

28.11.2023

Research Data Management – The Basics

Bachofner, Anusch
Cantini, Federico
Felder, Fabian
Förster, Christian

These are your trainers today!



Federico Cantini

- Software Developer
- Technical Lead at Lib4RI



Fabian Felder

- Open Science specialist
- Group Leader IT services and E-resources at Lib4RI

Who are you and why are you here?

Copyright protected material.

<https://www.pexels.com/photo/group-of-people-standing-indoors-3184396/>

Learning Aims

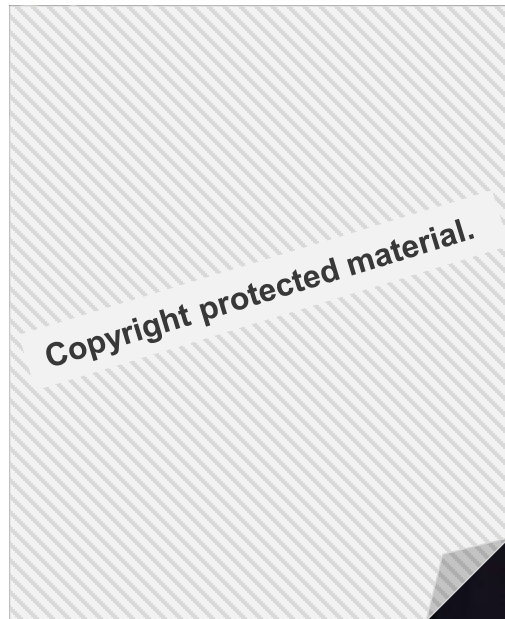
- Life cycle of research data
- Adequate metadata documentation for your code and data
- Storing and publishing data
- Using OpenBIS (ELN) and writing Data Management Plans (DMP)

Program

Topic	Speaker	Time
Introduction	Fabian Felder	9.00 - 9.15
Policies, Incentives & the Research Data Life Cycle	Fabian Felder	9.15 - 9.30
Collect & Store	Federico Cantini	9.30 - 10.00
Evaluate & Archive Share & Disseminate	Fabian Felder	10.00 - 10.10
Break		10.10 - 10.25
RDM Services & Support at Eawag	Christian Förster	10.25 - 10.45
RDM Services & Support at Empa	Anusch Bachofner	10.45 - 11.05
Plan & Design	Everyone	11.05 - 11.45

Why is data and associated metadata important?

Why is data and associated metadata important?



«We kill people based on metadata» (2014),
Michael V. Haden, director of CIA 2006-2009



Cham, J. G., «Scratch: A context-changing framework for contextualizing nano informatic structures» (2014), International Journal of Temporal Deflective Behaviour, 4 (1689), p. 432.

Why is data and associated metadata important?

No clean metadata

=

Limited access to Data



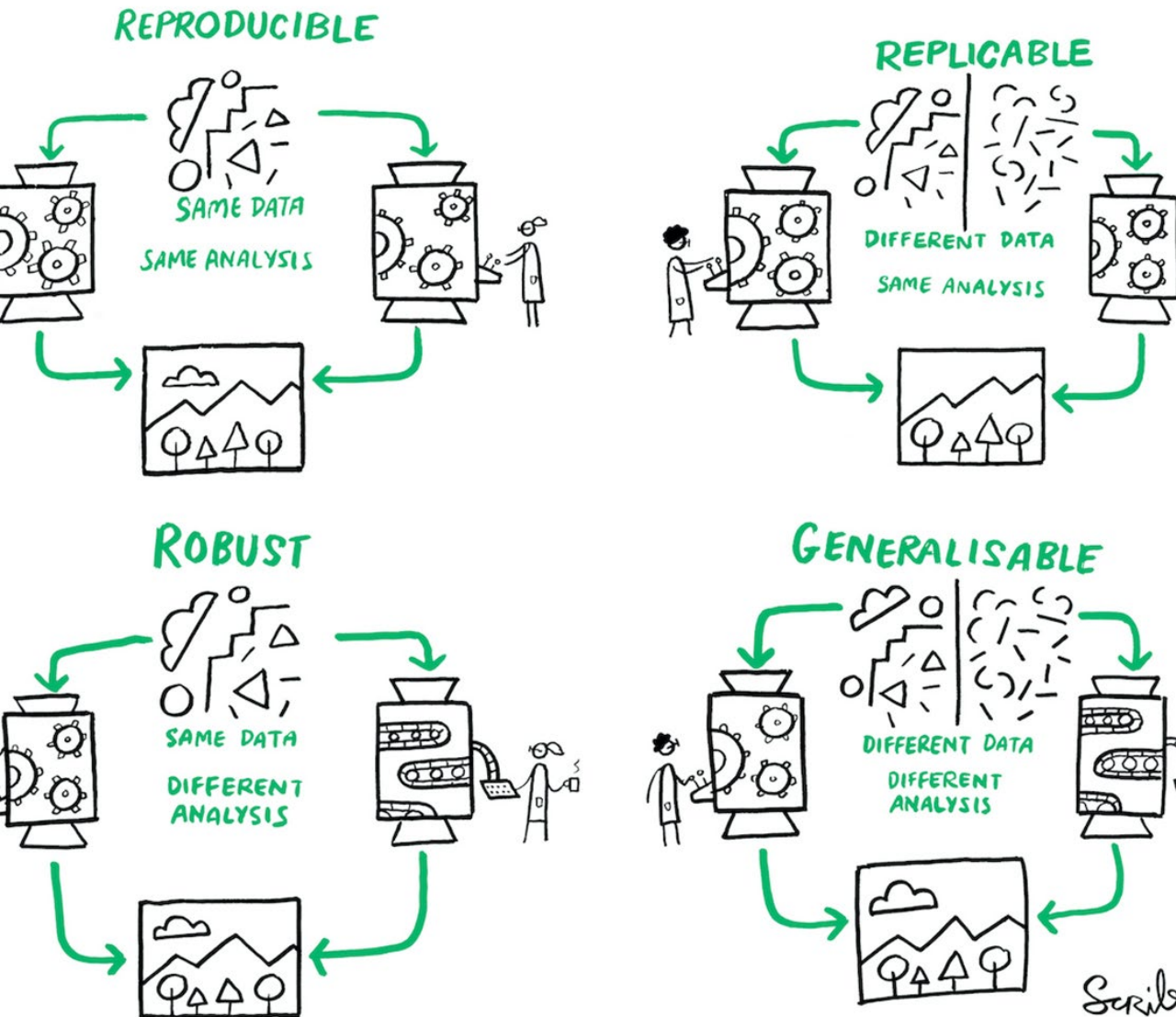
Source: www.fosteropenscience.eu/project

Why is data and associated metadata important?

**Proper metadata tagging
and
research description
is
time consuming**



Reproducibility



Reproducibility

Scriberia, “Reproducible Research”, *The Turing Way*, CC-BY, DOI: [10.5281/zenodo.3332807](https://doi.org/10.5281/zenodo.3332807)



A Handbook for Reproducible Data Science, <https://the-turing-way.netlify.app/welcome.html>

LEGO® Metadata for Reproducibility



Group A builds Car



**Group A documents
build**



**Group B rebuilds the
car**

Donaldson, Mary and Matt Mahon, «Lego: Metadata for reproducibility», 10.5281/zenodo.3685685.

