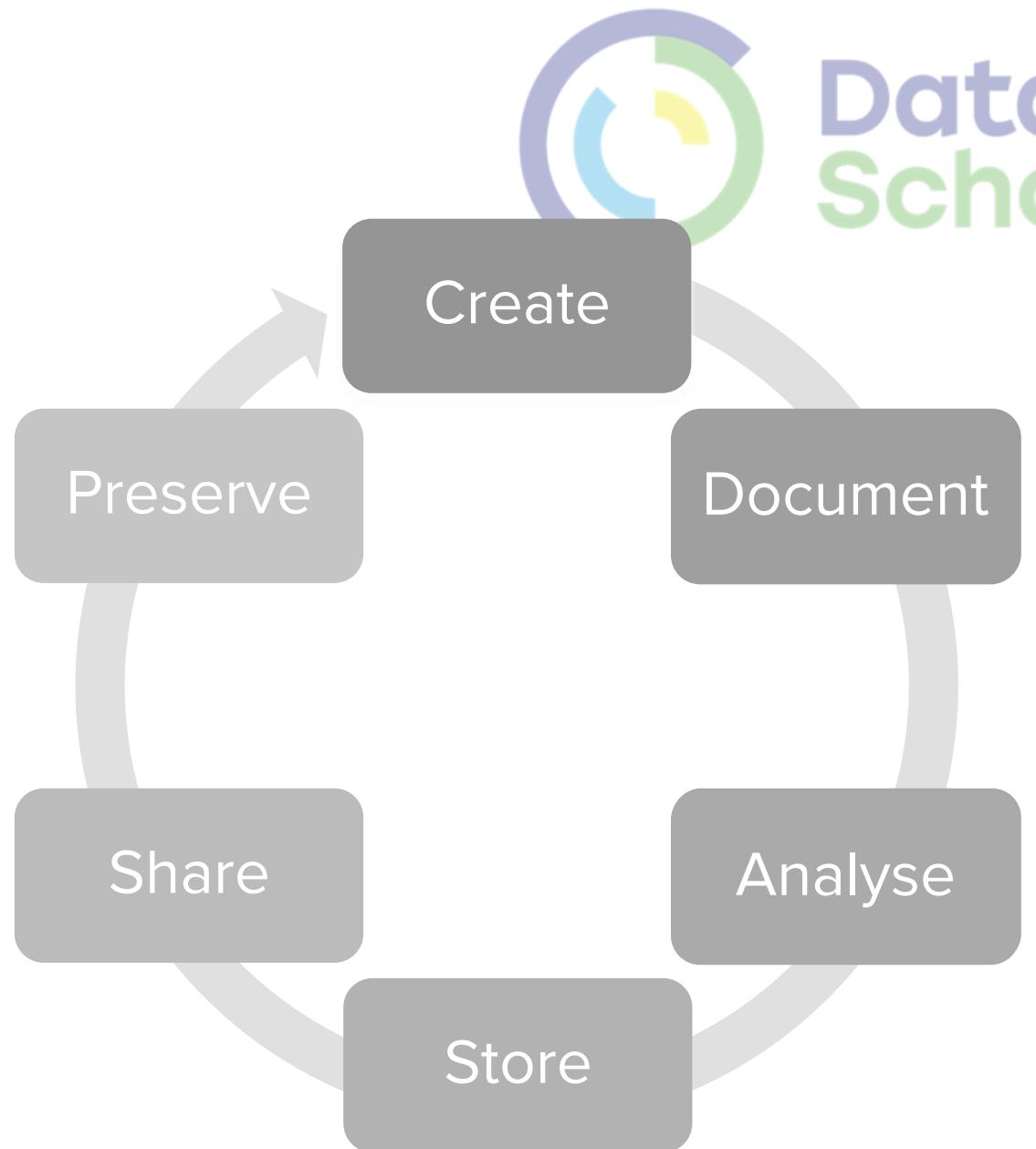


Example from ARM Climate Research Facility [www.arm.gov/data/docs/plan](http://www.arm.gov/data/docs/plan)

# Documentation

Think about what is needed in order to evaluate, understand, and reuse the data.

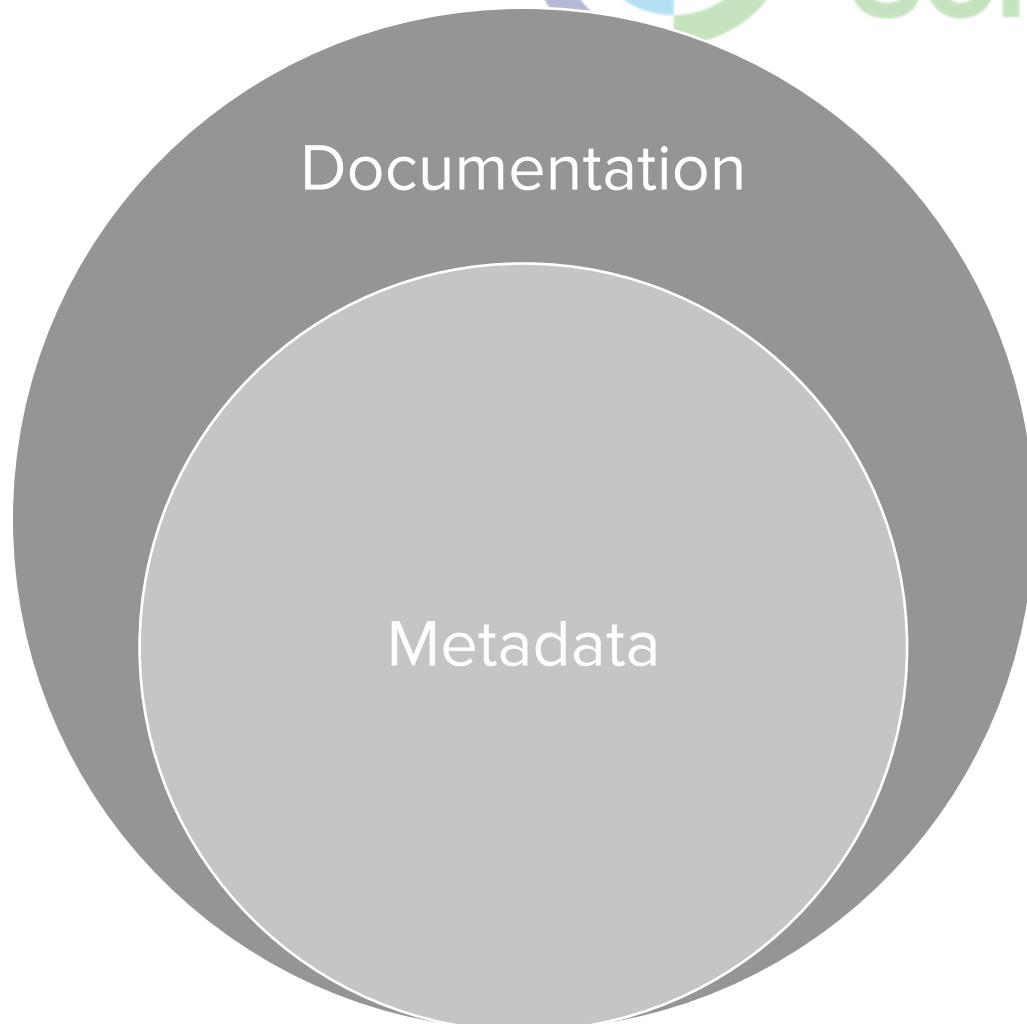
- Why was the data created?
- Have you documented what you did and how?
- Did you develop code to run analyses? If so, this should be kept and shared too.
- Important to provide wider context for trust.



# What are metadata?

---

- Metadata
    - Standardised
    - Structured
    - Machine and human readable
  - Metadata helps to cite and disambiguate data.
  - Documentation aids reuse.
- 





# Metadata standards

---

These can be general – such as Dublin Core

Or discipline specific

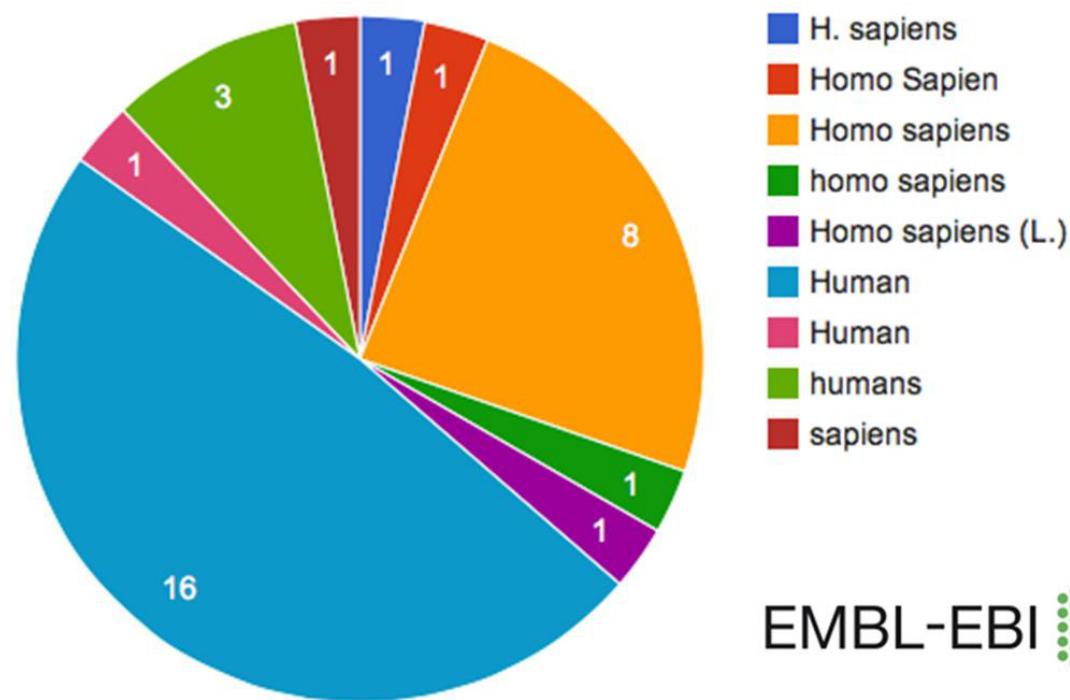
- Data Documentation Initiative (DDI) – social science
- Ecological Metadata Language (EML) - ecology
- Flexible Image Transport System (FITS) – astronomy

