

Introduction to Research Data Management

... and how not to get overwhelmed by data

Lecturer: Jan Vališ 

Authors of the presentation: Adéla Jílková, Martin Schätz, Jan Vališ

April 8, 2026

Content of this presentation is licensed via [CC BY 4.0](#),
except where otherwise noted for content created
by third-parties.



How about you?

- What is your academic affiliation?
- What is your level of study?
- What is your main research field?

Learning objectives

Knowledge:

- Define what research data are
- Identify key stages of the research data lifecycle
- Understand the FAIR principles
- Distinguish FAIR from Open data

Skills

- Outline a DMP, incl. administrative, legal, ethical requirements.
- Estimate storage costs
- Apply basic metadata

Attitudes:

- Research data management is a good scientific practice

What is research data and why manage it?

Research data – theoretically

- Any information **collected, observed, generated, or created** during the research process to produce and support research findings
- Primary data
 - Directly from a given research
 - Quantitative (volume) vs. Qualitative (interview)
 - Experimental (pH) vs. Observational (bird migration)
- Secondary data
 - From other research used for different reason

Research data – practically

- Pictures (SEM pictures, photos of birds etc.)
- Videos (movement of particles, interviews etc.)
- Sound records (birds singing, interviews etc.)
- Tables (spectra, observations etc.)
- Texts
- Statistics
- Electronic patient records
- ...

Research data are not:

- Project proposal
- Project budget
- Evaluation report
- Conference proceedings
- Journal article
- Electronic supplementary information

Research data management

- A set of practices, strategies, and activities, including data:
 - **organization,**
 - **documentation,**
 - **storage, and**
 - **sharing**
- Covers all stages of the research process
- Ensures the effectiveness, reproducibility, and reuse of research data

Why manage research data?

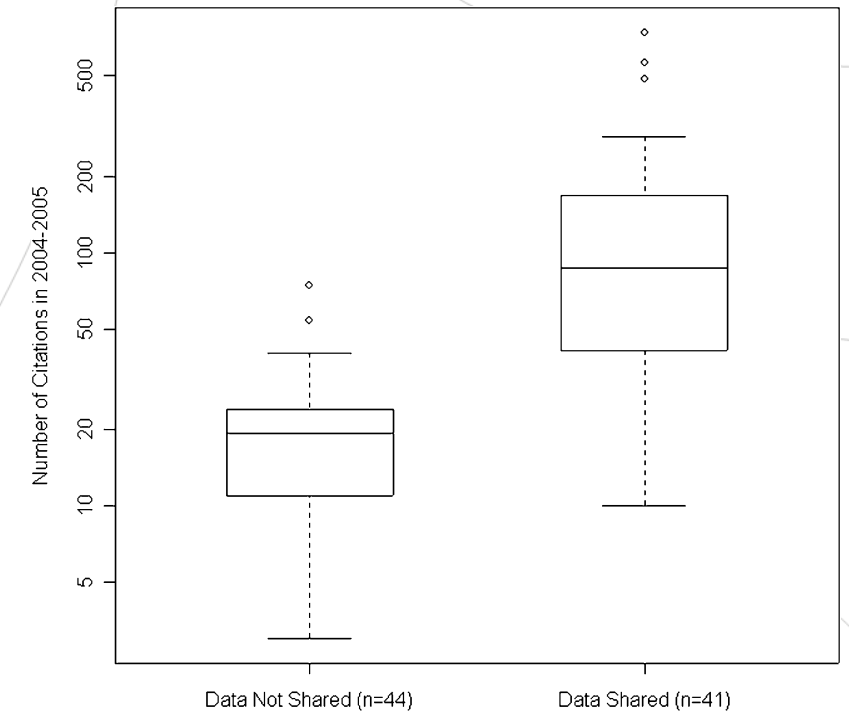
It may be mandatory (institutional, publisher, or research funder requirements)

Keep the research process secure and organized

- Increase efficiency, save time and resources
- Share data with colleagues
- Reduce risk of data loss and improve data security

Enhance global data sharing

- Enable data reuse and enhance collaboration
- Increase the visibility and impact of research
- Increase transparency & improve trust in findings
- Support research integrity & validation of results



2004–2005 citation counts of 85 trials by data availability.
Heather A. et al. 2007. PLOS ONE. License: [CC BY 4.0](https://creativecommons.org/licenses/by/4.0/)
[10.1371/journal.pone.0000308](https://doi.org/10.1371/journal.pone.0000308).

Research data

Different fields and disciplines

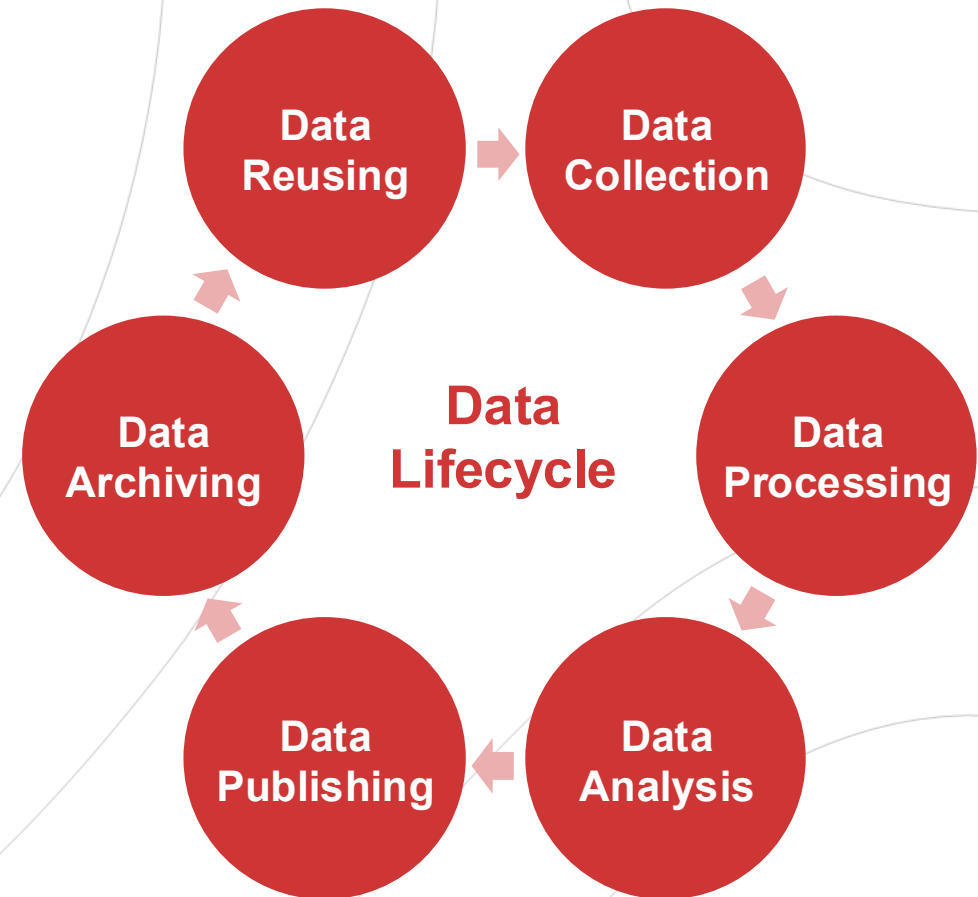
- Natural and life sciences
- Medical and health sciences
- Engineering and technology
- Social sciences
- Arts and humanities

Research data

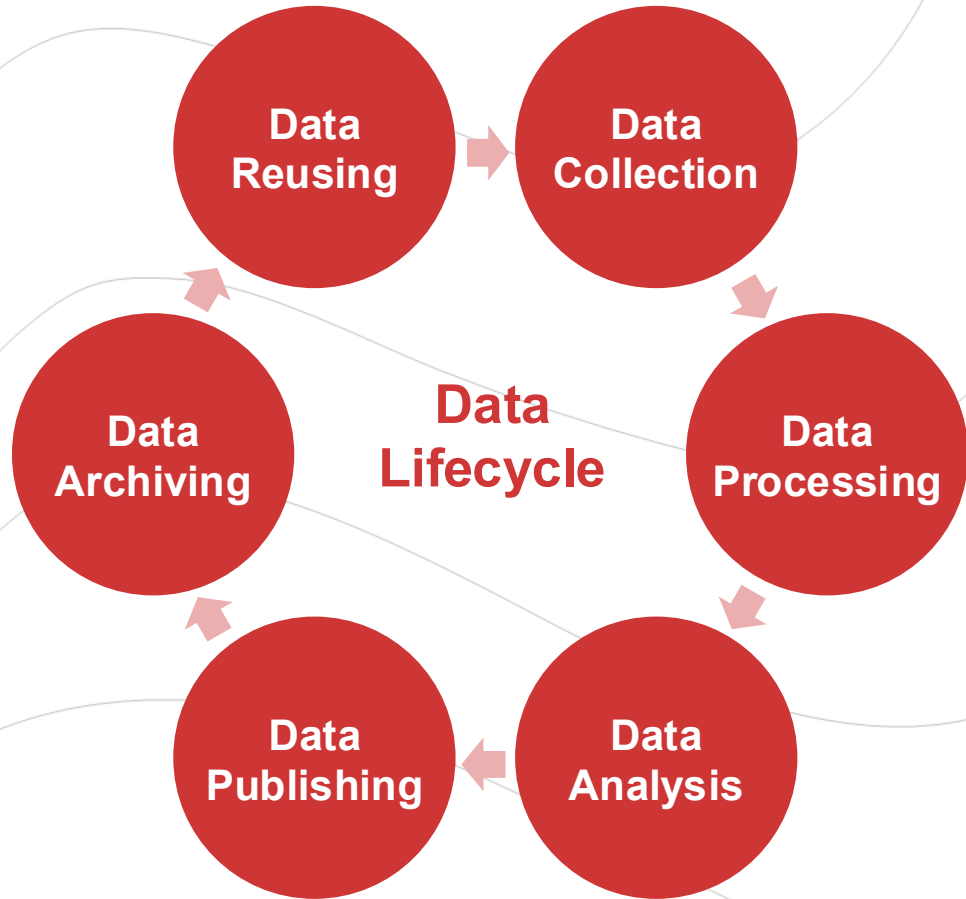
Different fields and disciplines

- Natural and life sciences
- Medical and health sciences
- Engineering and technology
- Social sciences
- Arts and humanities