LEGO® Metadata for Reproducibility

What matters?

What will you need to record?

Is there a way to record it automatically?

Which structure do you use? Or do you rely on a narrative expression?

How do you describe your materials?

Which formats do you use?

Is there a standard?

Donaldson, Mary and Matt Mahon, «Lego: Metadata for reproducibility», 10.5281/zenodo.3685685.

FAIR principles – A lot of Metadata

Findable

- F1 (Meta)data are assigned a globally unique and persistent identifier
- F2 Data are described with rich metadata
- F3 Metadata clearly and explicitly includes the identifier
- F4 (Meta)data are registered or indexed in a searchable resource

Accessible

- 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, when necessary
- A2. Metadata are accessible, even when the data are no longer available

Interoperable

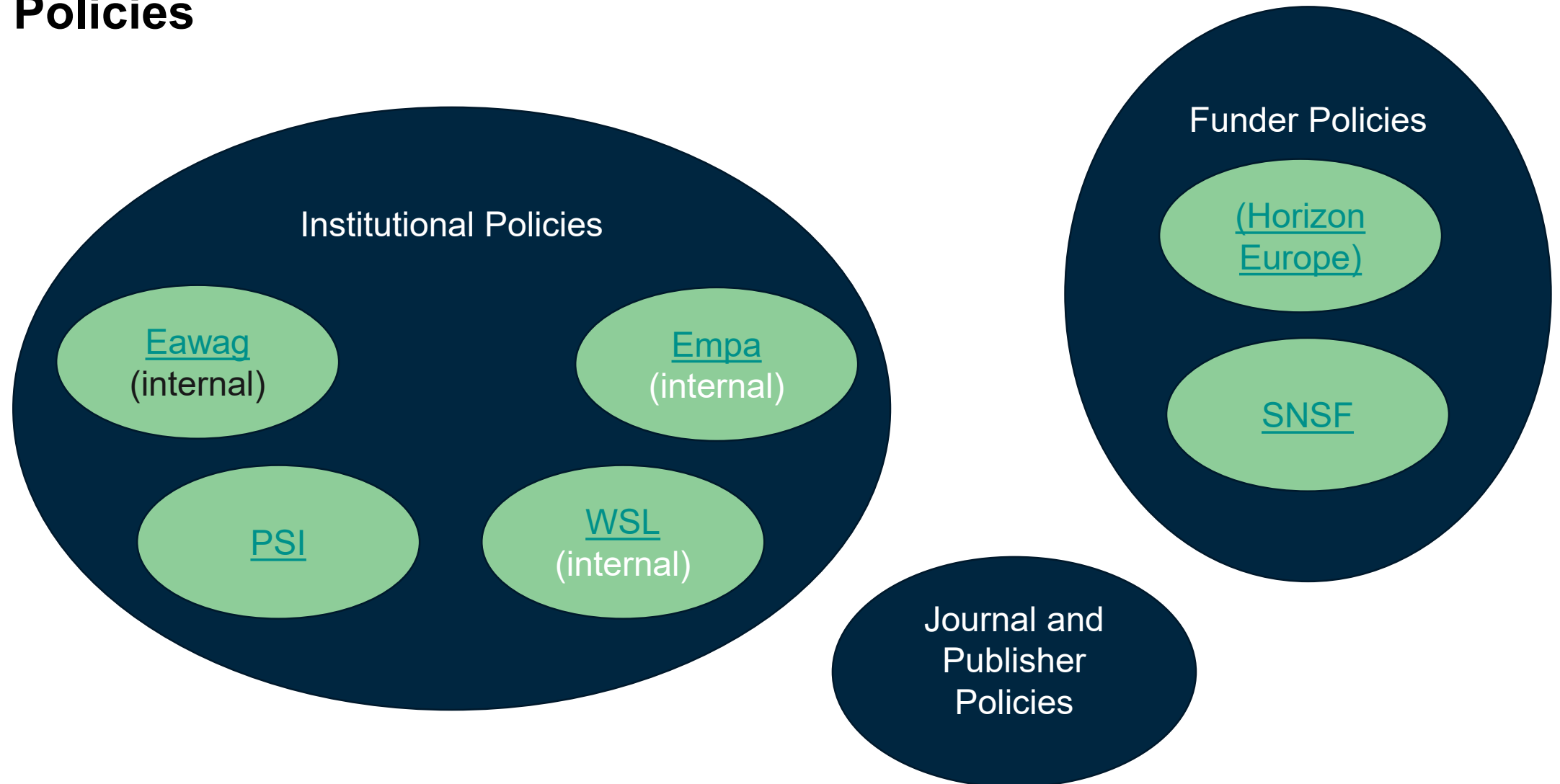
- I1 (Meta)data use formal, accessible, shared, and broadly applicable language for knowledge representation
- I2 (Meta)data use vocabularies that follow FAIR principles
- I3 (Meta)data include qualified references to other (meta)data

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 a detailed provenance
 - R1.3 (Meta)data meet domain-relevant community standards

Policies

Policies




Policies

Compliance




Project Manager/
Group Leader

DMP

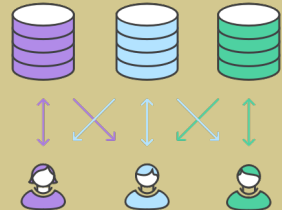


Required for
Funders

As open as possible,
as closed as necessary.