Search for standards in catalogues like:

- <http://rd-alliance.github.io/metadata-directory/>
  - <https://rdamsc.dcc.ac.uk/>
  - <http://www.fairsharing.org>
-

# Controlled vocabularies

*"MTBLS1: A metabolomic study of urinary changes in type 2 diabetes in....."*



# ...and ontologies?

- e.g. SNOMED CT (clinical terms) or MeSH
- Defined terms + taxonomy.
- Useful for selecting keywords to tag datasets.
- You can find many ontologies in the [BARTOC catalogue](#) and elsewhere.

➤ **Organism A**

- Term A1
- Term A2
- Term A3
  - Term B1
  - Term B2
- Term C4
- .
- .
- .
- Term n



► **Organism B**

- Term A1
- Term A2
- Term A3
  - Term B1
  - Term B2
- Term C4
- .
- .
- .
- Term n



❖ **Organism n**

- ❖ Term A1
- ❖ Term A2
- ❖ Term A3
  - ❖ Term B1
  - ❖ Term B2
- ❖ Term C4
- ❖ .
- ❖ .
- ❖ .
- ❖ Term n

# Where will you store the data?

- Your own device (laptop, flash drive, server etc.)
  - And if you lose it? Or it breaks?
- Departmental drives or university servers.
- “Cloud” storage.
- Do they care as much about your data?

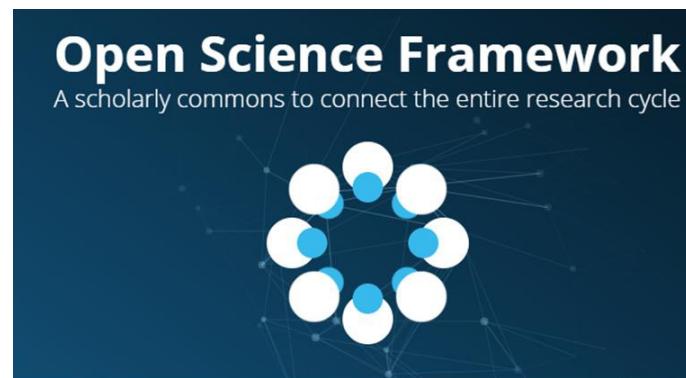
The decision will be based on how sensitive your data are, how robust you need the storage to be, and who needs access to the data and when.



# Collaborative platforms e.g. OSF

---

Open platform for sharing data in active phase with fellow researchers and others in secure environment.



Structured projects  
Keep all your files, data, and protocols in **one centralized location**. No more trawling emails to find files or scrambling to recover from lost data. [SECURE CLOUD](#)

Control access  
**You control which parts of your project are public or private** making it easy to collaborate with the worldwide community or just your team. [PROJECT-LEVEL PERMISSIONS](#)

Respect for your workflow  
Connect your favorite third party services directly to the Open Science Framework. [3RD PARTY INTEGRATIONS](#)

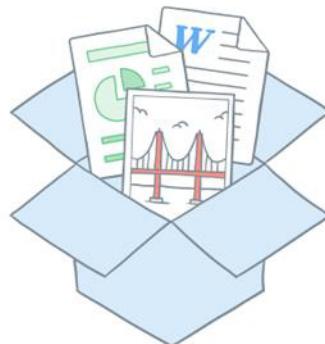
<https://osf.io>

# Third-party tools for collaboration

---

Dropbox, Google Drive, OneDrive and other cloud services

- Commercial
- Who owns your data?



ownCloud

- Open source product with Dropbox-like functionality.
- Used by many universities and service providers to offer 'approved' solution.



<https://owncloud.org>

# Backup and preservation—not the same thing!

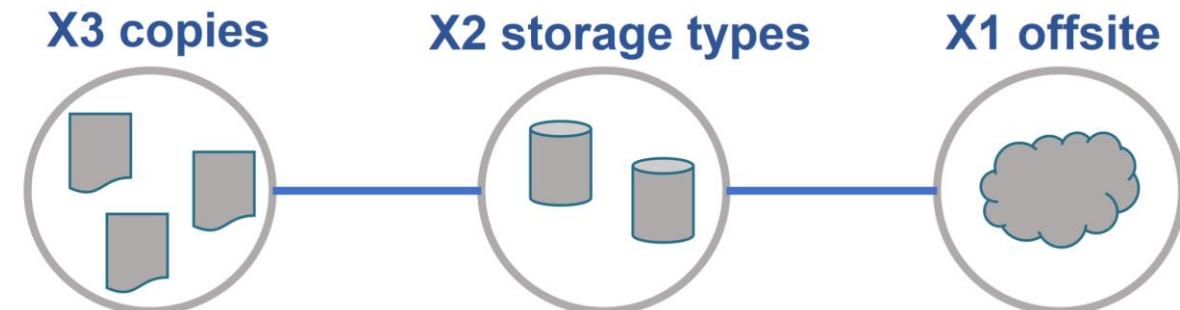


## Backups

- Used to take periodic snapshots of data in case the current version is destroyed or lost.
- Backups are copies of files stored for short or near-long-term.
- Often performed on a somewhat frequent schedule.

## Archiving

- Used to preserve data for historical reference or potentially during disasters.
- Archives are usually the final version, stored for long-term, and generally not copied over.
- Often performed at the end of a project or during major milestones.



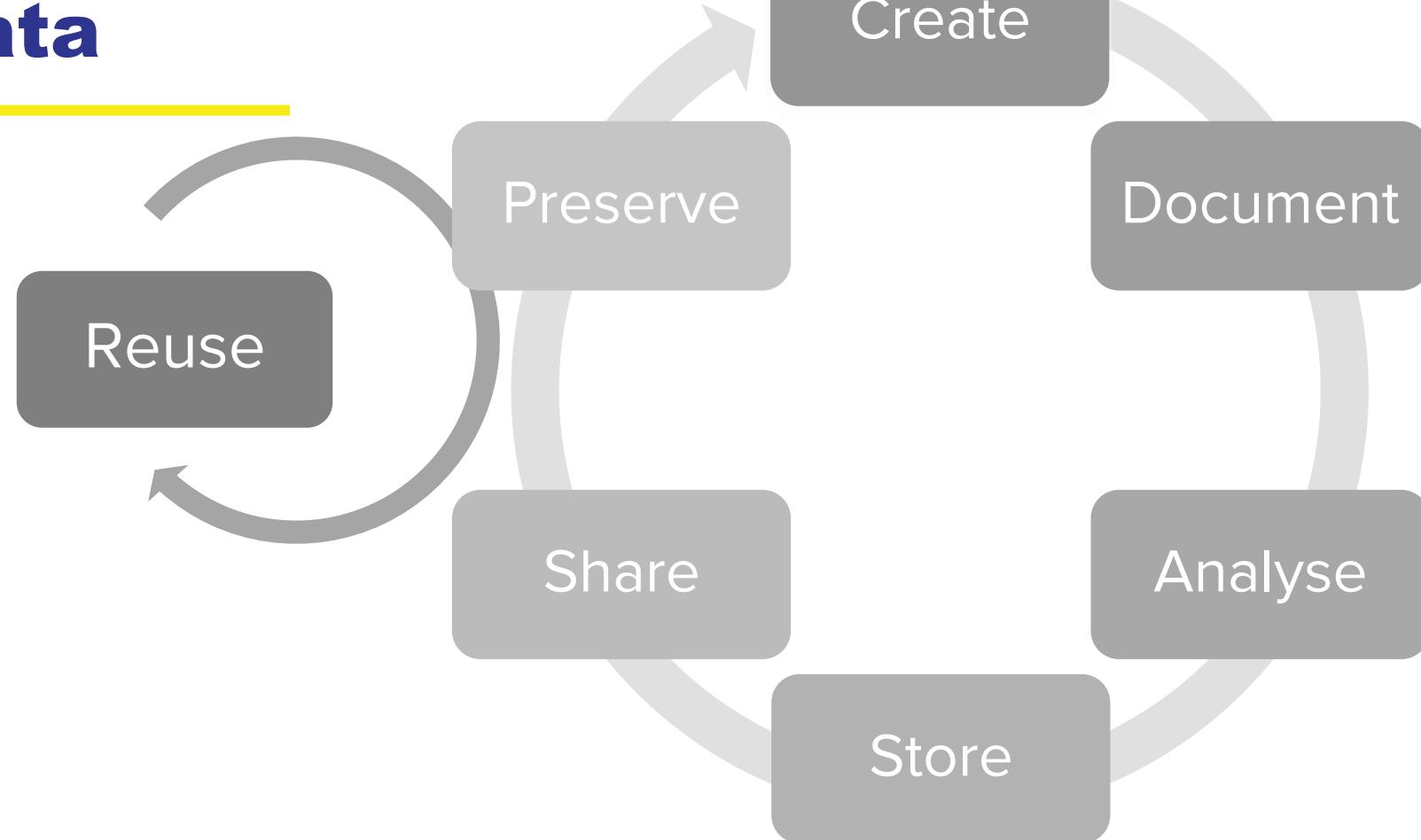
# How will you allow others to use your data?