Different stages of research data lifecycle



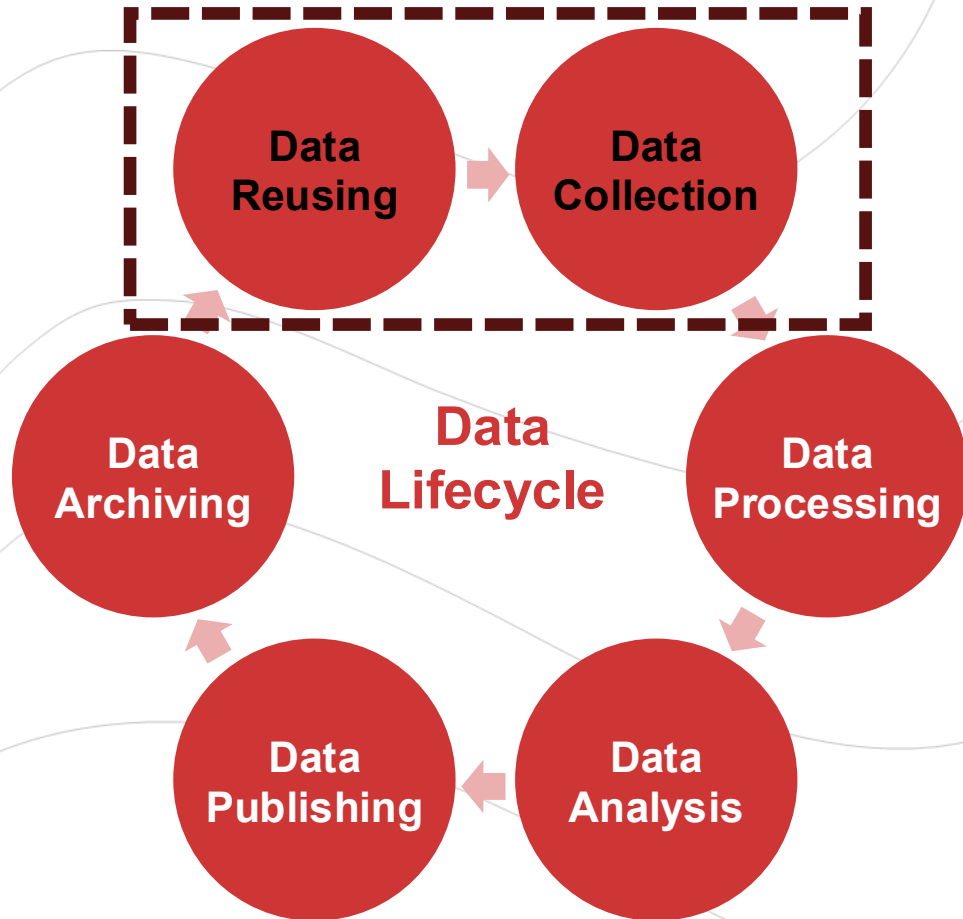
Research data lifecycle



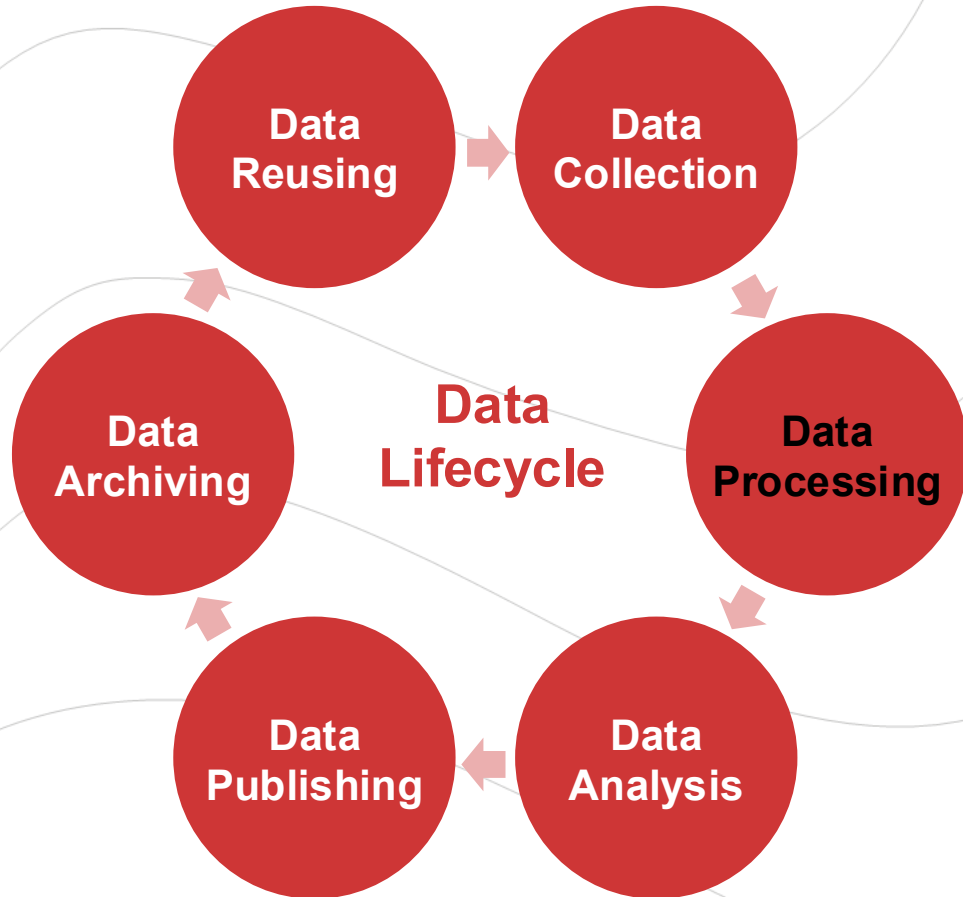
Research data lifecycle

Source Data

Collected/produced “raw data”
Reused data from a database/repository



Research data lifecycle



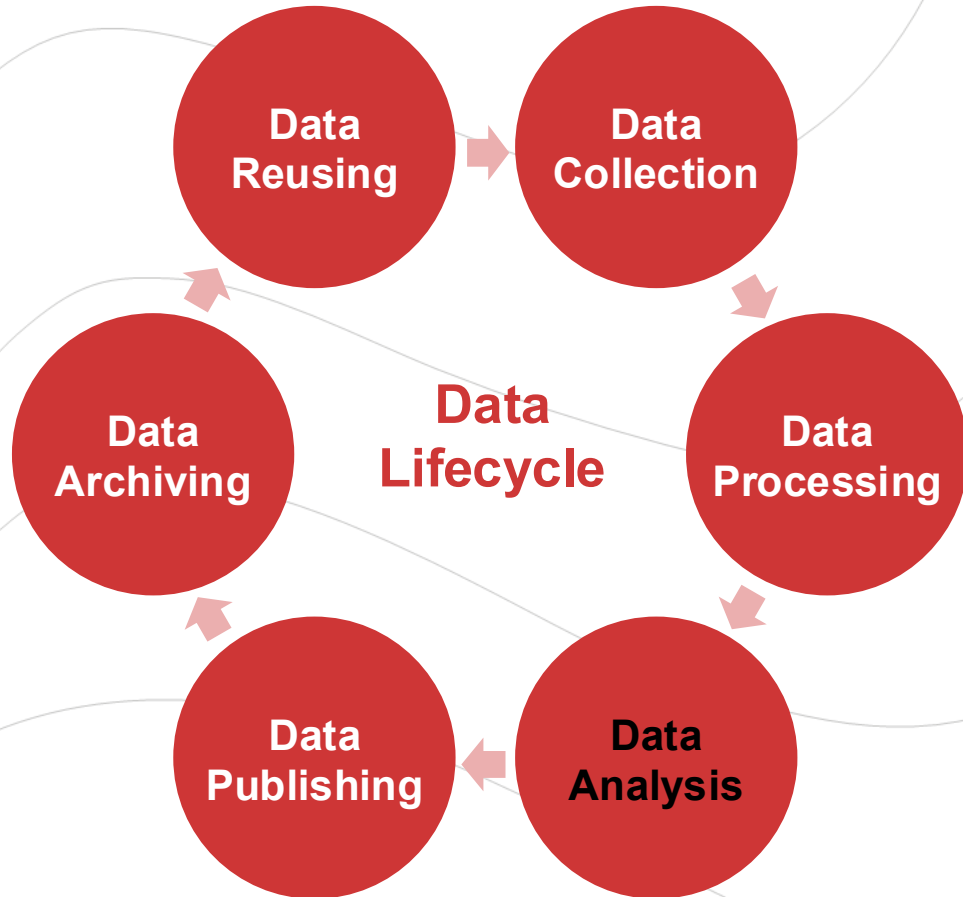
Source Data

Collected/produced “raw data”
Reused data from a database/repository

Data Processing

Transformation of raw data

Research data lifecycle



Source Data

Collected/produced “raw data”
Reused data from a database/repository

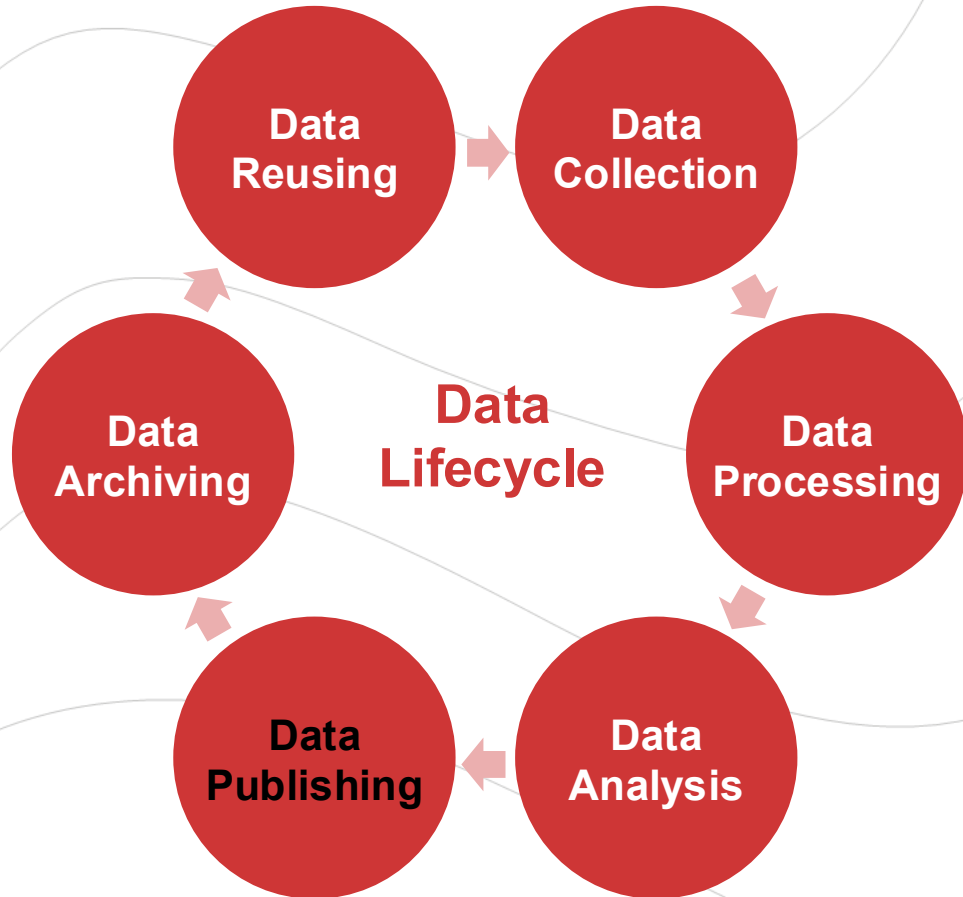
Data Processing

Transformation of raw data

Data Analysis

Data interpretation
Generation of results and outputs

Research data lifecycle



Source Data

Collected/produced “raw data”
Reused data from a database/repository

Data Processing

Transformation of raw data

Data Analysis

Data interpretation
Generation of results and outputs

Data Publishing

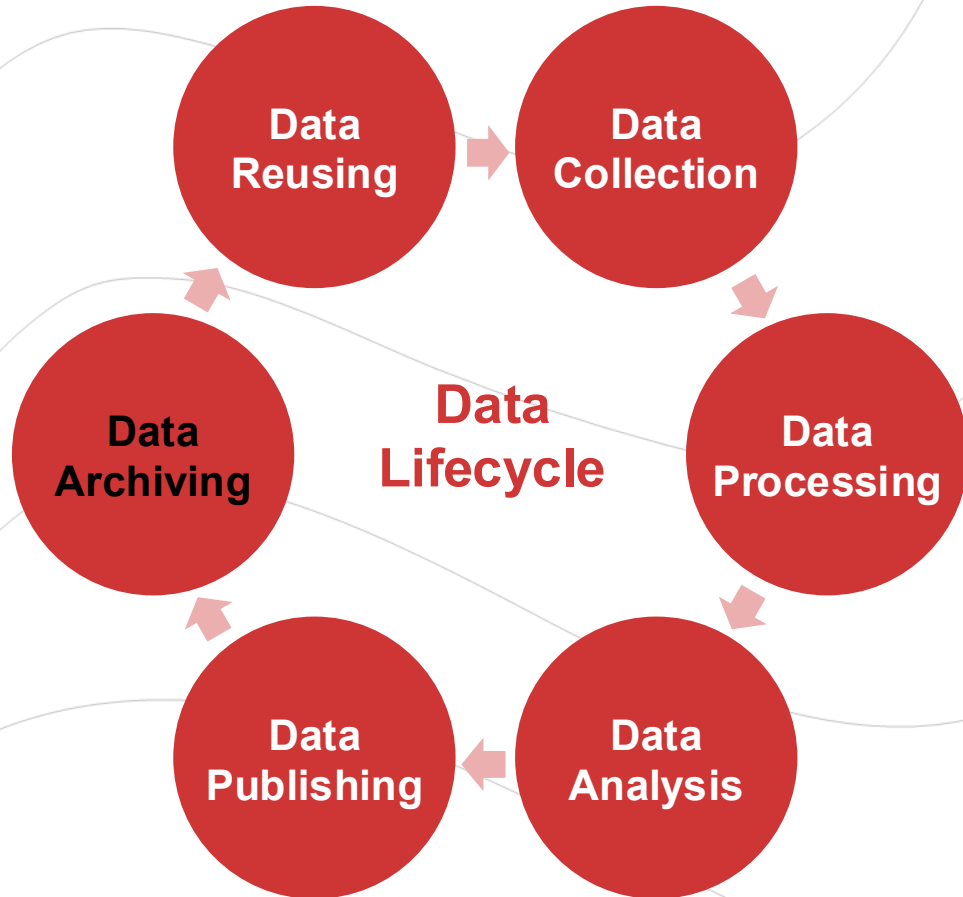