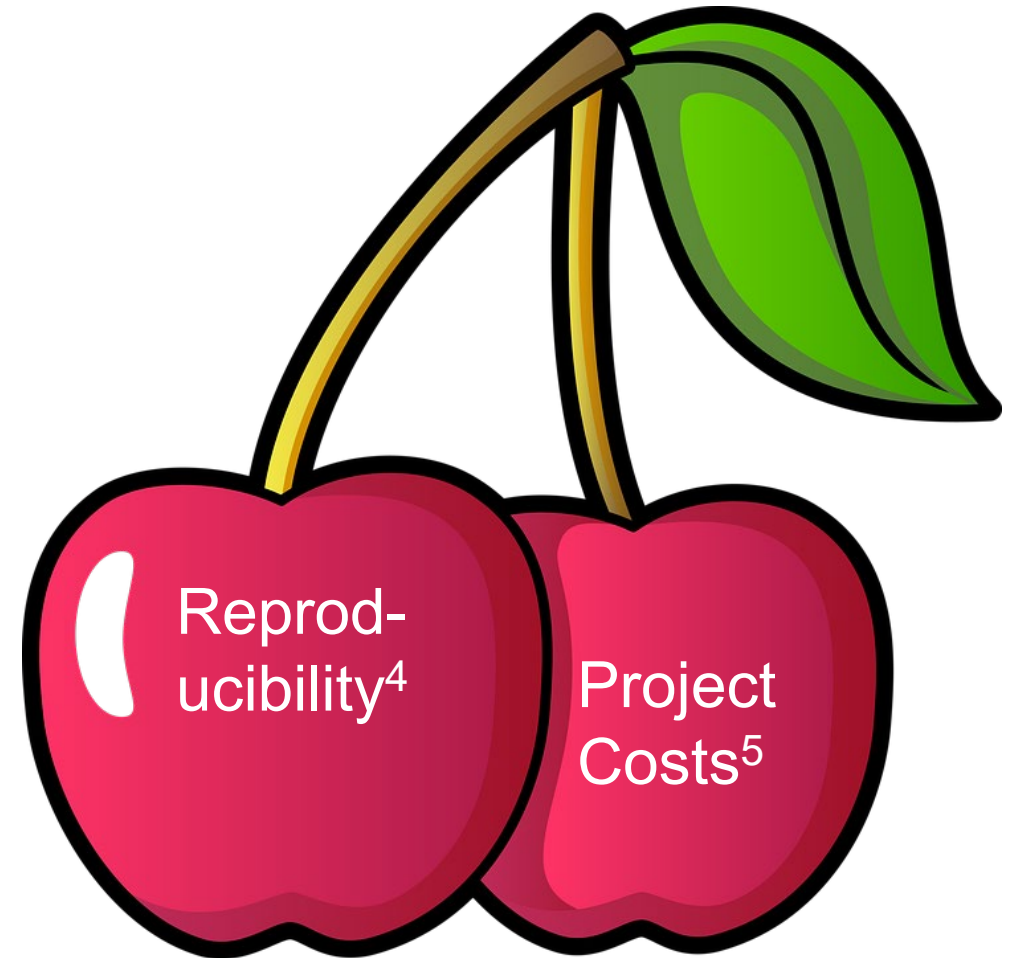
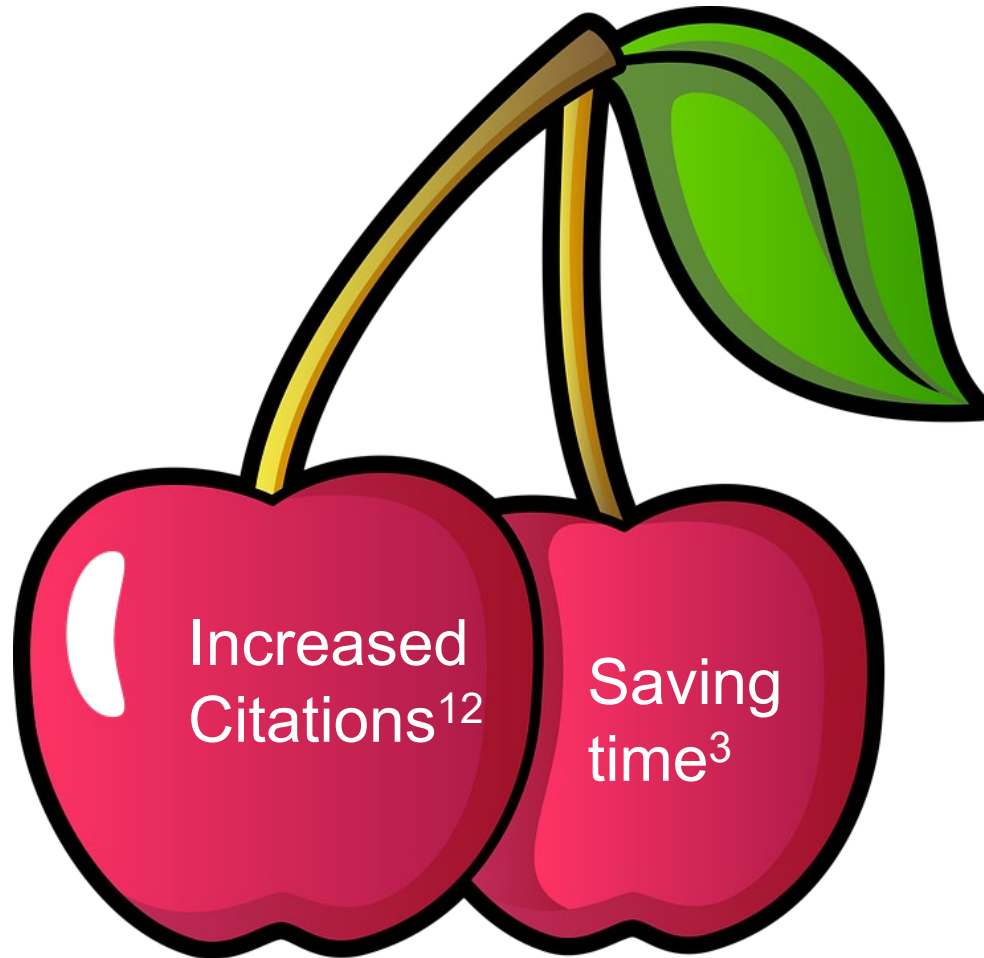