Apply licences to disambiguate reuse restrictions.

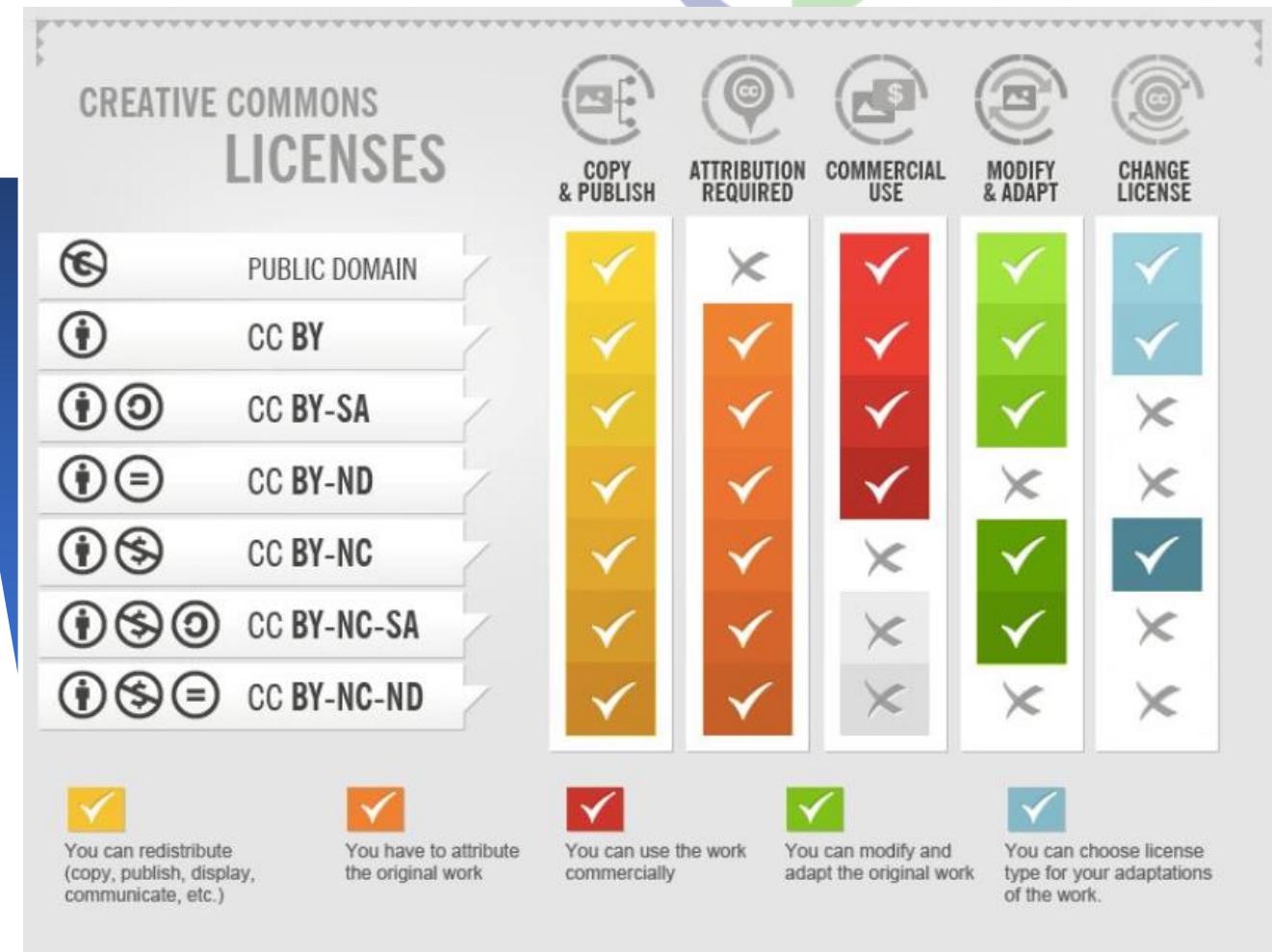


## Secondary vs primary data

---



# License research data openly



Part of [How To Attribute Creative Commons Photos by Foter](#), licensed CC BY SA 3.0

# Tools to decide which license to choose

Choose a license for your data

Check other researchers' license to know how to re-use their work

<https://chooser-beta.creativecommons.org/>



## SELECT YOUR LICENSE

Follow the steps to select the appropriate license for your work.

### 1 Do you know which license you need?

- Yes, I know which license I need.  
 No, I need help selecting a license.

NEXT STEP

2 Attribution

3 Commercial Use

4 Derivative Works

5 Sharing Requirements

6 Attribution Details



# Deposit in a data repository

Long-term  
preservation of data.



# Deposit in a data repository

The Re3data catalogue can be searched to find a home for data.

[www.fosteropenscience.eu/  
content/re3data-demo](http://www.fosteropenscience.eu/content/re3data-demo)

The screenshot shows the re3data.org search interface. On the left, there is a filter sidebar with various categories like Subjects, Content Types, Countries, AID systems, API, Certificates, Data access, Data access restrictions, Database access, Database access restrictions, Database licenses, Data licenses, Data upload, Data upload restrictions, Enhanced publication, Institution responsibility type, Institution type, Keywords, Metadata standards, PID systems, Provider types, Quality management, Repository languages, Software, Syndications, Repository types, and Versioning. The main search results area displays two entries: 'UniProtKB/Swiss-Prot' and 'Khazar University Institutional Repository'. Each entry includes basic metadata such as Subject(s), Content type(s), and Country. Below the entries, there is a world map titled 'Browse by country' with a play button overlay, indicating a video demo.

[www.re3data.org](http://www.re3data.org)

# Criteria for selecting a repository

- Better to use a domain specific repository if available.
- Check they match particular data needs e.g. formats accepted, mixture of Open and Restricted Access.
- Do they assign a persistent and globally unique identifier for sustainable citations and to links back to particular researchers and grants?
- Look for certification as a '*Trustworthy Digital Repository*' with an explicit ambition to keep the data available in long term.

The screenshot shows the EASY DANS-EASY repository interface. It includes sections for Subject(s) (History, Ancient Cultures, Social and Behavioural Sciences, Geosciences (including Geography), Humanities, Humanities and Social Sciences, Natural Sciences, Economics, Life Sciences), Content type(s) (Standard office documents, Images, Structured graphics, Audiovisual data, Raw data, Databases, Plain text, Structured text, Scientific and statistical data formats), and Country (Netherlands). A note at the bottom states: "EASY is the online archiving system of Data Archiving and Networked Services (DANS). EASY offers you access to thousands of datasets in the humanities, the social sciences and other disciplines. EASY can also be used for the online depositing of research data."



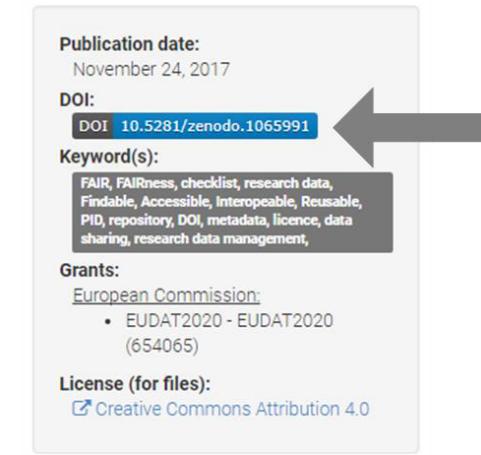
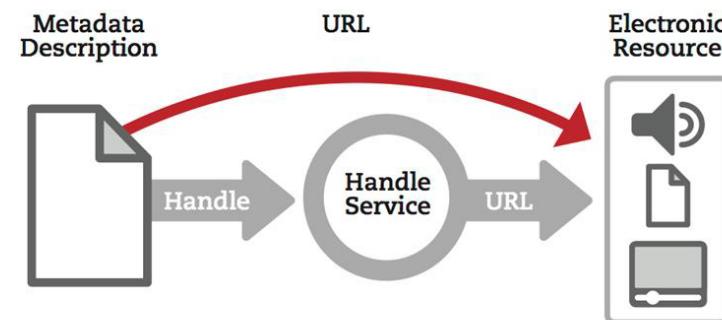
Icons to note open access, licences, PIDs, certificates...