Journal article

Manuscript + supplementary information

Databases/repositories

Data underlying publication
Separate datasets

Research data lifecycle



Source Data

Collected/produced “raw data”
Reused data from a database/repository

Data Processing

Transformation of raw data

Data Analysis

Data interpretation
Generation of results and outputs

Data Publishing

Journal article

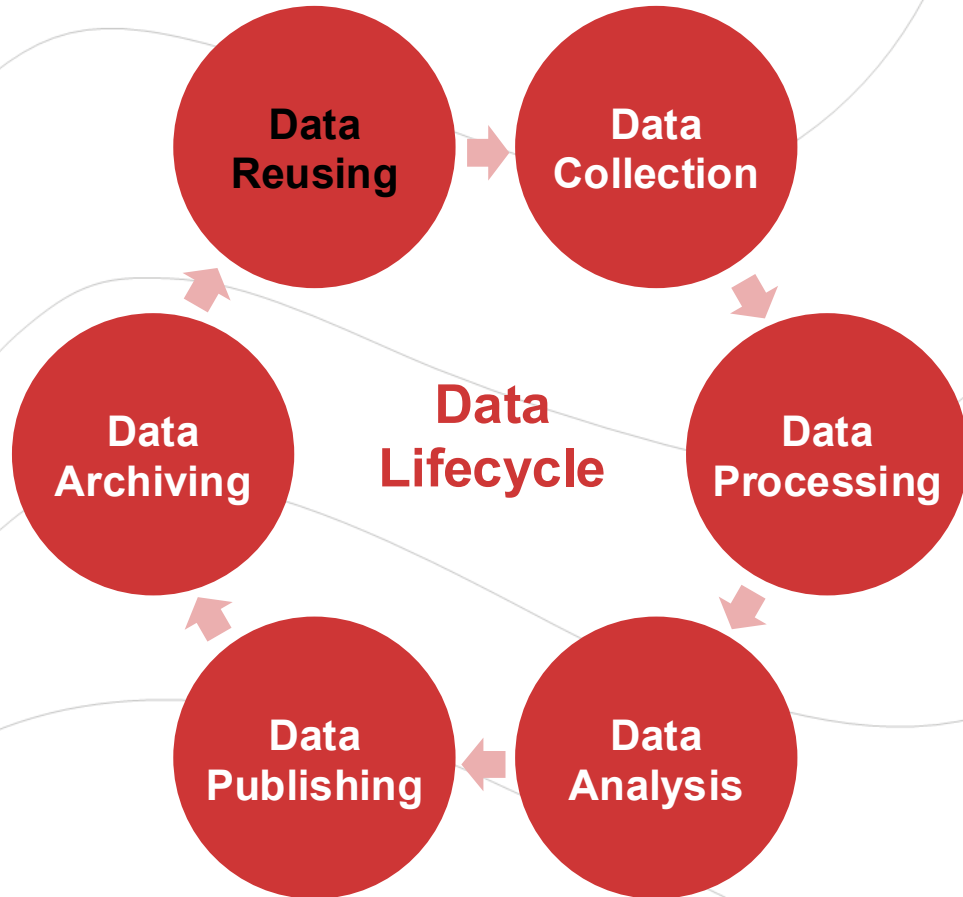
Manuscript + supplementary information

Databases/repositories

Data underlying publication
Separate datasets

Data Archiving

Research data lifecycle



Source Data

Collected/produced “raw data”
Reused data from a database/repository

Data Processing

Transformation of raw data

Data Analysis

Data interpretation
Generation of results and outputs

Data Publishing

Journal article

Manuscript + supplementary information

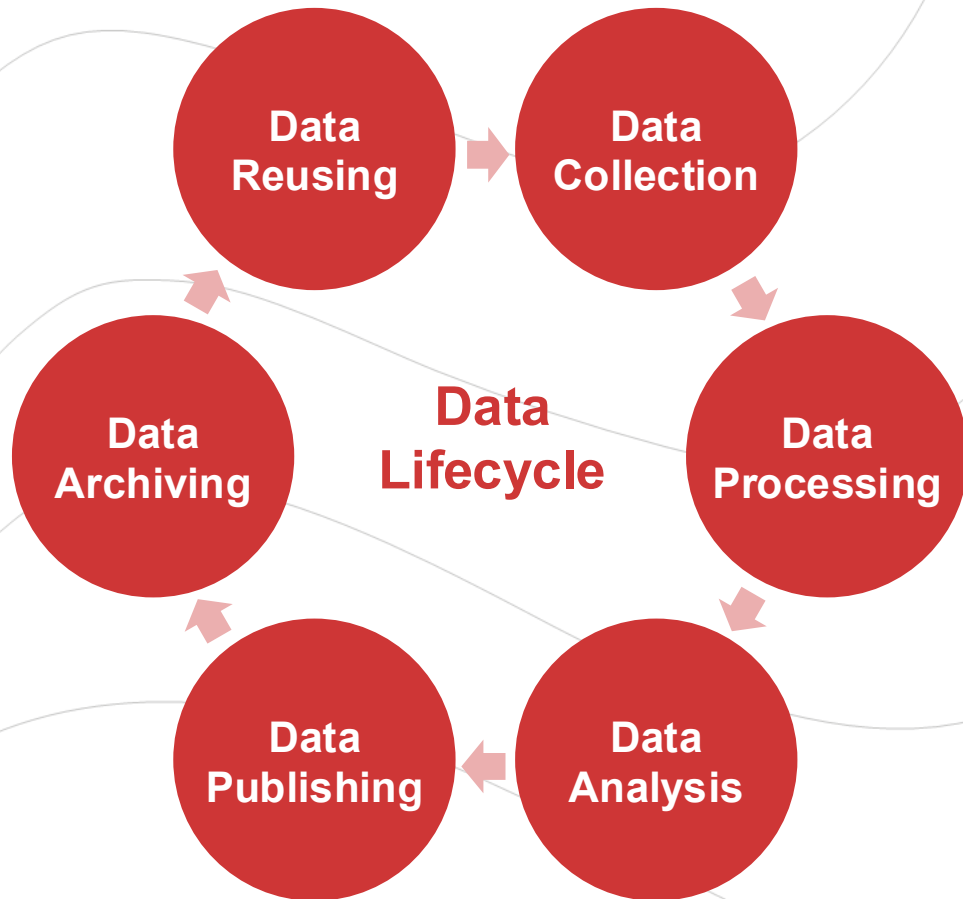
Databases/repositories

Data underlying publication
Separate datasets

Data Archiving