Repositories



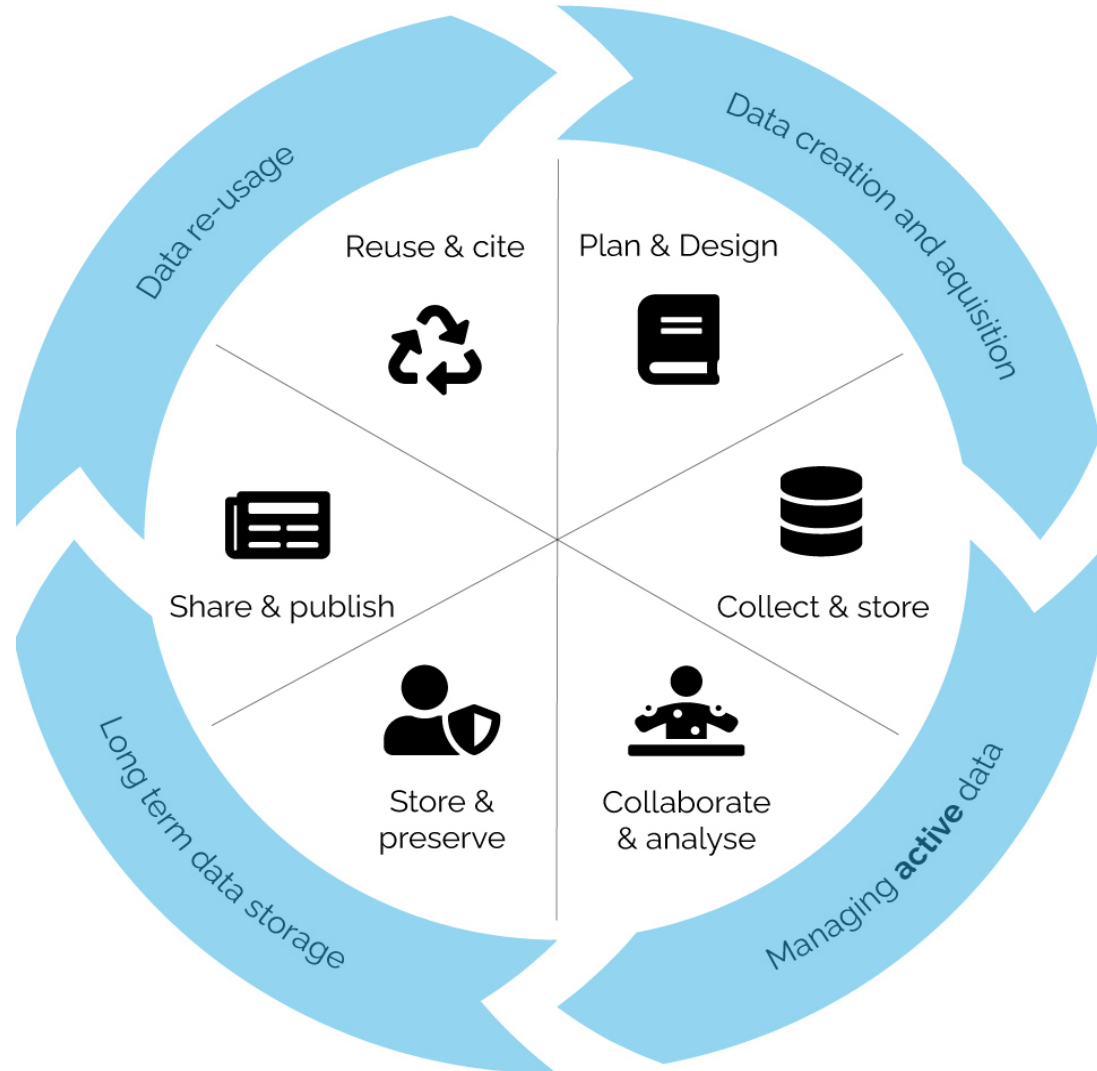
FAIR and non-commercial

Policies



Research Data Life Cycle

Research Data Life Cycle



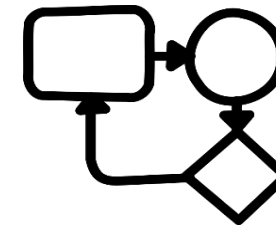
Collect & Store

Collect & Store

```
01010100 01101000
01101001 01101110
01101011 00100000
01100100 01101001
01100110 01100110
01100101 01110010
01100101 01101110
01110100 00101110
```

Data

observational, experimental, simulation...



Code

Applications, scripts...

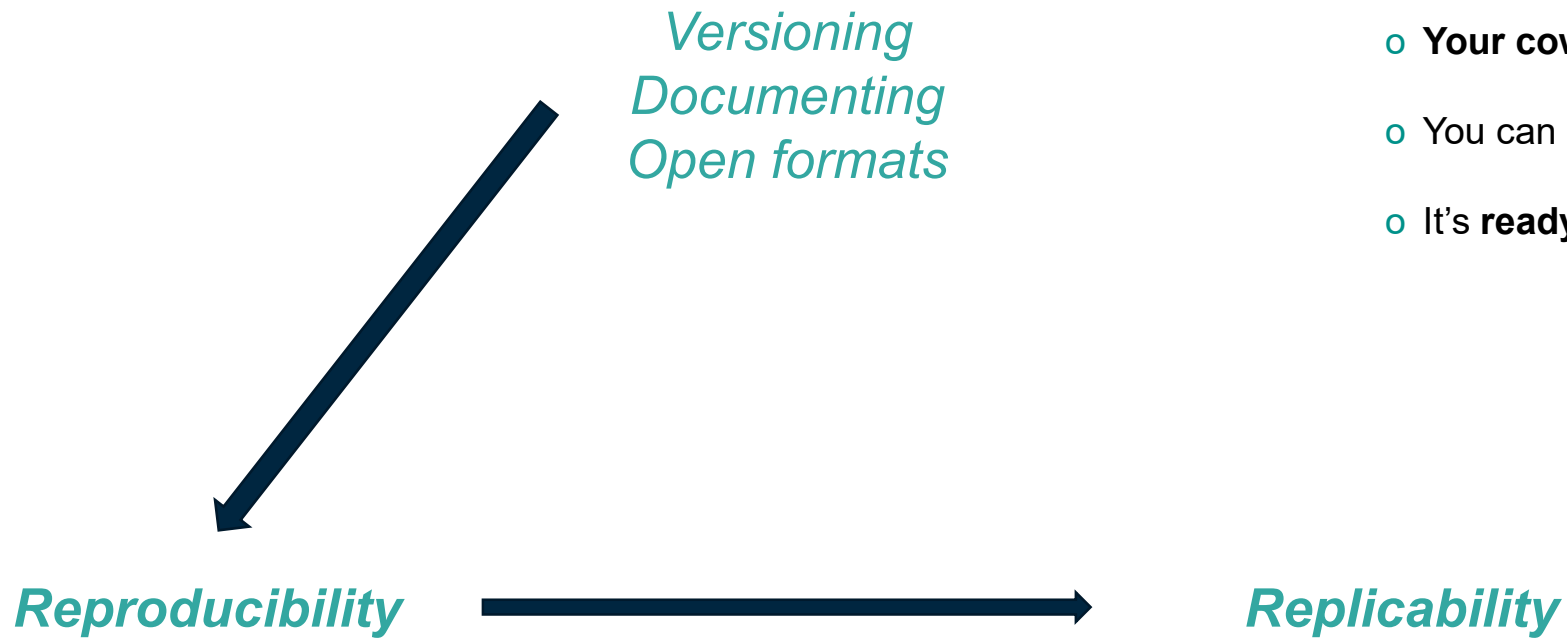


Metadata

Structured information associated with data (and code)