[www.re3data.org](http://www.re3data.org)

# What is a Persistent Identifier (PID)?

*a long-lasting reference to a document, file or other object*

- PIDs come in various forms e.g. ORCID, DOI, ISBN...
- Typically they're actionable i.e. type it into web browser to access.
- Many repositories will assign them on deposit.



[www.re3data.org](http://www.re3data.org)

# Sensitive data

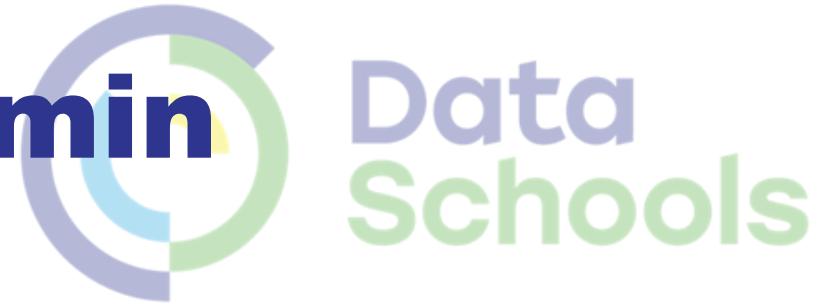
- Personal data (and metadata)
- Confidential data (trade secrets, investigations,...)
- Security data (passwords, financial information, national safety, military,...)
- Data protected by Intellectual Property Rights (IPR)
- Location Data/GPS/mobile phones
- Endangered (plant or animal) species, where their survival is dependent on the protection of their location data (biodiversity community)
- Combination of different datasets could lead to sensitive data?
  - racial or ethnic origin
  - political opinions
  - religious or philosophical beliefs
  - trade union membership
  - genetic data, biometric data
  - physical or mental health
  - sex life or sexual orientation
  - criminal offences

# Sensitive data best practices

---

- Access controls  
passwords, firewall (viruses, hacking)
  - Anonymisation  
removing or aggregating variables or reducing the precision or detailed textual meaning of a variable
  - Encryption  
encoded digital information
  - Share in a secure place  
no cloud drives
  - Store in an isolated machine  
server not connected to Internet
  - Secure disposal  
no data recovery is possible (uninstall)
-

# Exercise - 25 min (+ 20 min discussion)



---

Imagine you are a biologist who is doing microscopy experiments imaging tissue specimens. The data captured by the imaging is 100s of GB in size and is then cleaned and analysed to produce derivatives of the original captured data. Some of these derivatives may eventually be published. In preparation for publication, the data will also be segmented and annotated using standard ontologies. Documentation will also include metadata standards that will sufficiently describe the experimental procedure to allow reproducibility. Publication of the data is mandatory due to funder policy and must be deposited in a repository within 3 years of data production and must use an open licence without restrictions on reuse.

Now...please split into groups and see if you can answer the following questions using the tools and guidelines that have been described:

- What **file format(s)** should data be captured/preserved in?
  - Which **metadata standard(s)** should be used?
  - What **ontology(ies)** should be used?
  - Which **licence(s)** should be used?
  - Which **repository** would be the best fit for these data?
  - Do you foresee any problems with the data?
  - (Hint: not all the questions can be answered definitively! – but why not?)
-



Data  
Schools

**Data sharing and  
openness**

## Give us back our crown jewels

Our taxes fund the collection of public data - yet we pay again to access it. Make the data freely available to stimulate innovation, argue Charles Arthur and Michael Cross

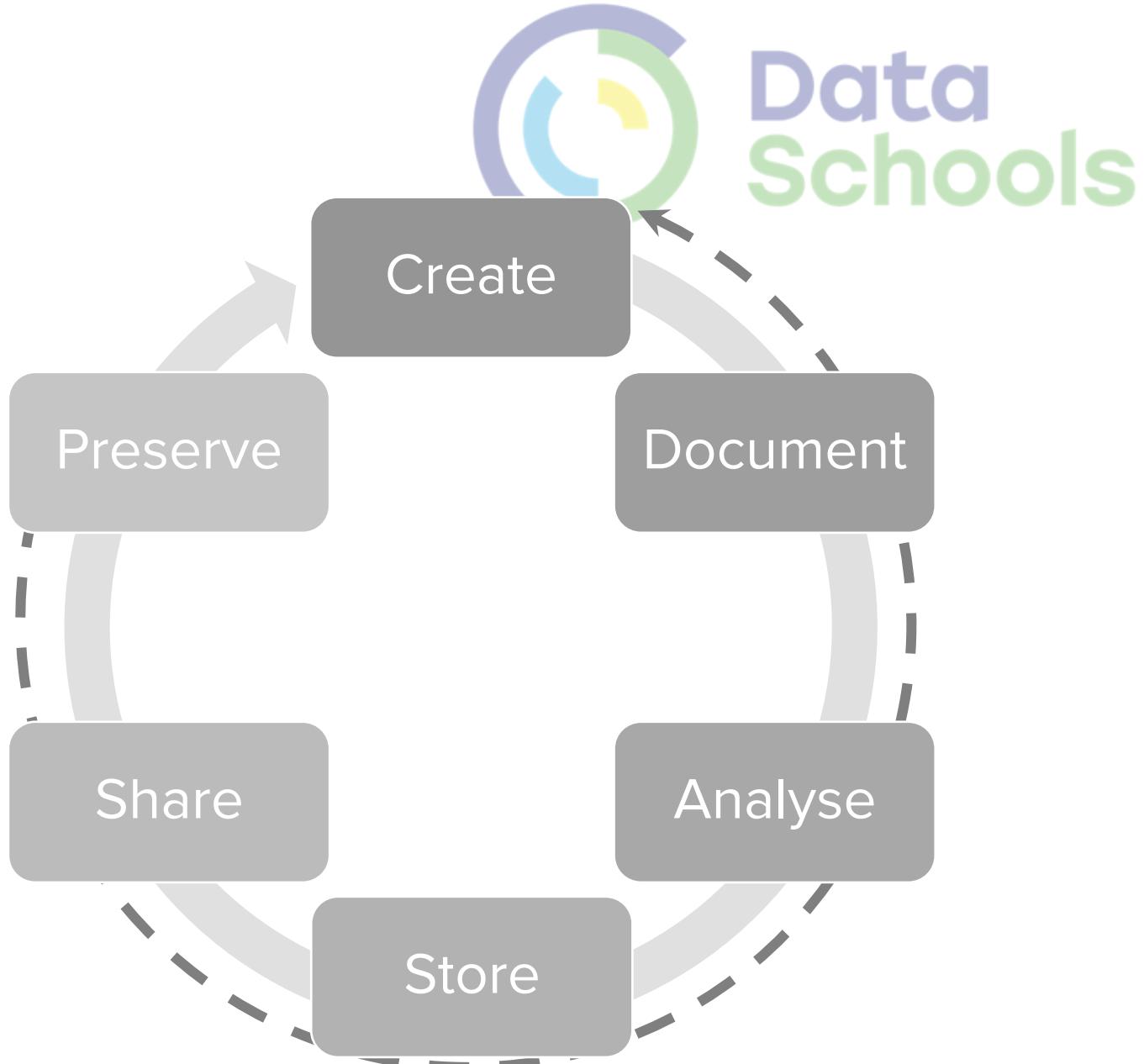
Charles Arthur and Michael Cross  
The Guardian, Thursday 9 March 2006  
[Article history](#)



## And open research...

---

- Change the typical lifecycle.
- Publish earlier and release more.
- Papers + Data + Methods + Code...
- Support reproducibility.



# Why make data available?

---

"It was \*never\* acceptable to publish papers without making data available."

- Ewan Birney

#OpenData  
#OpenScience



Original image via doi:10.1038/461145a. "Research cannot flourish if data are not preserved and made accessible. Data management should be woven into every course in science." - Nature 461, 145

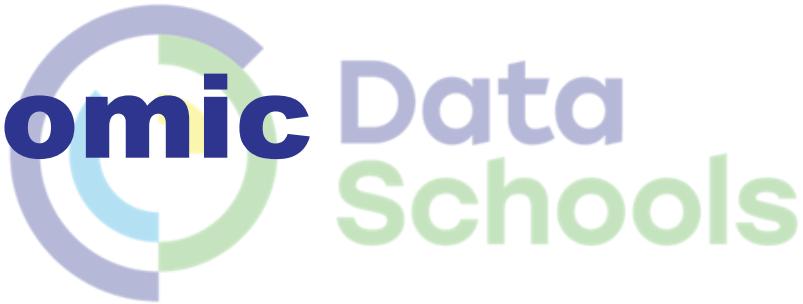
---

# The Old Weather Project

# Data for research, not from research



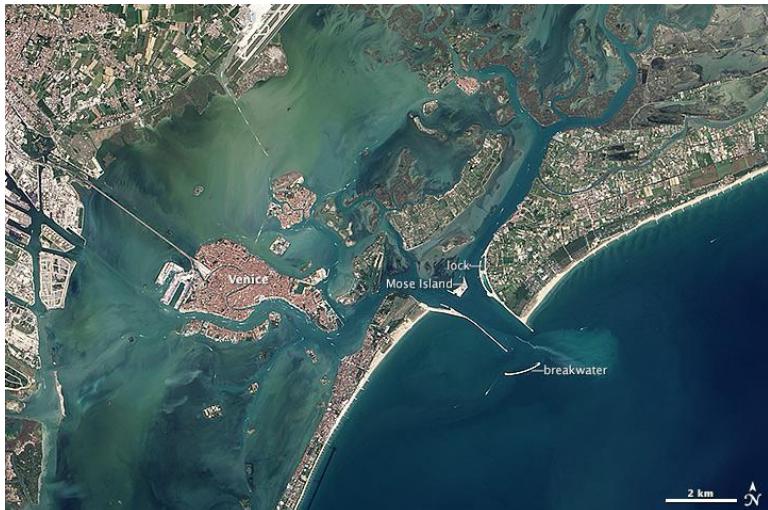
# Increased use and economic benefit



## The case of NASA Landsat satellite imagery of the Earth's surface

Up to 2008

- Sold through the US Geological Survey for US\$600 per scene
- Sales of 19,000 scenes per year
- Annual revenue of \$11.4 million