Data Reusing (registries, repositories)

RDM strategies



Organizing

- Directory structure
- Formats, names, versions

Documentation

- Data description
- Experimental details
- Decisions made
- Metadata

Storage

- Backup
- Long-term preservation

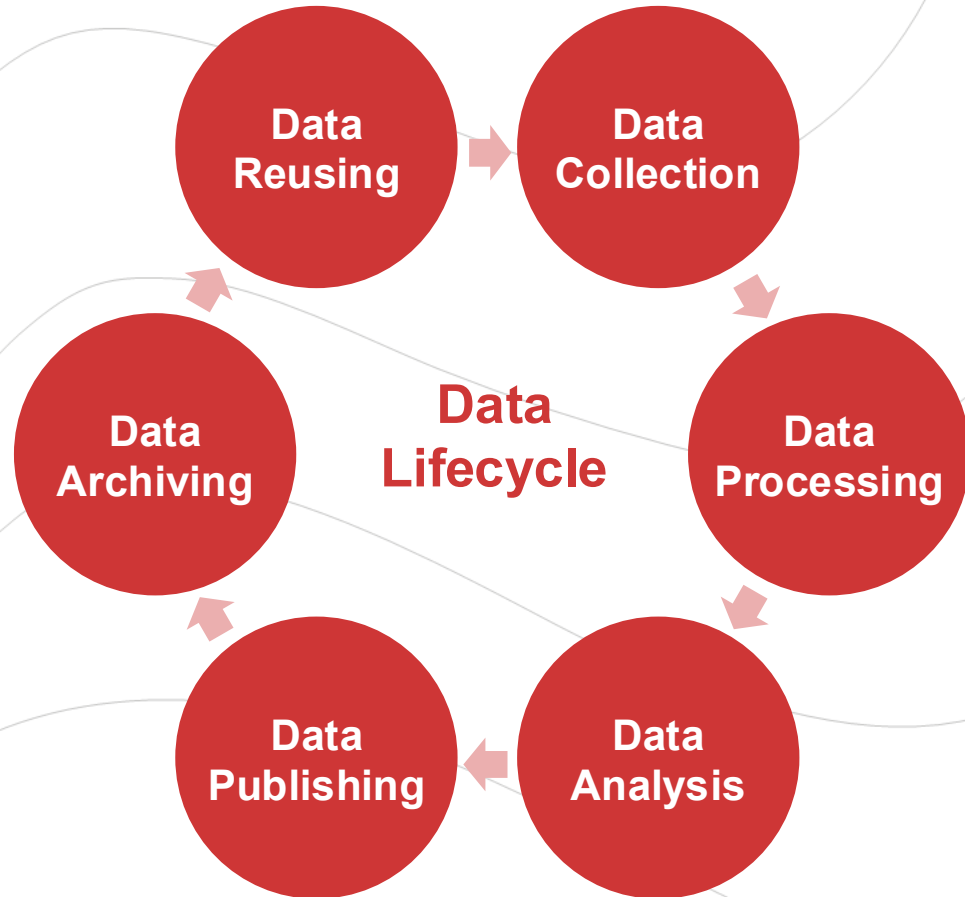
Data access

- Access rights (open, restricted)
- Licenses

RDM strategies

Plan

Generate ideas
Design research
Funding proposal



Organizing

Directory structure
Formats, names, versions

Documentation

Data description
Experimental details
Decisions made
Metadata

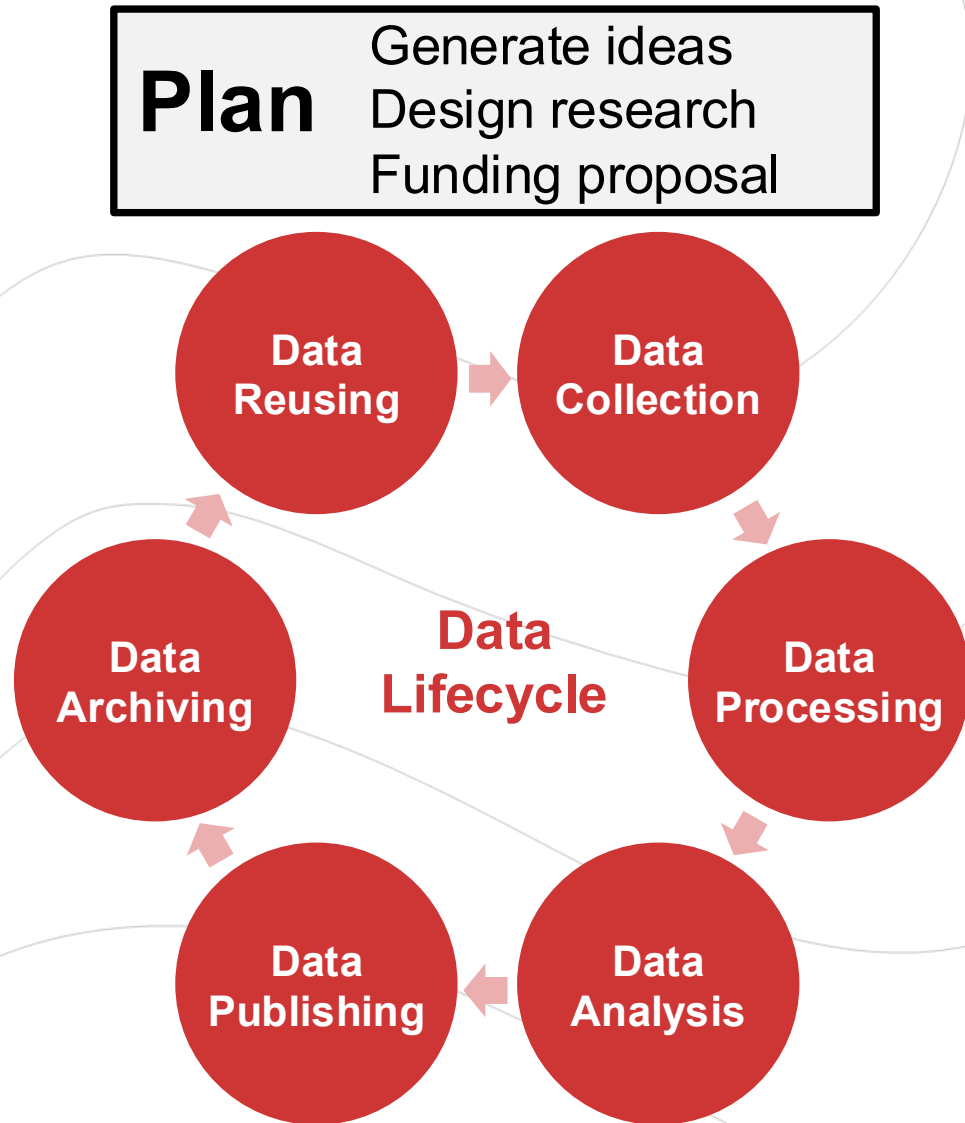
Storage

Backup
Long-term preservation

Data access

Access rights (open, restricted)
Licenses

Examples of research data requirements and policies



Funding agency policies

- Open Access policy
- Data management plan

Legal and ethical requirements

- National and European legislation
- Ethical framework for researchers
- Personal data protection
- Intellectual property rights
- Commercial use of data