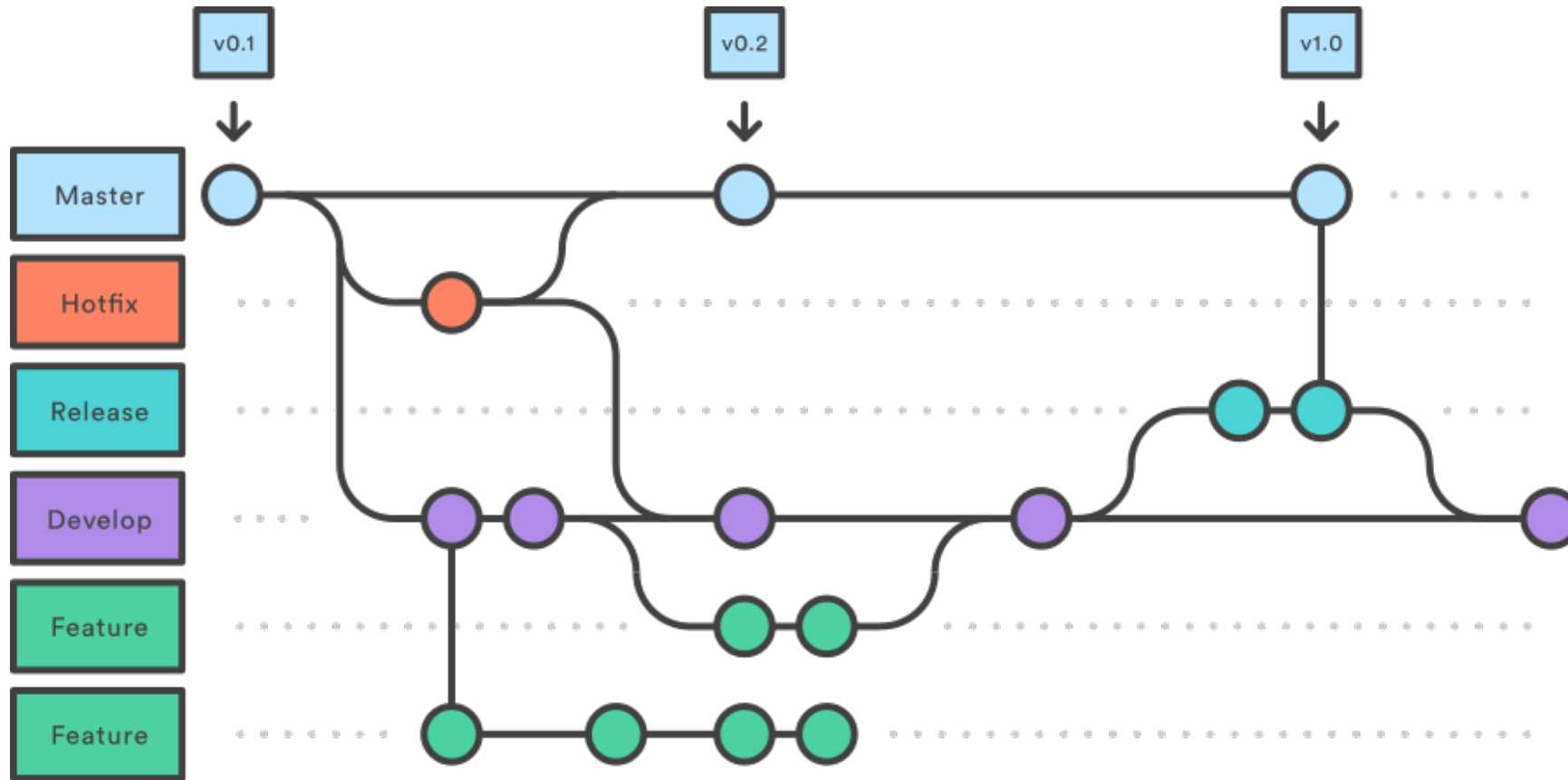
The Who, What, Where, Why & How of data

Collect and Store

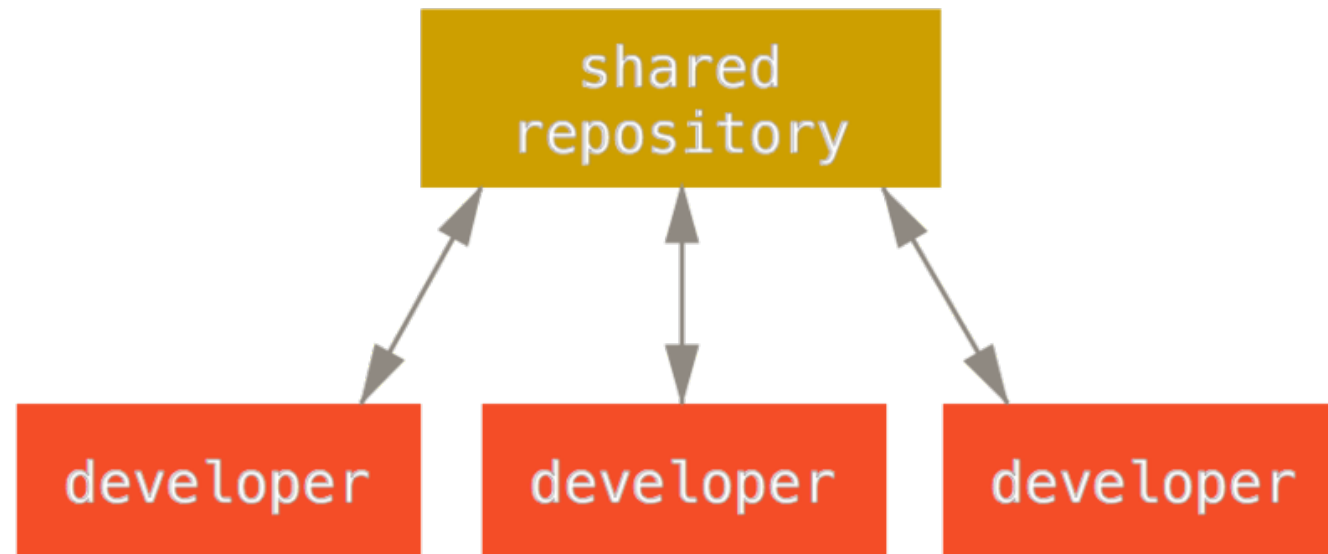


- **You** can find it
- **Your coworkers** can find it
- You can easily **share** it
- It's **ready** for **archiving/publishing**

Collect and Store: Software version control



Collect and Store: Software version control



Collect and Store: Software version control



<https://git-scm.com/>

- o CLI (*Command Line interface*)
- o GUIs (*Graphical User Interfaces*)
<https://git-scm.com/downloads/guis>