Since 2009

- Freely available over the internet.
- Google Earth now uses the images.
- Transmission of 2,100,000 scenes per year.
- Estimated to have created value for the environmental management industry of \$935 million, with direct benefit of more than \$100 million per year to the US economy.
- Has stimulated the development of applications from a large number of companies worldwide.
- <http://earthobservatory.nasa.gov/IOTD/view.php?id=83394&src=ve>

# Validation of results

---

*"It was a mistake in a spreadsheet that could have been easily overlooked: a few rows left out of an equation to average the values in a column.*

*The spreadsheet was used to draw the conclusion of an influential 2010 economics paper: that public debt of more than 90% of GDP slows down growth. This conclusion was later cited by the International Monetary Fund and the UK Treasury to justify programmes of austerity that have arguably led to riots, poverty and lost jobs."*

## The error that could subvert George Osborne's austerity programme

The theories on which the chancellor based his cuts policies have been shown to be based on an embarrassing mistake

---

Charles Arthur and Phillip Inman

The Guardian, Thursday 18 April 2013 21.10 BST



George Osborne says that Ken Rogoff, the man whose economic error has been uncovered, has strongly influenced his thinking. Photograph: Stefan Wermuth/PA

# Cut down on academic fraud

Stapel – 55 publications – “fictitious data”



# Data Schools

nature International weekly journal of science

nature news home news archive specials opinion features news blog nature journal

comments on this story

Stories by subject

- Brain and behaviour
- Lab life

Stories by keywords

- Diederik Stapel
- Tilburg University
- Academic fraud
- Retractions
- Social psychology

This article elsewhere

-  Blogs linking to this article
-  Add to Digg
-  Add to Facebook
-  Add to Newsvine
-  Add to Del.icio.us
-  Add to Twitter

Published online 1 November 2011 | Nature 479, 15 (2011) | doi:10.1038/479015a  
Updated online: 1 November 2011  
Updated online: 8 December 2011

News

## Report finds massive fraud at Dutch universities

Investigation claims dozens of social-psychology papers contain faked data.

Ewen Callaway

When colleagues called the work of Dutch psychologist Diederik Stapel too good to be true, they meant it as a compliment. But a preliminary investigative report ([go.nature.com/tqmp5c](http://go.nature.com/tqmp5c)) released on 31 October gives literal meaning to the phrase, detailing years of data manipulation and blatant fabrication by the prominent Tilburg University researcher.

"We have some 30 papers in peer-reviewed journals where we are actually sure that they are fake, and there are more to come," says Pim Levelt, chair of the committee that investigated Stapel's work at the university.

Stapel's eye-catching studies on aspects of social behaviour such as power and stereotyping garnered wide press coverage. For example, in a recent *Science* paper (which the investigation has not identified as fraudulent), Stapel reported that untidy environments encouraged discrimination ([Science 332, 251–253; 2011](http://science.sciencemag.org/content/332/6028/251.full)).

  
Dutch psychologist Diederik Stapel.  
Persbureau van Eindhoven

Related stories

- Seven days: 9–15 September 2011  
14 September 2011
- Chaos promotes stereotyping  
07 April 2011

Naturejobs

Tenure-Track Faculty Positions (Assistant / Associate / Full Professor) Yale University, Department of Genetics  
Yale University School of Medicine

Assistant Professor  
Harvard Medical School

 More science jobs

 Post a job for free

Resources

 PDF Format

 Send to a Friend

 Reprints & Permissions

 RSS Feeds

external links

- Tilburg University
- Interim investigation report

[www.nature.com/news/2011/111101/full/479015a.html](http://www.nature.com/news/2011/111101/full/479015a.html)



# Sharing leads to breakthroughs!

...and increases the speed of discovery

*"It was unbelievable. It's not science the way most of us have practiced in our careers. But we all realized that we would never get biomarkers unless all of us parked our egos and intellectual property noses outside the door and agreed that all of our data would be public immediately."*

*Dr John Trojanowski, University of Pennsylvania*

[http://www.nytimes.com/2010/08/13/health/research/13alzheimer.html?pagewanted=all&\\_r=0](http://www.nytimes.com/2010/08/13/health/research/13alzheimer.html?pagewanted=all&_r=0)

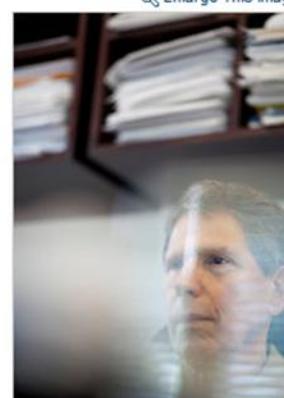
## Sharing of Data Leads to Progress on Alzheimer's

By GINA KOLATA

Published: August 12, 2010

In 2003, a group of scientists and executives from the [National Institutes of Health](#), the [Food and Drug Administration](#), the drug and medical-imaging industries, universities and nonprofit groups joined in a project that experts say had no precedent: a collaborative effort to find the biological markers that show the progression of [Alzheimer's disease](#) in the human brain.

[Enlarge This Image](#)



Now, the effort is bearing fruit with a wealth of recent scientific papers on the early diagnosis of Alzheimer's using methods like PET scans and tests of spinal fluid. More than 100 studies are under way to test drugs that might slow or stop the disease.

And the collaboration is already serving as a model for similar efforts against [Parkinson's disease](#). A \$40 million project to look for biomarkers for Parkinson's, sponsored by the [Michael J. Fox Foundation](#), plans to enroll 600 study subjects in the United States and Europe.

# How do you share data effectively?

---

- Use appropriate repositories, this catalogue is a good place to start:  
<http://www.re3data.org>
  - Document and describe it enough for others to understand, use and cite:  
<http://www.dcc.ac.uk/resources/how-guides/cite-datasets>
  - License it so others can reuse:  
[www.dcc.ac.uk/resources/how-guides/license-research-data](http://www.dcc.ac.uk/resources/how-guides/license-research-data)
- 



Data  
Schools



# Who has heard of this before...?

---

**F**indable    **A**ccessible    **I**nteroperable    **R**eusable

- Metadata
- PIDs
- Repositories

- Metadata
- Open file formats and software

- Metadata
- Ontologies
- Repositories

- Metadata
- Licences

# European perspective...

<https://publications.europa.eu/en/publication-detail/-/publication/7769a148-f1f6-11e8-9982-01aa75ed71a1/language-en/format-PDF/source-80611283>



# What FAIR means: 15 principles

## Findable:

- F1. (meta)data are assigned a globally unique and persistent identifier;
- F2. data are described with rich metadata;
- F3. metadata clearly and explicitly include the identifier of the data it describes;
- F4. (meta)data are registered or indexed in a searchable resource;

## Interoperable:

- 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;

## Accessible:

- A1. (meta)data are retrievable by their identifier using a standardized communications protocol;
  - A1.1 the protocol is open, free, and universally implementable;
  - A1.2. the protocol allows for an authentication and authorization procedure, where necessary;
- A2. metadata are accessible, even when the data are no longer available;

## Reusable:

- 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;

doi: [10.1038/sdata.2016.18](https://doi.org/10.1038/sdata.2016.18)