Institutional policies

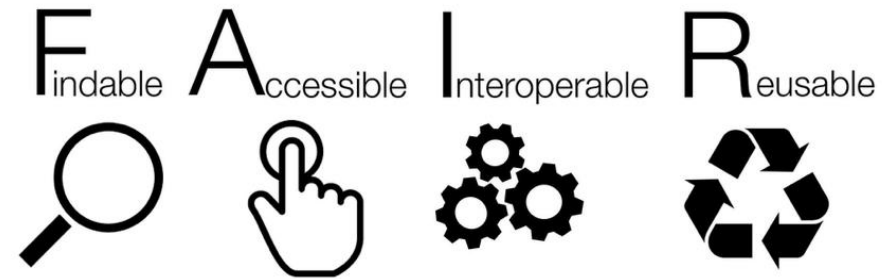
- RDM policy
- Codes of conduct and ethics
- Data protection
- Partnership agreement (for collaboration)

Journal & Publisher policies

- Data sharing policy

FAIR principles

FAIR principles



Source: [SangyaPundir, FAIR data principles, CC BY-SA 4.0](#)

Findable

- Metadata
- Persistent Identifiers (DOI, ORCID, ROR, IGSN...)
- Registration and indexing in searchable repository

Accessible

- Free and open metadata
- Metadata available even when data are not available

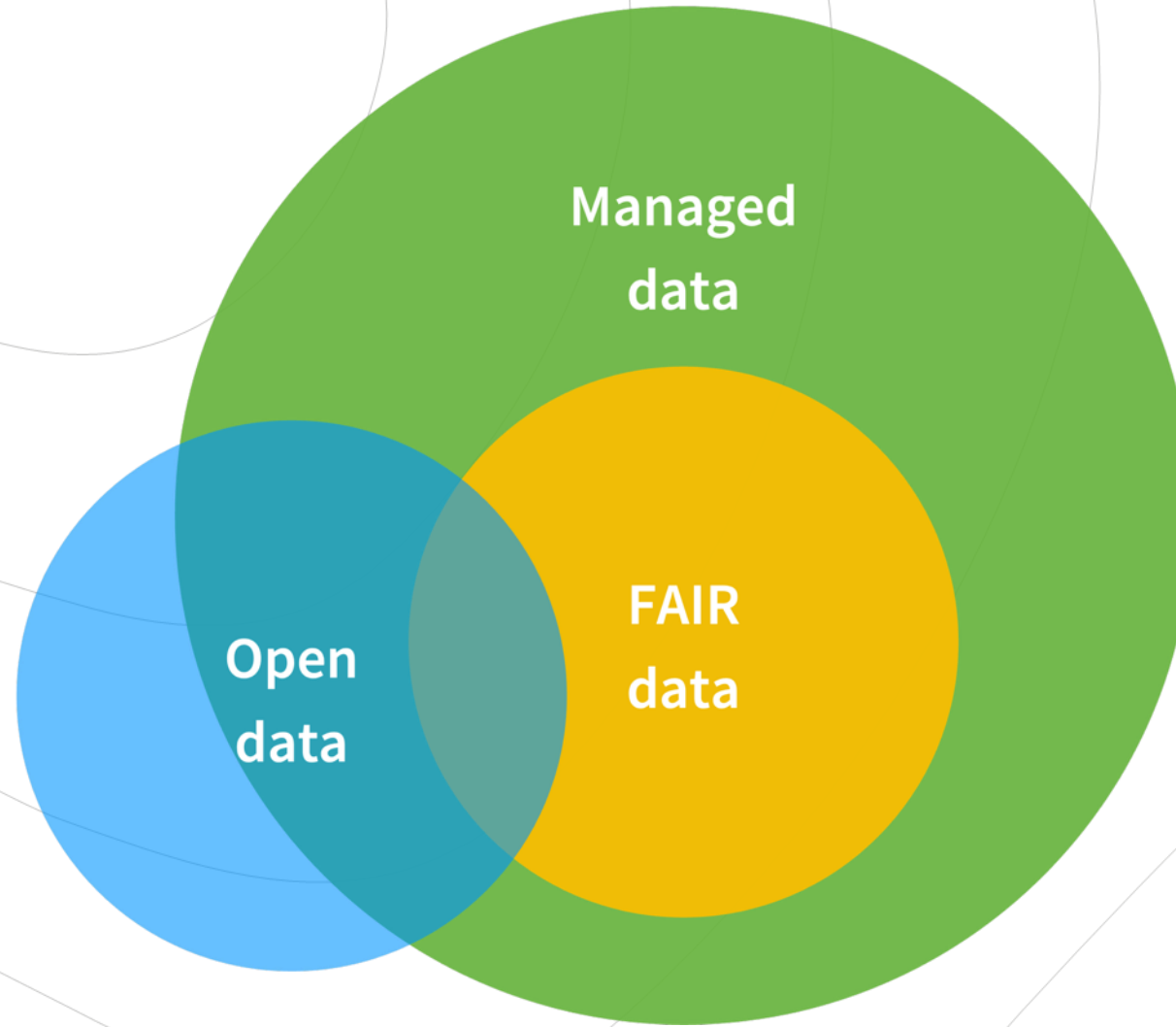
Interoperable

- Widely used language
- Preferred formats
- Vocabularies and ontologies

Reusable

- Rich description (Read Me File)
- License
- Field/Community standards

Open vs. FAIR Data



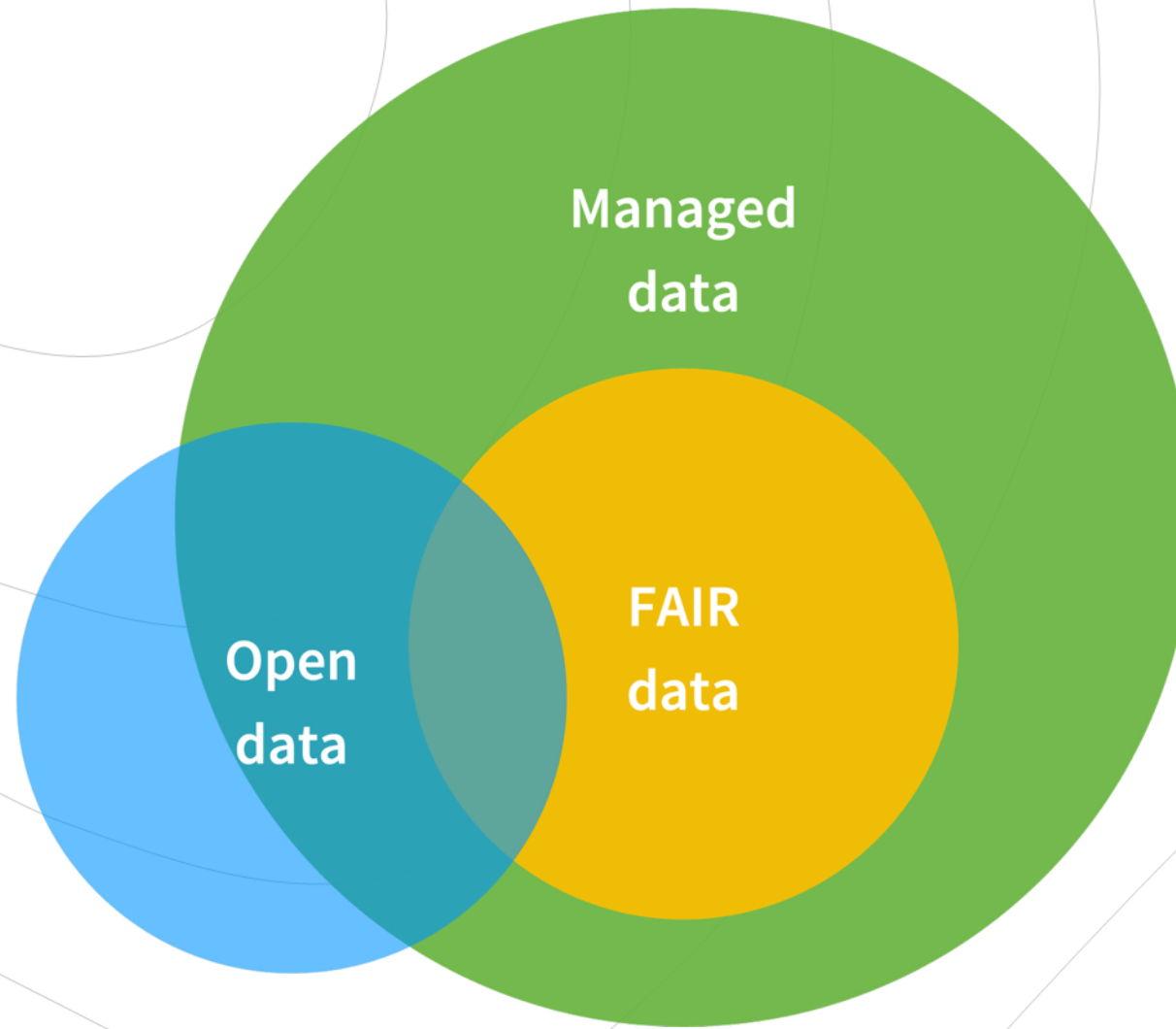
What (not) to publish?

- Personal data
- Sensitive data
- Protected by legitimate interest:
 - Intellectual Property
 - Commercial Interests

What to do with these data?