Workstation



Your own server



Gitea



GitLab



Internet



Codeberg



GitLab



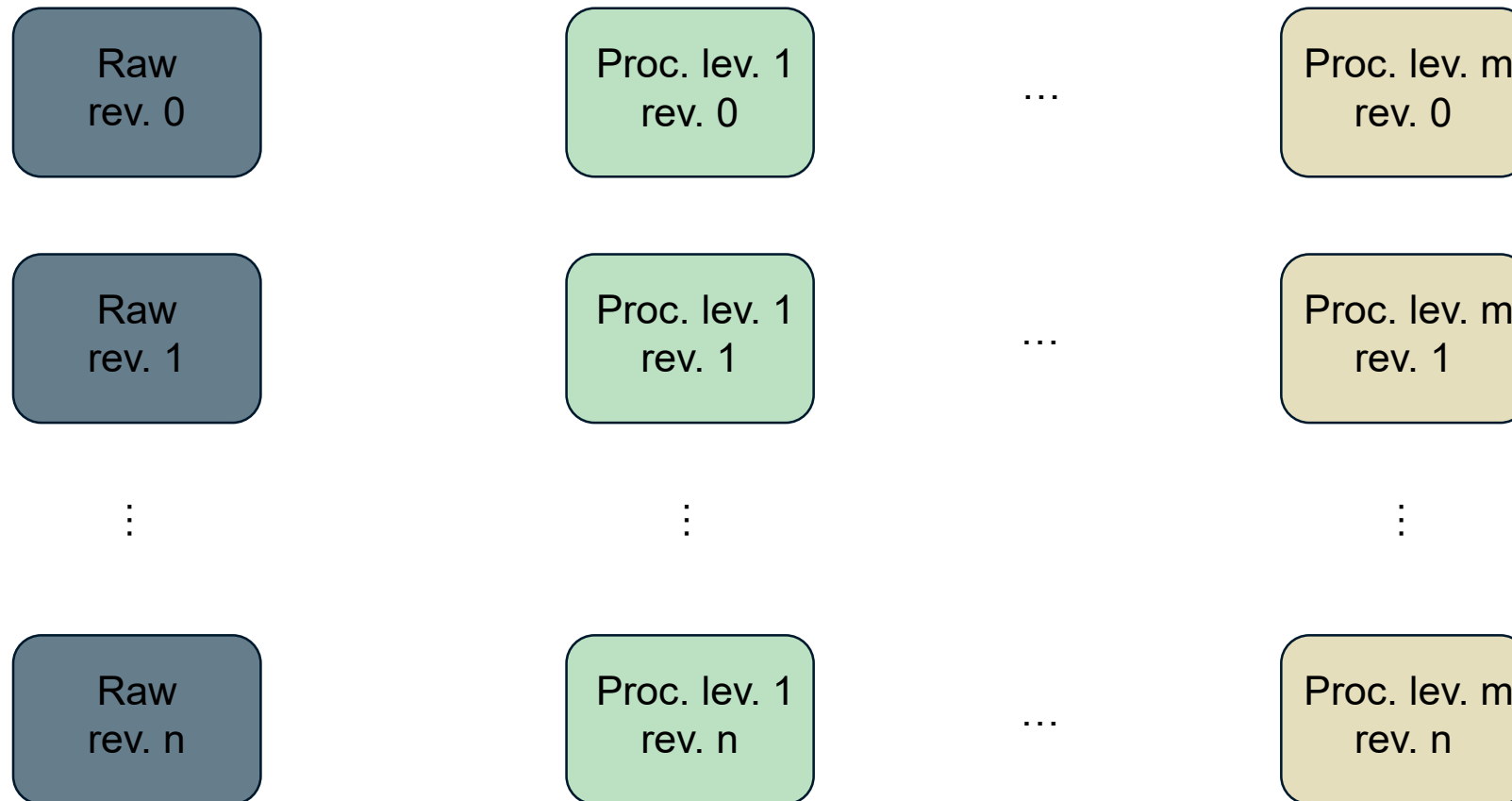
GitHub



Bitbucket



Collect and Store: Data versioning



Collect and Store: Data versioning tools



Renku (<https://renku.readthedocs.io/en/stable/index.html>)



Data Version Control (<https://dvc.org>)



Git Large File Storage (<https://git-lfs.com>)



Lake FS (<https://docs.lakefs.io>)

Collect and Store: File Naming

- Use unique names referencing content
- Limit to 42 characters (preferably less)
- Use ASCII characters, no spaces, points or special characters, e.g. ~!@#\$%^&*()[]{}<>';,»/
- Include dates and label versions
- Use names to order files:
 - Either, use Dates YYYY-MM-DD or YYYYMMDD (according to ISO 8601) at the beginning to enable chronological order
 - Or, use Versioning with leading zeroes to enable numerical order (enables versions to go beyond 9 without disrupting order)
- If you have started with your project use *Bulk Rename Utility* (Windows) or *Renamer 6* (Mac), *Rename/Thunar Bulk Rename* (GNU/Linux)

Collect and Store: File Formats (recommendation)

Data type	Recommended file formats
Text	<ul style="list-style-type: none"> • PDF/A • Plain Text coded as ASCII. UTF-8 or UTF-16 • XML
Spreadsheet	<ul style="list-style-type: none"> • CSV (NEAD)
Images	<ul style="list-style-type: none"> • TIFF (uncompressed or lossless compressed) • PNG
Code	<ul style="list-style-type: none"> • Languages with free environments (e.g. Py or R UTF-8 format of ASCII text)
Audio	<ul style="list-style-type: none"> • FLAC • Wav

Open and lossless formats

If you are using a proprietary format, think about adding an additional format

Collect & Store: Metadata Standards

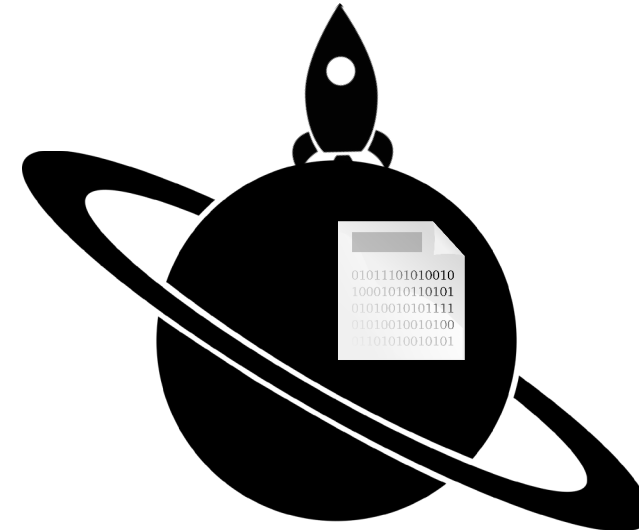
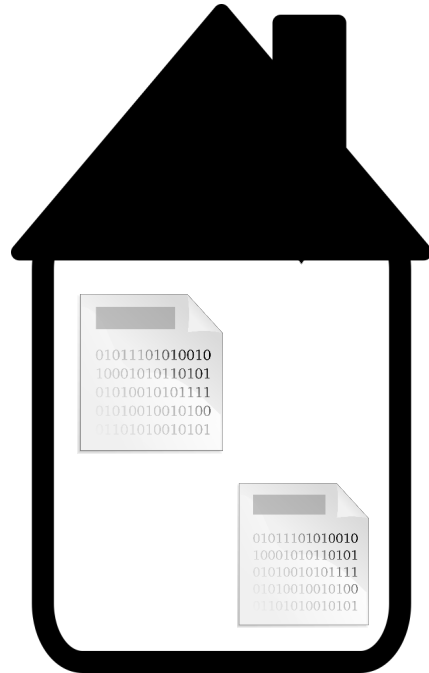
- Definition: Structured data that contains information about other data, but is not the content of the data.
- Metadata is very subject specific. The following directories are helpful:
 - Digital Curation Centre (<https://www.dcc.ac.uk/guidance/standards>)
 - RDA Metadata Standards (<https://rdamsc.bath.ac.uk>)
 - Fairsharing (<https://fairsharing.org>)
- Recommendation: Stick to a list of defined terms (controlled vocabulary) and don't use synonyms to describe the same object (e.g. picture or image)

Collect & Store: README File

General information	<ul style="list-style-type: none"> • Title of the dataset • Contact information principal investigator • Date of data collection • Geographic location
Data and file overview	<ul style="list-style-type: none"> • Short discription for each file name • Date
Sharing and access informations	<ul style="list-style-type: none"> • Licenses or restrictions
Methodological information	<ul style="list-style-type: none"> • Description of methods for data collection or generation • Description of methods used for data processing
Data specific information (repeat for each dataset)	<ul style="list-style-type: none"> • Variable list, including names and definitions • Units of measuments • Definition for codes or symbols to record missing data

Cornell University: Minimal viable content. For recommended visit: <https://data.research.cornell.edu/content/readme>