Slide CC-BY by Erik Schultes, Leiden UMC

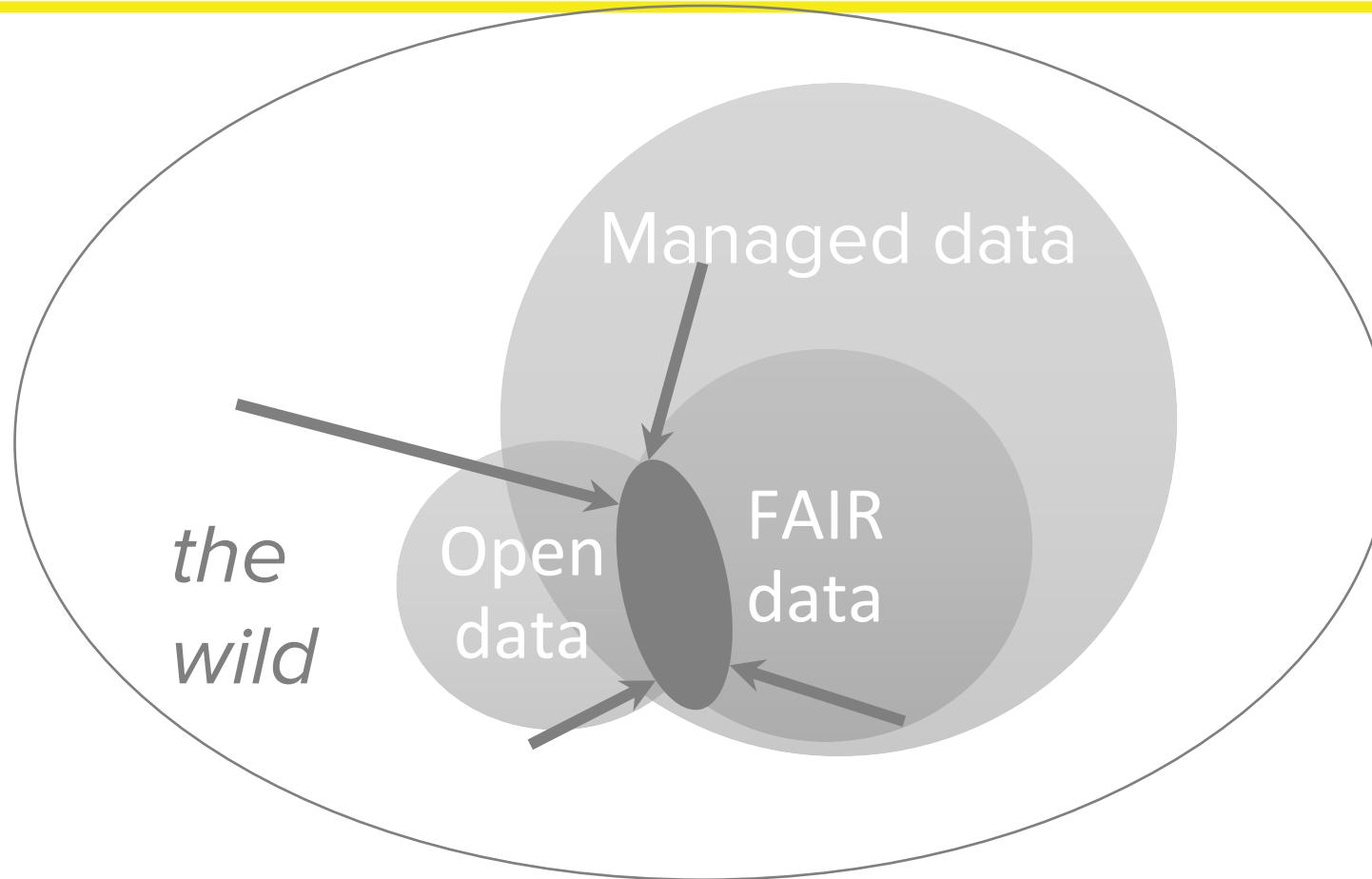
Comprehensive descriptions can be found at <https://www.go-fair.org/fair-principles/>

# Common misconceptions

---

- FAIR data does not have to be open.
  - The principles do not specify particular technologies or implementations e.g. semantic web.
  - FAIR is not a standard to be followed or strict criteria – it's a spectrum/continuum.
  - It doesn't only apply to the life sciences.
-

# Increasing that which is FAIR & open



Adapted from [DCC](#).

## FAIR ≠ Open

---

as open as  
possible, as  
closed as  
necessary



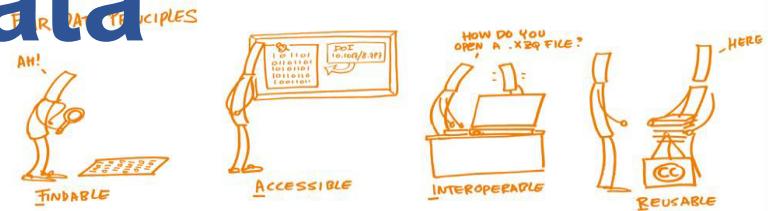
Image: 'Balancing rocks' by Viewminder CC-BY-SA-ND [www.flickr.com/photos/light\\_seeker/7780857224](http://www.flickr.com/photos/light_seeker/7780857224)

---



# Check how FAIR is your data

F-UJI Home Assess About Methods Docs Code



Automated FAIR Data Assessment Tool

F-UJI is a web service to programmatically assess FAIRness of research data objects at the dataset level based on the FAIRsFAIR Data Object Assessment Metrics ☺

[Click here to assess a dataset.](#)

F-UJI was developed by Anusuriya Devaraju & Robert Huber ([PANGAEA](#)) under the umbrella of the [FAIRsFAIR](#) project.

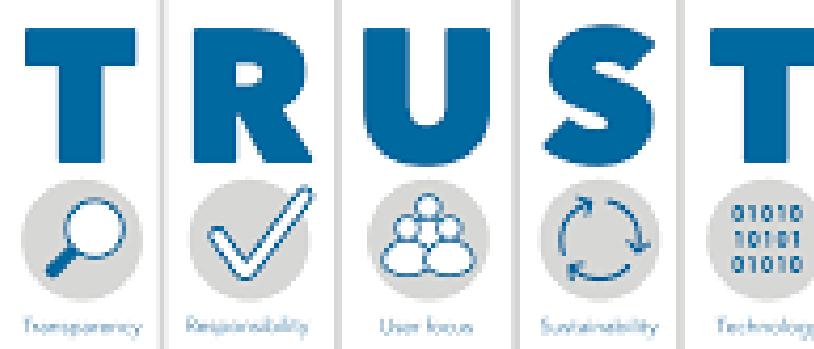
[About](#) [Feedback](#) [Privacy Policy](#) [Terms of Use](#) [Legal Notice](#)

<https://ardc.edu.au/resources/working-with-data/fair-data/fair-self-assessment-tool/>

<https://www.f-oji.net/>

**FAIR isn't the  
only  
consideration...**

---



# New(ish) frontiers...

<https://unesdoc.unesco.org/ark:/48223/pf0000387324>

- Collaboration
- Reproducibility
- Transparency
- Trust

**"Open science practices are on the rise but access to, participation in and sharing of the benefits from open science are uneven across the world."**



## Open Science Outlook 1

Status and trends around the world



# New(ish) frontiers...

<https://www.springernature.com/gp/researchers/campaigns/state-of-open-data>

- **Support is not making its way to those who need it**  
Almost three-quarters of respondents had never received support with making their data openly available.
- **One size does not fit all**  
Variations in responses from different subject expertise and geographies highlight a need for a more nuanced approach to research data management support globally.
- **Challenging stereotypes**  
Are later career academics really opposed to progress? The results of the 2023 survey indicate that career stage is not a significant factor in open data awareness or support levels.
- **Credit is an ongoing issue**  
For eight years running, our survey has revealed a recurring concern among researchers: the perception that they don't receive sufficient recognition for openly sharing their data.
- **AI awareness hasn't translated to action**  
For the first time, this year we asked survey respondents to indicate if they were using ChatGPT or similar AI tools for data collection, processing and metadata creation.

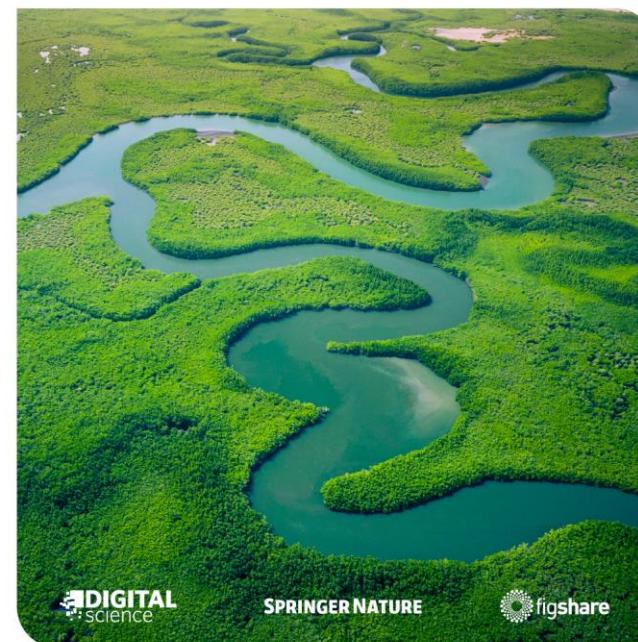
A Digital Science Report

November 2023

## The State of Open Data 2023

The longest-running longitudinal survey and analysis on open data.

With opening remarks from Springer Nature's CPO, Harsh Jegadeesan, and Digital Science's CEO, Daniel Hook. Authors Mark Hahnel, Graham Smith, Niki Scaplehorn, Henning Schoenenberger and Laura Day.