- Can the data be:
 - anonymized?
 - shared with informed consent?
 - shared later (embargo)?
 - shared only with selected researchers (controlled access)?
- If not, we can still usually share at least **metadata**

Open vs. FAIR Data

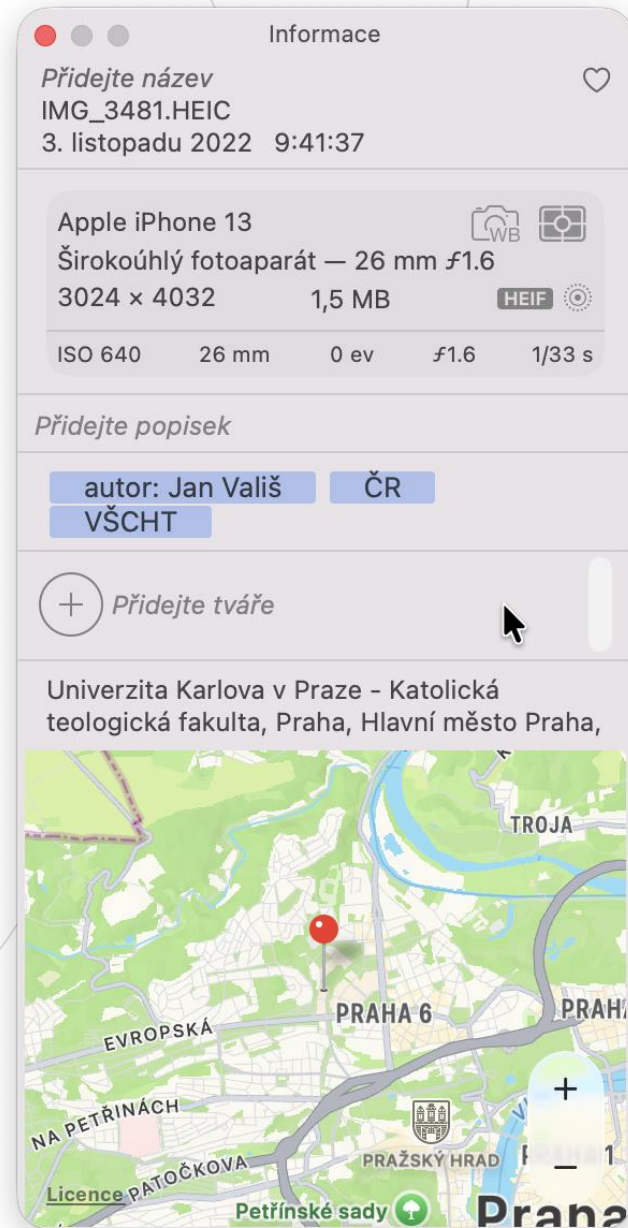


Metadata

Metadata – theoretically

- Documentation describing data = “**data about data**”
- **Machine readability** → search & retrieval of data from repository
- Metadata: who, what, when, where, why, how
- Different types
 - descriptive (title, author, date, keywords ...)
 - administrative (file formats, owner, license ...)
 - field-specific (spectral resolution, CAS number of molecule, species ...)

Metadata – practically



Metadata – models

- Universal
 - Dublin Core
 - DataCite
 - ...
- Field-specific
 - Open Geospatial Metadata Standard
 - Chemical Research Object Framework
 - ...

DataCite

- Mandatory
ID, Creator, Title, Publisher, Year, Type
- Recommended:
Subject, Contributor, Date, Related ID, Description, Geo Location
- Elective:
Language, Alternative ID, Size, Format, Version, Rights, Funding, Related Item

Quiz time!

Select examples of research data:

- Infrared spectrum of a sample.
- Grant budget submitted to grant provider.
- Recording of a structured interview.

Select truthful statements regarding metadata:

- Metadata help in finding research data.
- Metadata should be machine-readable.
- Must be human readable.

Data Management Plan (DMP)

DMP – plan before project starts

Administrative

- People involved (qualifications, training, roles & access, ORCID)
- Requirements (institutional, funder ...) & Support
- Budget

Instruments

- Access?
- Documentation?

Software

- Capture/processing/analysis workflows?
- Access to proper software?
- Use of open file formats?

DMP – plan before project starts

Size

- Enough storage for (captured/processed/analyzed) data?
- Many small files or fewer large files?

Backup

- How and where? (HDD, NAS, on/off-campus server)
- Encryption and access control?

Archiving

- What to archive?
- For how long?

DMP – plan before project starts

Legal & Ethical aspects

- Collaboration and services
- Personal/sensitive data
- Ethical committee approval?
- Informed consent?

Copyright License

- How are we legally bound?
- How do we want to license our results?

Publishing

- Can we publish data?
- Is there any domain-specific repository?

Tools to help you with DMP

- Data Stewardship Wizard (FAIR Wizard, e.g. CAS) (Free-to-use on-line instance for researchers with Czech eduID)
- DMP Online
- Argos

The screenshot displays the DS Wizard interface for a user named Albert Einstein. The main content area is titled 'My Experiment' and shows a questionnaire for 'IV. Processing data'. The questionnaire includes a 'Current Phase' dropdown set to 'Before Submitting the Proposal' and a list of chapters. The current chapter, 'IV. Processing data', contains several questions, including 'Will you be using a shared working space to work with your data?' and 'Will this work space be run by dedicated specialists?'. The interface also features a sidebar with navigation options like 'Users', 'Knowledge Models', 'Projects', 'Documents', and 'Settings'. A right-hand panel shows a version history for the current question, with the current version being 1.0.0. The user's profile and name are visible at the bottom left of the interface.

source: <https://ds-wizard.org/>

Data Reuse

Benefits

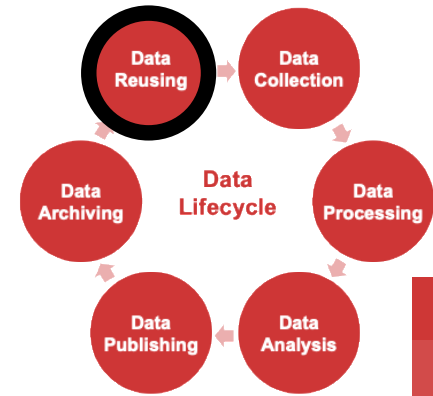
- Resource efficiency (time, money, equipment)
- Access to hard-to get data
- Larger datasets

Perils

- Different collection methodologies
- Quality?
- Incompatible formats
- Ethical aspects

Solutions

- Implementation of FAIR principles
 - Community standards
 - Preferred formats
 - Unified ontologies/controlled vocabularies
 - Permissive licensing
- Data harmonization
 - Cleaning
 - Processing
 - Transformation
 - Normalization



Data Reuse

Good practice

- Use of persistent identifiers (DOI, IGSN)
- Use of ethically/legally indisputable
- Reusing established workflows
- Reusing data from other files

Questionable practices

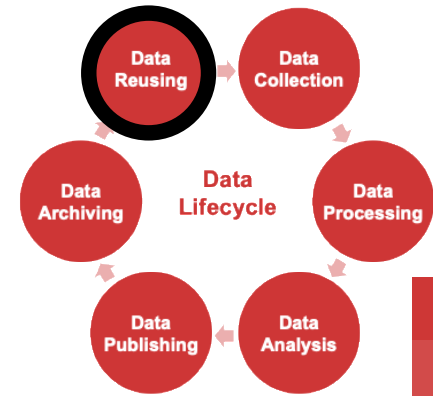
- Assuming data are correct
- Use without/in breach of license

Reference datasets

- Trusted data
- Peer-reviewed, validated (e.g. NIST Webbook of chemistry)

Non-reference datasets

- Experimental/observational data
- Variable quality



Data Collection

Organization

File naming

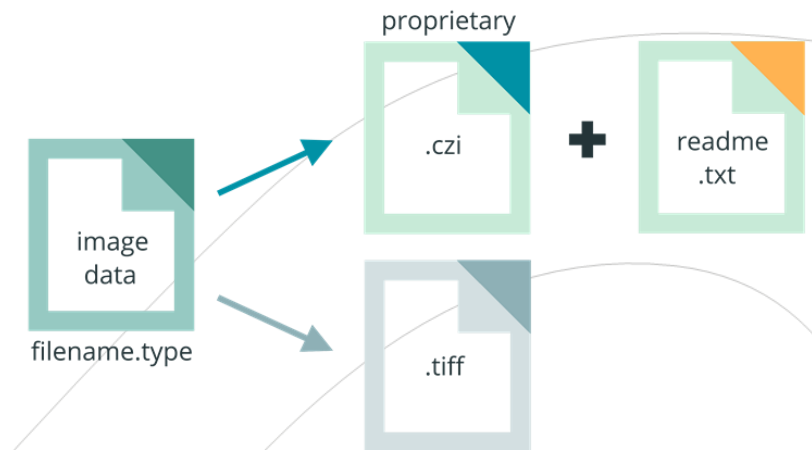
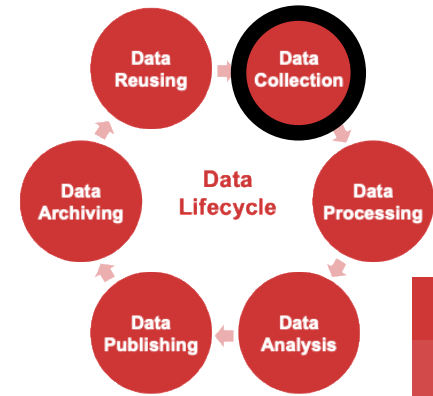
Metadata