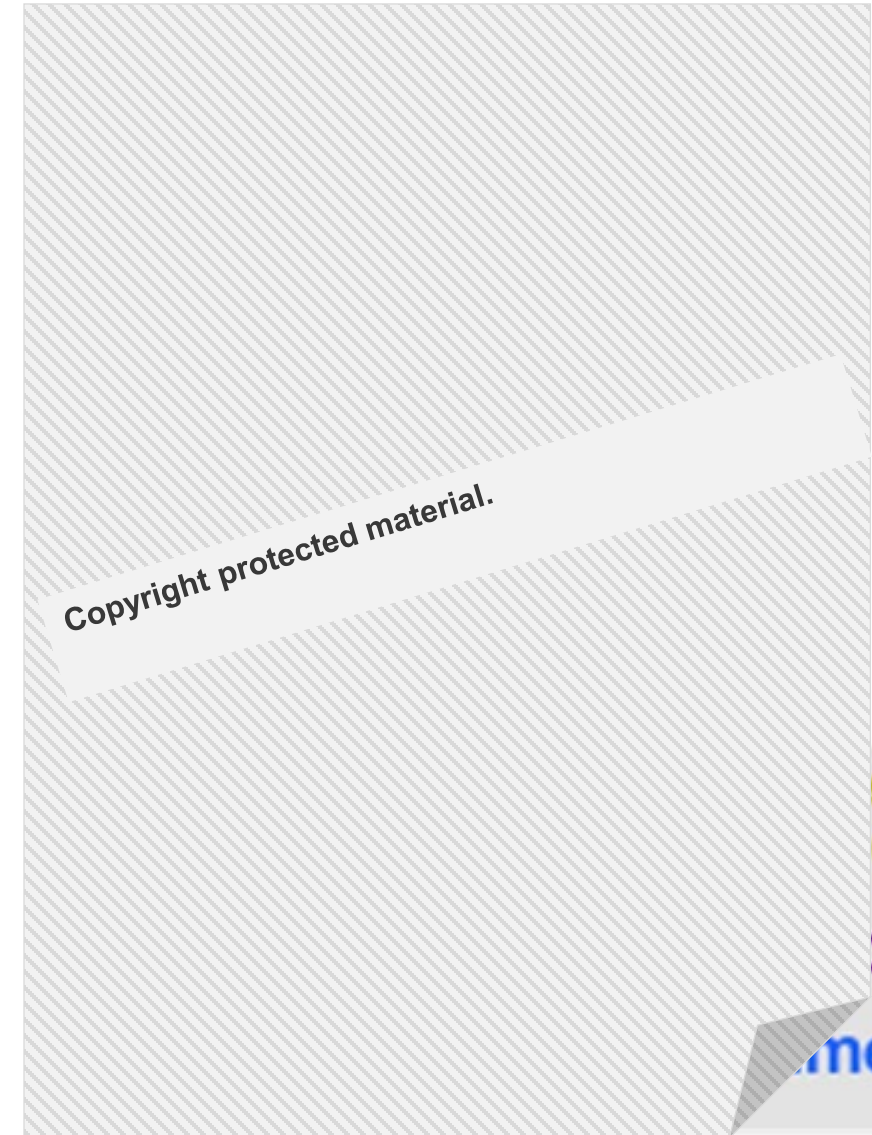
Collect and Store: 3 – 2 – 1 backup



Evaluate & Archive

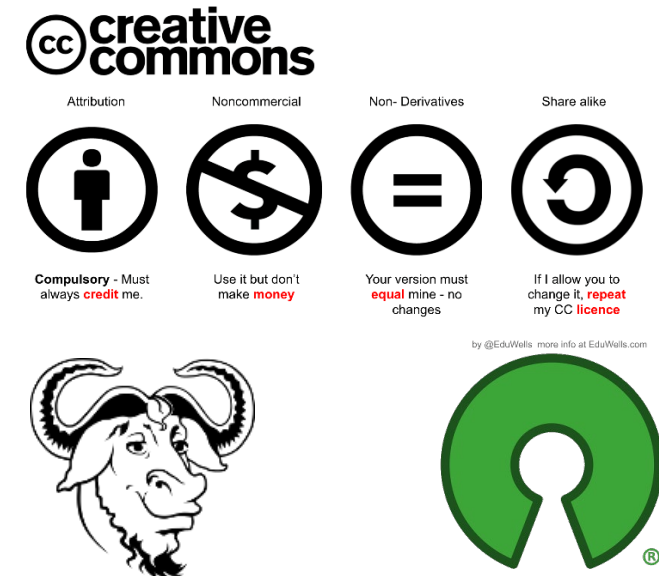
Evaluate & Archive: Data Protection

- **Relates to identified or identifiable person**
 - **Solutions (<https://dmlawtool.ccdigitallaw.ch/>) :**
 - Identity irrelevant -> anonymisation
 - Identity relevant -> Ask for consent
- > Pseudoanomization
- > Manage access rights
- > Ability to address subject's rights
- **Always contact Data Protection Officers at your Research Institute if your research involves personal data**



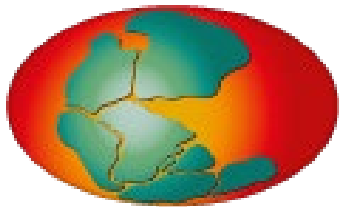
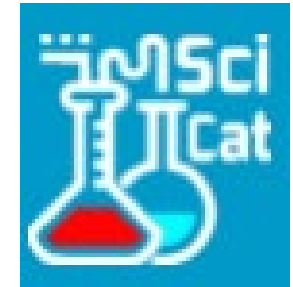
Evaluate & Archive: Data Protection

- Processed Data has copyright according to Swiss law
- Use CC licences when publishing factual data on data repositories (ideally CC 0)
- For software use licences specifically designed for software:
- Free Software (Open Source) licences like GPL, Apache, BSD and MIT.
- **Exceptions!** If you collaborated with external partners in your research project, you need to clarify together with them how and if data can be published.
- Contact the legal teams at your research institute if you feel lost.