SPRINGER NATURE



The State of Open Data

April 2024

## The Global Lens:

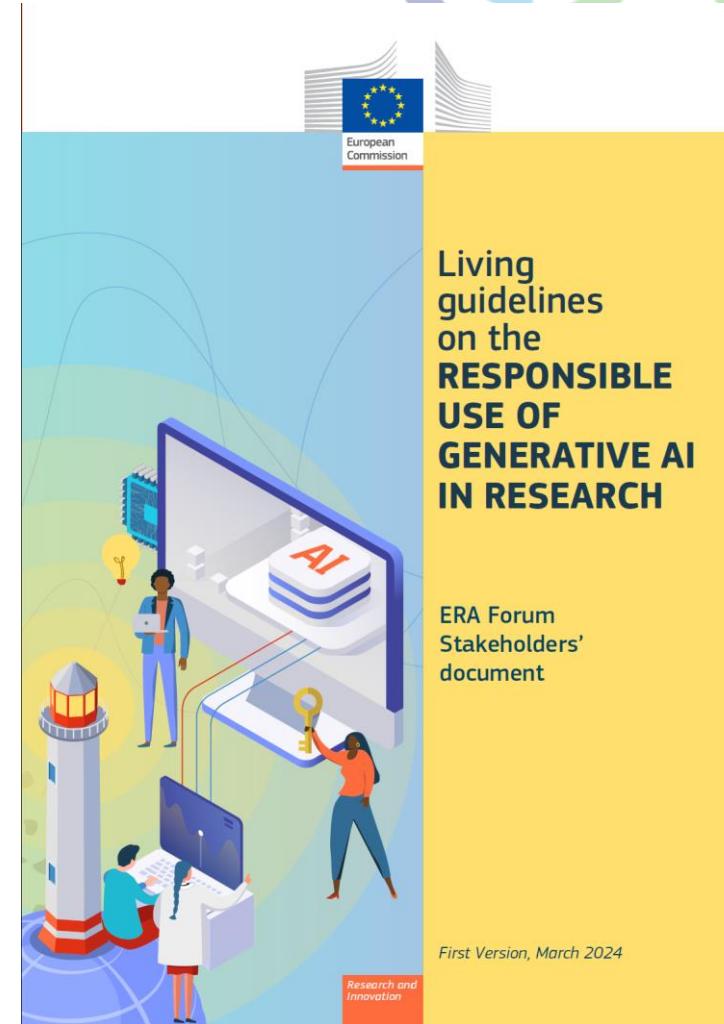
Highlighting national nuances in researchers' attitudes to open data



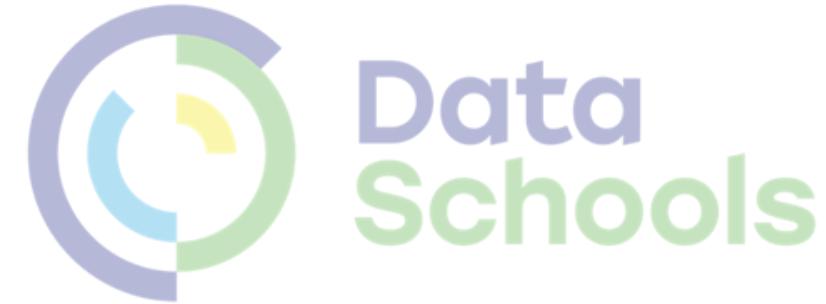
# New(ish) frontiers...

---

<https://european-research-area.ec.europa.eu/news/living-guidelines-responsible-use-generative-ai-research-published>



# FOSTER Open Science



What is Open Science?	Best Practice in Open Research	Open Access Publishing	Open Peer Review	Sharing Preprints
Data Protection & Ethics	Open Source Software & Workflows	Managing & Sharing Research Data	Open Science & Innovation	Open Licensing

<https://www.fosteropenscience.eu/toolkit>

# Research Data Alliance



A screenshot of the RDA website homepage. The header includes the RDA logo, navigation links for "O&amp;A Members", "MEMBERSHIP", "RDA Groups", and social media links. Below the header is a search bar and a main banner with the text "FIND YOUR GROUP by topic or discipline" and a magnifying glass icon. The main content area is divided into sections: "NEWS &amp; EVENTS" (with a "Submit your news" button), "RECENT BLOGS" (listing posts like "The Greek effect in enhancing RDA outputs adoption" and "RDA Secretariat Face-to-face Meeting, July 2019"), and "FOLLOW US" (a Twitter feed from @rdalliance). At the bottom, there are two sections: "The Value of RDA for" (with icons for individuals, organizations, students, funders, libraries, European Open Science Cloud, and regions) and "RDA in one Word" (a video thumbnail).

<https://www.rd-alliance.org>



Data  
Schools

# Data Management Plans

# Bringing together what you've learnt

- Make informed decisions to anticipate and avoid problems.
- Avoid duplication, data loss and security breaches.
- Develop procedures early on for consistency.
- Ensure data are accurate, complete, reliable and secure.
- Save time and effort to make your life easier!
- Useful both to researchers and institutions

## Making plans

*They sound dull, but data-management plans are essential, and funders must explain why.*

**D**ata are the lifeblood of scientific and social research. A versatile tool, they can both be raw material for producing knowledge and, when processed and interpreted with an expert eye, the end product of the exercise.

So it might sound perverse that researchers should be cautious about how they store and share their data — and the data they generate and use. The problem is that this can be hard to do.

As science produces day by day a huge volume of data, it's a growing challenge to manage them. In response, many funding agencies now ask applicants to submit a concrete data-management plan with their grant proposals. Effectively, a to-do list that details how they will handle their data from start to finish.

Such plans are important, and are something that *Nature* supports (we discuss them in detail in a Careers article on page 403). But to accelerate acceptance of what seems like just another administrative burden, we must make them meaningful and useful.

First, rigorously collected, well-presented data sets — including maps, documents and tables — will give the data owners so much solid, meaningful results. Second, they will help future investigators to make sense of and reuse data, thereby enhancing utility and reproducibility. Providing comprehensive data, ideally for many years after the work is completed, will facilitate the process of peer review. There is no single recipe for proper data management. The task varies according to the field of science, project size and the specific types

of data in question. That makes cross-disciplinary common standards unlikely, so research agencies need to engage with different scientific communities to find the best way forward. In addition, to provide a broadset of standards, formats and data protocols — indispensable in our increasingly global scientific enterprise — research agencies in all parts of the world must engage.

An excellent recent international alignment of research data-management policies, launched in January by Science Europe and the Netherlands Organisation for Scientific Research, is an important step in the right direction. A study of data management in physics and genetics shows that internationally aligned data governance not only is probably doable, but also has a positive impact on collaborative research. A good example is the agreement, setting up a centre in Paris, to specifically align the data used in the Large Hadron Collider at CERN.

The message must now be passed on to scientists who work in fields less familiar with big data. Many of these, at all career stages, are worried about the new requirements. Some have never been asked to provide a data-management plan, and that most are unaware of policies and guidelines already in place to encourage data sharing. In response, the European Commission, the European Research Council and the European Council of Doctoral Candidates and Junior Researchers, had actually written a data-management plan, with another quarter saying they didn't even know what such a plan was or what it entailed, and that they were not sure what to do.

Funders and universities, then, must ensure that the rationale of data management, and the basic skills of executing it properly, become part of the culture of research. Data management and support must go further and be offered at every career level. The laudable move towards open science — under which data are shared — makes the need for good data management and processing that ever more pressing. If sharing your data if they're clean and annotated enough to be reused. If you haven't got a plan for your data, you need one now. ■

384 | NATURE | VOL 555 | 15 MARCH 2018

© 2018 Macmillan Publishers Limited, part of Springer Nature. All rights reserved.



# Data Schools

## CAREERS

**PERSONAL ETIQUETTE** How a vegetarian biologist balances his beliefs with his work p.405 | <http://tinyurl.com/naturejobs>

**NATUREJOBS** For the latest career listings and advice [www.naturejobs.com](http://www.naturejobs.com)



## For the record

*Making project data freely available is vital for open science.*

BY QUIRIN SCHIERMEIER

**W**hen Marterie Etique learnt she had to create a data-management plan for her next research project, she was not sure exactly what to do.

The geochemist, based at the Swiss Federal Institute of Technology (ETH) in Zurich, studies the interaction of trace elements in sediments and water. While preparing a grant proposal for the Swiss National Science Foundation last October, she learnt of the funder's new data rules. These require applicants to provide a written plan for the organization and long-term storage of their research data, to help minimize the risk of data

loss and provide guidance for other scientists on how to use the data in the future.

Etique found the task daunting. "Data management is really not my primary skill," she says. "I had absolutely no idea how to go about it." She was able to get advice from her supervisor and from others in her field. Other researchers might not be so lucky, and may not even know what a data-management plan is — let alone why they would need one and how to produce it. Here, we answer these questions.

### WHAT ARE DATA-MANAGEMENT PLANS?

A data-management plan explains how researchers will handle their data during and after a project, and encompasses creating,

sharing and preserving research data of any type, including text, spreadsheets, images, recordings, models, algorithms and software. It does not matter whether the data are generated by large pieces of research equipment, such as telescopes, particle accelerators, or from straightforward field surveys.

Many funders are asking grant applicants to provide data plans. Requirements vary from one discipline to another. But in general, scientists will need to describe — before they begin any research — what data they will generate; how the data will be disseminated, described, secured and curated; and who will have access to those data after the research is completed. They must also explain any data sharing and reuse restrictions, such as legal and confidentiality issues. Researchers can consult their funder and their host institution's digital library services for assistance. Colleagues who have previously produced data plans may also be able to help (see *'Keeping stock'*).

### WHO NEEDS THEM?

Data management is one example of the way in which public research sponsors and research institutions are implementing open science, the push to make research outputs and data freely accessible. Many funding agencies have made data-management plans mandatory for grant applicants in the past decade or so. All US federal agencies, including the National Science Foundation and the National Institutes of Health, now require data-management plans must also be included in grant proposals to the European Research Council and other European Union–funded research programmes. And many national funding agencies in Europe — including the UK research councils and the London-based Wellcome Trust, which funds the medical research charity — also ask for data plans.