File formats

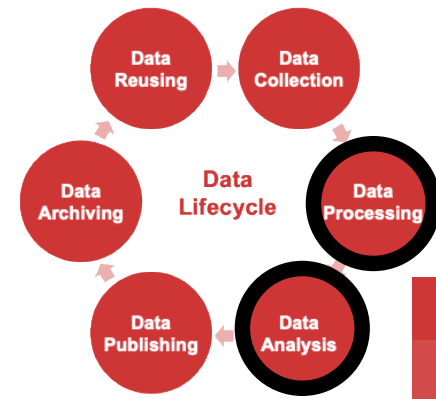
- Preferred vs. popular
- Open vs. proprietary
- If proprietary necessary, include info in ReadMe on how to open/use the file (SW, workflow etc.)
- E.g., TIFF (~~JPEG~~), PDF/A (~~DOCX~~), CSV (~~XLSX~~)

Preferred format

- archiving-friendly
- open
- well-documented
- human & machine readable
- lossless compression
- not dependent on a specific SW

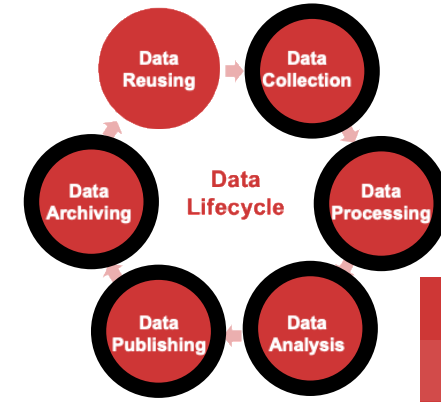


Data Processing & Analysis



- Algorithmic workflows preferred (reproducibility, efficiency)
- Always work on a copy: original raw data need to stay untouched
- Minimize work with sensitive data
(e.g. split sensitive/personal (meta)data and reconstitute them only using a key available to a small group of people)

Data Storage



By stage

- During collection/processing/analysis (on-campus/off-campus)
- After analysis – **publishing** (e.g. repository)
- After the project finishes – **archiving** (locally, repository, or cold-storage?)

Security aspects

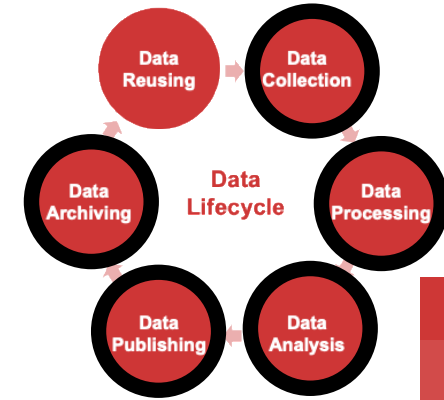
- Encryption
- Training
- Access restriction (by role, read/write etc.)

Financial aspect

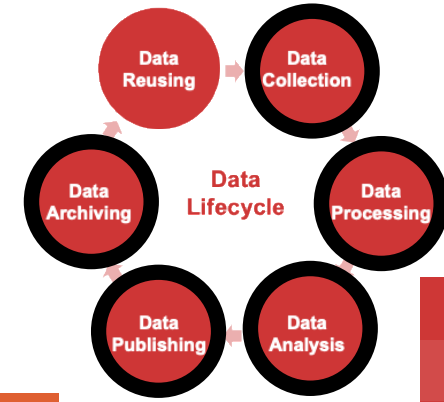
Try to guess: how much should your budget be if you need to store 1 TB of data for 10 years?

Assume that a large number of small files is going to be generated. And even though the data will only be archived at the storage facility, the researcher wants to be able to retrieve the whole dataset in one day. As a good custodian, the researcher requires 24/7 support and will tolerate only 1 day of downtime per year. Due to the sensitivity of the data, controlled access will be required

- A. 10 000 CZK
- B. 50 000 CZK
- C. 100 000 CZK
- D. 500 000 CZK
- E. 1 000 000 CZK
- F. 5 000 000 CZK
- G. 10 000 000 CZK



Data Storage – price estimation



DSW Storage Costs Evaluator

Total costs:
2 261 €

TB costs per year:
452 €

Result details
▼

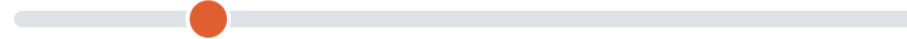
Volume



500

GB

Lifetime



10

years

Detailed storage properties ▼

Data Publishing

License

- Exclusive
- Non-exclusive
 - SW specific (GPL, MIT, Apache)
 - Creative commons (CC)

0 - Public domain

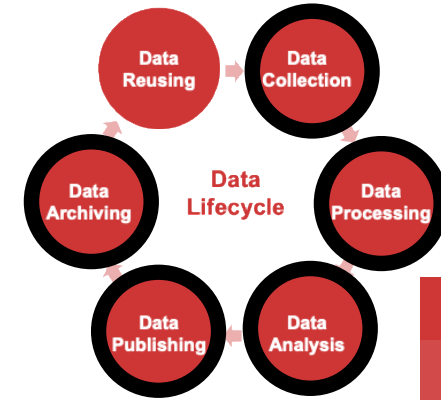
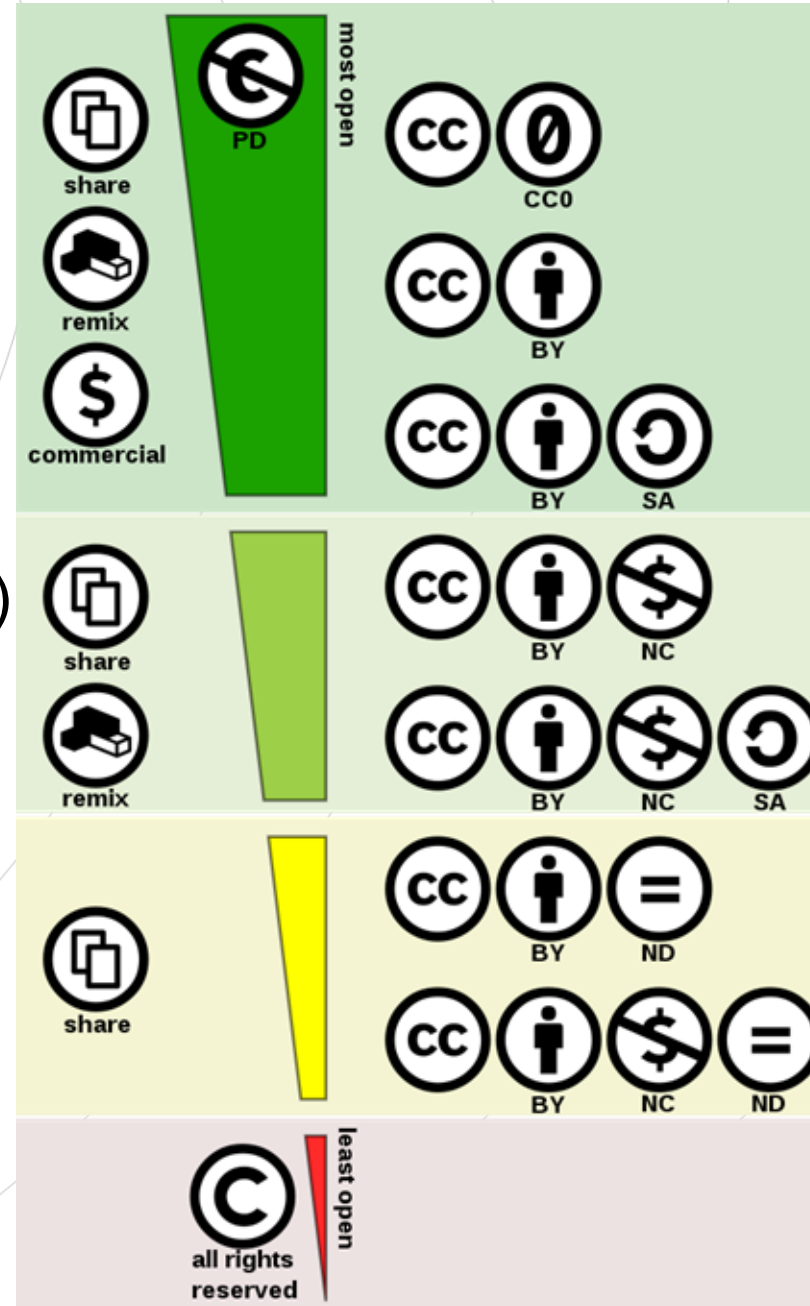
👤 **BY** - By Attribution

= **ND** - No Derivatives

🚫 **NC** - Non-Commercial

♻️ **SA** - Share Alike

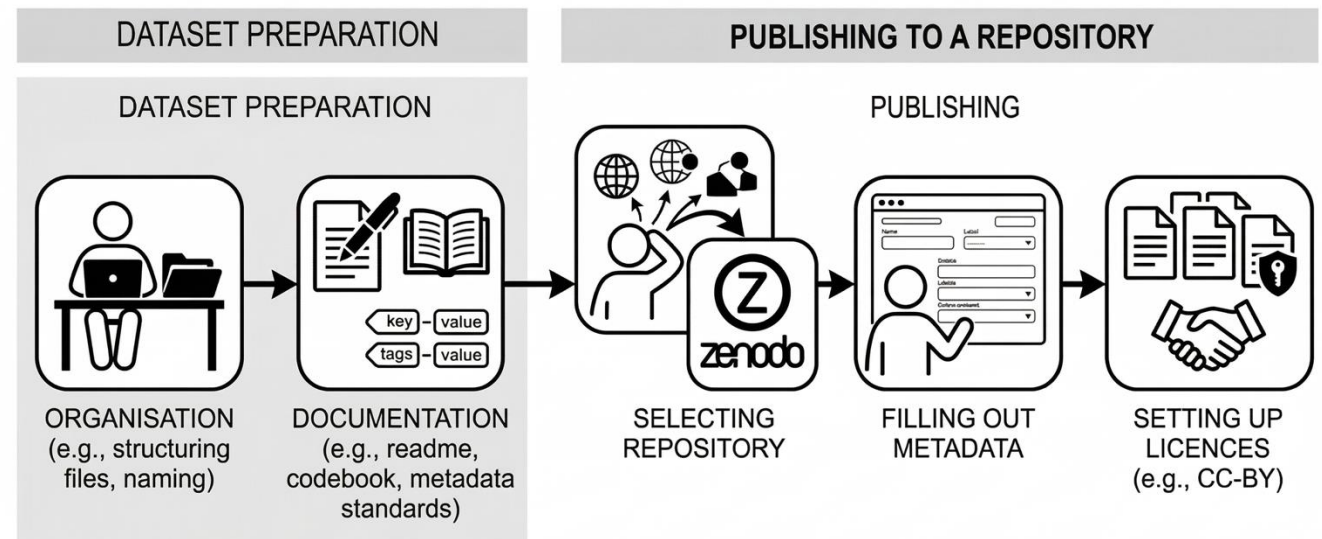
Source: creativecommons.org/



How to prepare a dataset and publish it?