Many scientists already practise data management by default. Astronomers, for example, have been doing so for decades when calibrating their observations and archiving huge amounts of telescope-survey data in standardized, machine-readable catalogues for reuse.

Geneticists, too, use special data repositories to archive the vast amounts of DNA and genome-sequencing data (see [go.nature.com/2om2lrb](http://go.nature.com/2om2lrb)). But less data-intensive fields of science and social research also benefit from data management. For example, geochemists analysing soil bacteria and mineral products in different environments can use it to ▶

15 MARCH 2018 | VOL 555 | NATURE | 403

Schiermeier, Q. "Data management made simple" *Nature* 555, 403–405 (2018).

<https://www.nature.com/articles/d41586-018-03071-1>

doi: 10.1038/d41586-018-03071-1

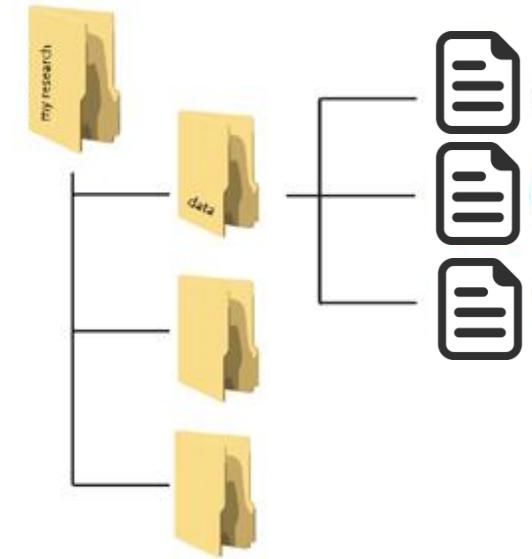
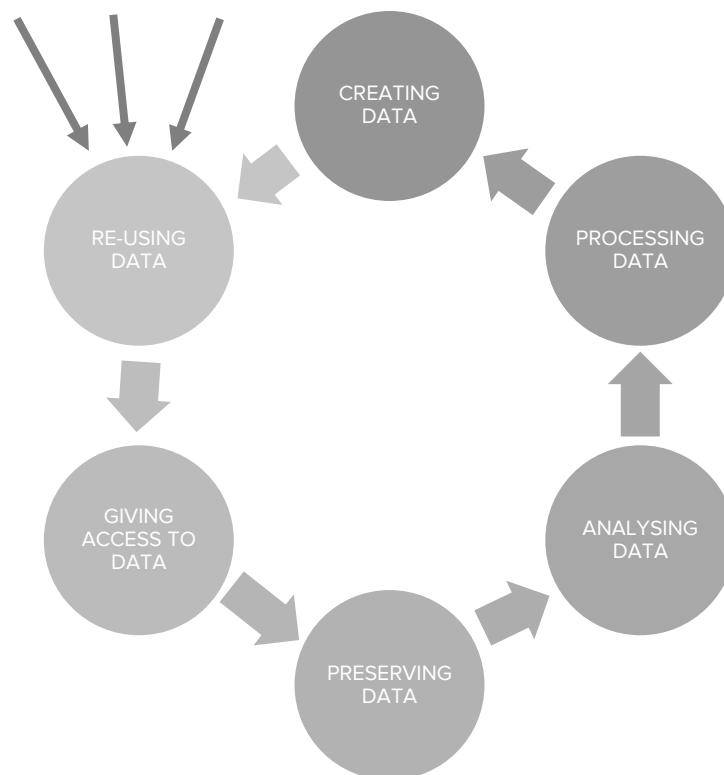
# Common themes in DMPs

---

1. Description of data to be collected / created (i.e. content, type, format, volume...).
  2. Standards/methodologies for data collection & management.
  3. Ethics and Intellectual Property (highlight any restrictions on data sharing e.g. embargoes, confidentiality).
  4. Plans for data sharing and access (i.e. how, when, to whom).
  5. Strategy for long-term preservation.
-

# Planning trick 1: think backwards

What data organisation would a re-user like?



Design how you will organise data in the project (folder structure, file naming convention, ...)

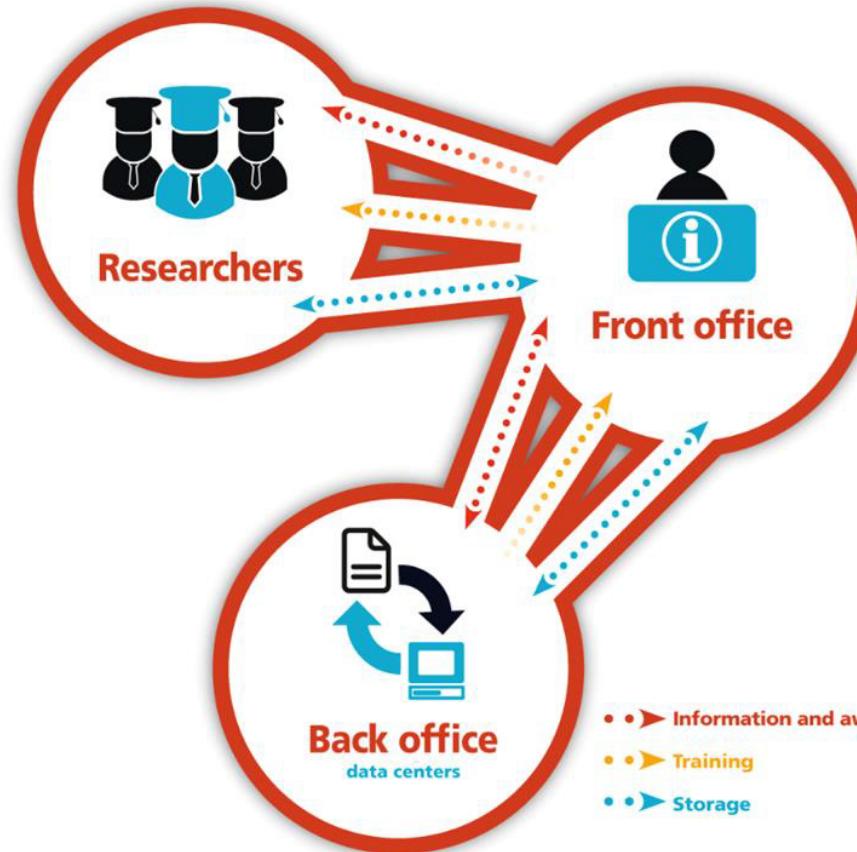
# Planning trick 2: include RDM stakeholders



Commercial  
partners



Publishers  
Data Availability  
policy



Institution  
RDM policy  
Facilities

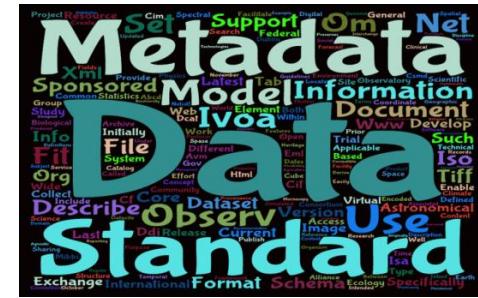


Research  
funders

# Planning trick 3: ground your plan in reality



Base plans on available skills, support and good practice for the field – show it's feasible to implement.





# What makes a good DMP?

---

- Clear, detailed information that is relevant to the science:
    - adopting recognised standards.
    - practices in line with norms for that field.
    - use of support services e.g. university storage, subject repositories...
  - Realistic approach that is feasible to implement.
  - Evidence of consultation and seeking advice.
  - Proper justification of restrictions and costs.
  - **Have you taken time to reflect on what to do?**
-

# Is the information specific enough?

---



*“we will use suitable formats to ensure that our data can be preserved and sustained over the long term”*

- Which standards? Name them!
  - Show that you know which are suitable.
  - Does your chosen repository have preferences?
-



# Are decisions justified?

---

*“data will be made available upon request to bona fide medieval historians”*

- Why is it restricted?
  - Could other communities not reuse the data?
  - Will the research team be around to handle access requests in the future?
-



# A better response...

---

*"We will provide MP3 audio files for online dissemination. While this is not an open format, it is well-established and the most widely supported. High-resolution WAV files will be used for the archival master recordings."*

- Be clear, specific and detailed.
  - Justify decisions.
-



# Example plans

---

- Plans from several funders and disciplines via DCC  
[www.dcc.ac.uk/resources/data-management-plans/guidance-examples](http://www.dcc.ac.uk/resources/data-management-plans/guidance-examples)
  - Scientific DMPs submitted to the NSF (USA) provided by DataOne  
<https://www.dataone.org/data-management-planning>
  - DMPs published in RIO journal  
[http://riojournal.com/browse\\_user\\_collection\\_documents.php?collection\\_id=3&journal\\_id=17](http://riojournal.com/browse_user_collection_documents.php?collection_id=3&journal_id=17)
  - Share yours! - [www.dcc.ac.uk/share-DMPs](http://www.dcc.ac.uk/share-DMPs)
-

# DCC Checklist for a DMP

---

- The DCC assessed existing funder requirements, DMP templates and other best practice to see what should be included in plans. This was synthesised down into common themes and questions.
- 13 questions on what's asked across the board.
- Prompts/pointers to help researchers get started.
- Guidance on how to answer.





# Thank you!

---

Questions?

(Please get in  
touch!)