Register for a brand-new webinar:

- April 22, 2026 (10:00 – 11:30)
- On-line
- Topics
 - Dataset preparation (e.g., organization, ReadMe)
 - Repository selection
 - Dataset deposition on Zenodo



Geminy (3 Flash, 2026-04-07), prompt: Create a flowchart depicting: dataset preparation (organization, documentation), and publishing to a repository (selecting a repository, filling out metadata, setting up licenses). Use Zenodo repository. The flowchart should be horizontal (left to right; 16:9) with minimal or no text - mainly pictograms. Use only black-and-white; simple lines. Style is professional/academic/technical.

Need to escape RDM?

Register for Data Horror Escape Room:

- April 29, 2026
(17:00 – 20:00)
- In-person, NTK

TALKS & SHARESPACE FOR PHD+

TOPIC DATA HORROR ESCAPE ROOM AT NTK

TIME APRIL 29, 2026, 17.00-20.00

OPEN DOOR 16.00

LOCATION NTK, PRAGUE

NTK
50°6'14.083"N, 14°23'26.365"E
Národní technická knihovna
National Library of Technology

CAD
THE CZECH ASSOCIATION OF DOCTORAL RESEARCHERS

SCREAM IF IT'S A DMP!

REGISTRATION

Where to learn more?

Your institution

- Library/Open Science Center
- Data Steward

NTK

- [RDM Guide](#)
- [Data Stewardship Course](#)
self-guided, finished by a test & certificate

Full courses

- [UISK, CUNI](#) (running, CZK 12 600)
- [Uni Wien](#) (running, CZK 72 000)
- [UCT Prague](#) (Apr 2026, CZK 1 000 + CZK 15 000)

The screenshot shows the NTK (National Technical University of Slovakia) website page for Research Data Management. The page features a search bar at the top, a navigation menu, and a main content area with the following sections:

- Research Data Management**: A section with an introductory paragraph about RDM and a list of links: [Research Data](#), [FAIR Principles](#), [Research Data Management](#), [Data Management Plan](#), [Data Resources](#), [Support](#), and [Resources](#).
- Research data**: A paragraph explaining that research data is information or material collected, used, or generated during the research process.
- Why Manage Research Data?**: A paragraph explaining the importance of managing research data properly to keep it secure and organized.

On the right side of the page, there is a "Your contact" section with two entries:

- Jan Veliš**: jan.velis@techib.cz, +420 220 442 072
- Karolina Podbucká**: datasupport@techib.cz, +420 771 269 626

Where to get help?

Self-study

- [EOSC CZ](#) – webinars (all)
- [NPOS](#) – Czech only (all)
- [RDMkit](#) (life-sciences)
- [OpenAIRE](#) (all)

Community

- Map
- Manuals
- [Discord](#)
- Meetings

MAPA DATASTEWARDŮ V ČESKÉ REPUBLICE

Statistiky Zapoj se! O aplikaci en

Hledat jméno NEBO instituce

Adrian Rosinec (Masarykova univerzita)	ⓘ
Aleš Mučka (Vysoká škola báňská - Technická)	ⓘ
Alexandra Šilerová (Masarykova univerzita)	ⓘ
Alžběta Karolyiová (Masarykova univerzita)	ⓘ
Aneta Pilátová (Masarykova univerzita)	ⓘ
Anna Kopecká (Západočeská univerzita v Plzni)	ⓘ
Anna Soldánová (Knihovna AV ČR, v. v. i.)	ⓘ
Antonio Feifar (Fyzikální ústav AV ČR, v. v. i.)	ⓘ

POČET DATASTEWARDŮ: 108
ZAPOJENÉ INSTITUCE: 43

Software a nástroje (skillset)
- Vše -

Agenda datastewarda
- Vše -

Vědní obor
- Vše -

Powered by Esri



Source: [EOSC CZ](#); Erik Dudinský & Lucie Skříčková

What to take home?

- Research data management:
 - is a helpful tool and a good scientific practice,
 - can save time and money and
 - ≠ Open Science, but sharing can create impact.
- Plan to get ahead of any issues, budget for expenses.
- Take advantage of community standards and on-line resources.
- Funding providers already demand and financially support OS and RDM practices, including DMP.
- If in doubt, reach out for support.

NTK Information Support Team

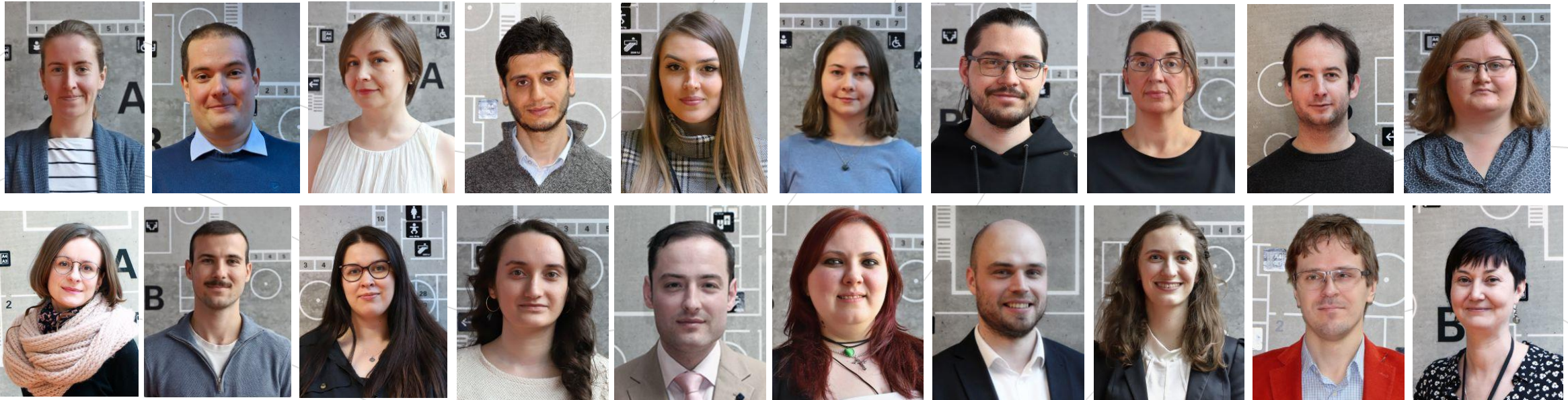
1) Schedule a free consultation with us

Don't be shy; our team includes doctoral candidates who understand the issues you face.

2) Attend another webinar

3) Explore on your own: Tutorials, AI tools for research or STEMskiller

4) Subscribe to our newsletter for updates on resources, writing support, publishing, research evaluation, and training opportunities.



Any questions? Contact us at info@techlib.cz

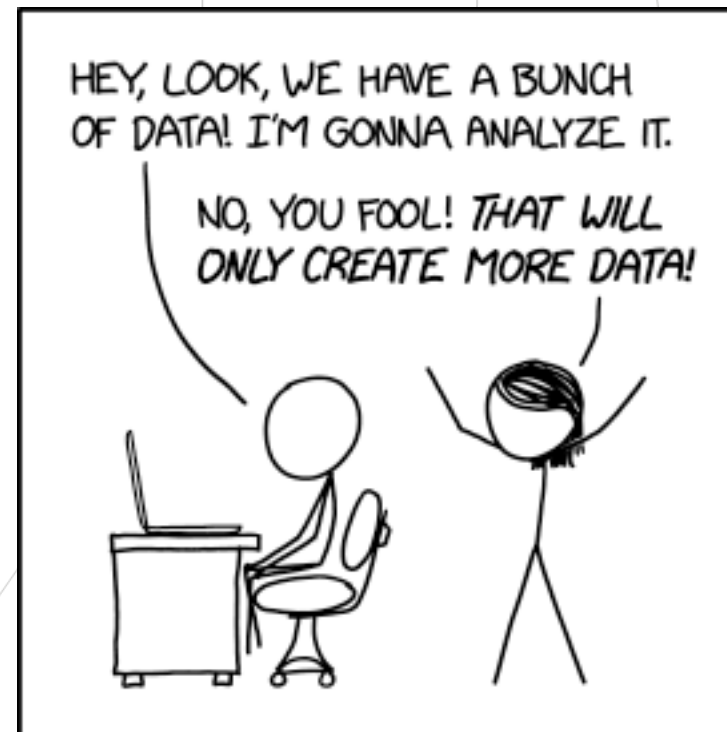
Contact

Jan Vališ ^{ID}

jan.valis@techlib.cz

NTK

50°6'14.083"N, 14°23'26.365"E
Národní technická knihovna
National Library of Technology



Source: <https://xkcd.com/2582/>;
Available via license: [CC BY NC 2.5](https://creativecommons.org/licenses/by-nc/2.5/)

Questions?