Share & Disseminate

Share & Disseminate: The Choice of Data Repository



For alternatives: <https://www.re3data.org/>

RDM Services and Support at Eawag and Empa

OPEN SCIENCE

OPEN ACCESS

OPEN DATA

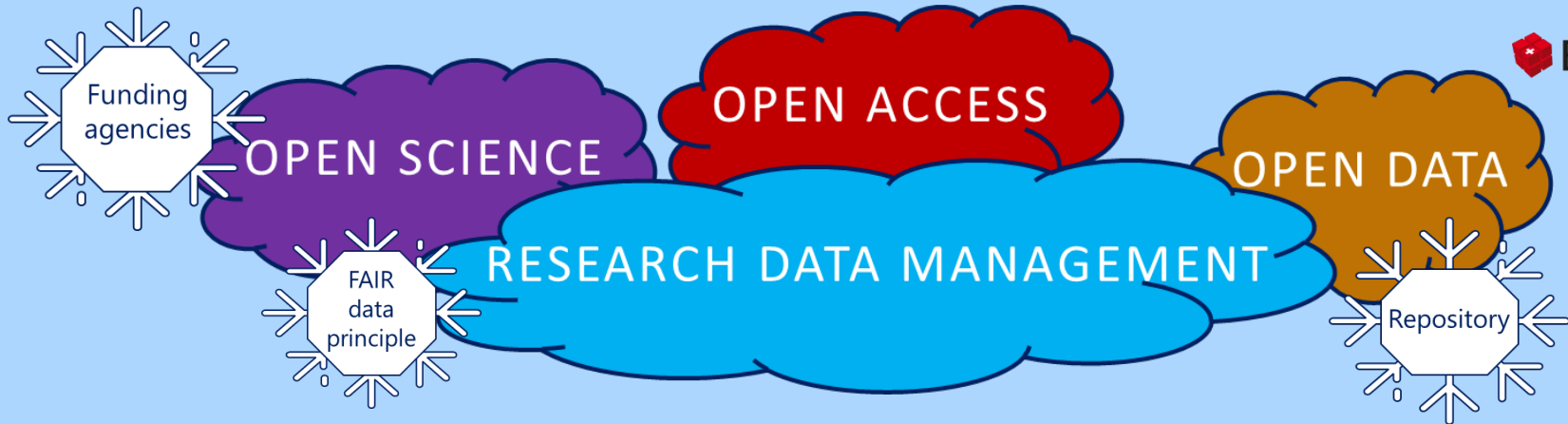
RESEARCH DATA MANAGEMENT

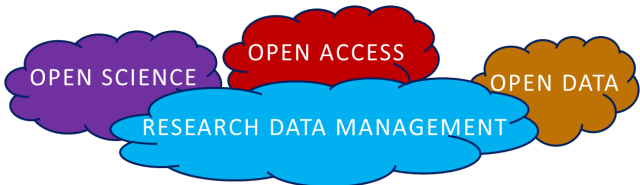
FAIR
data
principle

openBIS

Repository

Funding
agencies





DigitalScience@Empa

Intranet Plattform

On Intranet main page

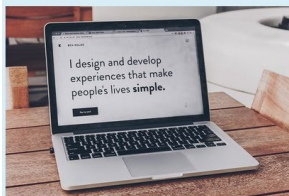
RESEARCH INFOS

Find [here](#) a list of Research Calls
 Find [here](#) a list of Research Awards
 DigitalScience@Empa - find [here](#):
 Tools and Platforms, Data Science, Modeling & Simulation, Open Science, Events & Trainings, Community

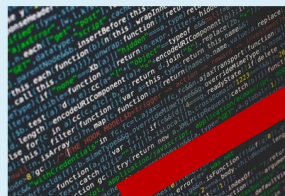
Where can I find Information & support

DigitalScience@Empa

Welcome to the platform about digitalization topics at Empa.



Tools & Platforms



Scientific IT



Modeling & Simulation



Open Science



Events & Training



Community

Our new documentation pages:

- Scientific IT team
- Computing & HPC
- Data Management & openBIS
- Data Science
- Software deployment & development

For input or feedback, contact scientificit@empa.ch.

<https://www.empa.ch/group/s909/overview>

RDM Services & Infrastructure @ Empa

DigitalScience@Empa

Welcome to the platform about digitalization topics at Empa.



Tools & Platforms Scientific IT Modeling & Simulation



Open Science Events & Training Community

Our new documentation pages:

[Scientific IT team](#)

[Computing & HPC](#)

[Data Management & openBIS](#)

[Data Science](#)

[Software deployment & development](#)

For input or feedback, contact scientificit@empa.ch.

How to manage your research data?

Check our information and support pages:

1. **RDM guidelines** to see the advantages of managing your data and possible horror stories if you do not: [here](#)
2. **openBIS** at Empa helps to manage your research data digitally, which includes save data storage, archiving and publishing in data repositories: [here](#)
3. **Data Management Plan (DMP)** to fulfill the requirements of funding agencies: [here](#)
4. **Open Access (OA)** to fulfill the requirement of publishing your paper openly: [here](#)
5. **Open Data Licenses** learn how licenses help to publish your research data openly: [here](#)
6. **Support & training** information [here](#) for:
 - openBIS & Data Management
 - DMP templates & how to get funding for RDM costs
 - Open Access by Lib4RI
 - Scientific IT weekly hours for Software development & Data Science questions
 - Python tutorial twice a year by Scientific IT

- ❑ Support for RDM
- ❑ SNF Data Management Plan template
- ❑ RDM guides, Best practice guide
- ❑ Handling of software licenses
- ❑ Open Access @ Empa
 - ❑ Policy
 - ❑ Publication fund

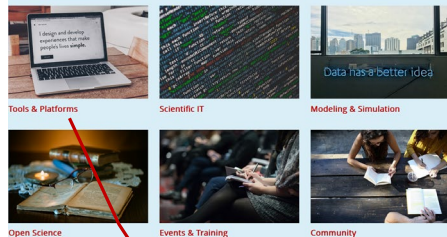
SNF Data Management Plan Template

- ❑ Sections analog to the SNF DMP
- ❑ Different categories for
 - ❑ Sensitive data
 - ❑ Using the data management system openBIS
 - ❑ Without using openBIS
- ❑ Description of backup system at Empa
- ❑ Recommendation of FAIR repository

Text
snippets
available

DigitalScience@Empa

Welcome to the platform about digitalization topics at Empa.



Our new documentation pages:

Scientific IT team

Computing & HPC

Data Management & openBIS

Data Science

Software deployment & development

For input or feedback, contact
scientificit@empa.ch.

Data Management System @ Empa

openBIS, the **Electronic Laboratory Notebook (ELN)** and **Inventory Management System (LIMS)** at Empa, which enables easy connection

- ✓ **Zenodo** – data repository according FAIR principles
- ✓ **Jupyter Hubs** – for programming in Python, R etc.
- ✓ **Longterm Archiving** – of research data at ETH

Collaboration platform – Code versioning

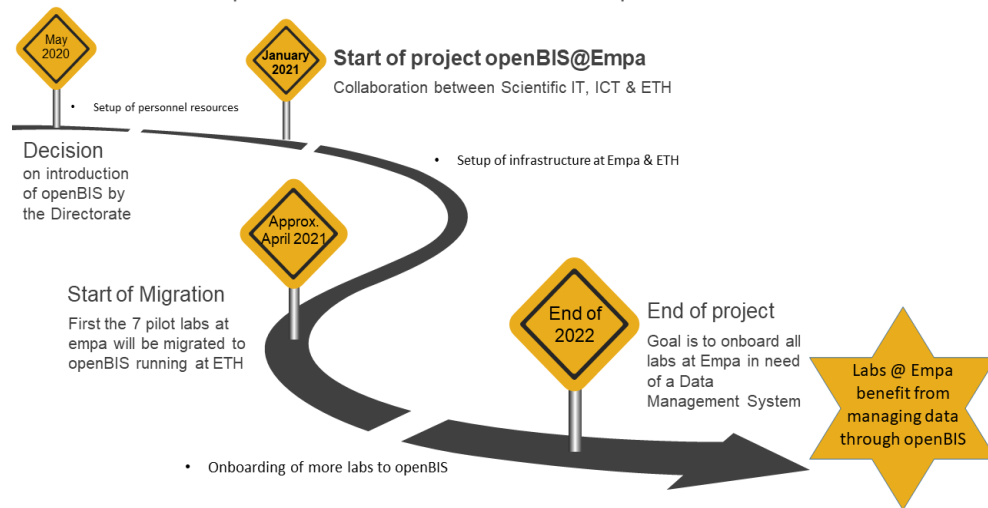


GitLab @Empa

High performance computing (HPC)

Provision of researchers with technical support regarding HPC-related common problems and long-term projects.

Timeline of openBIS introduction@Empa

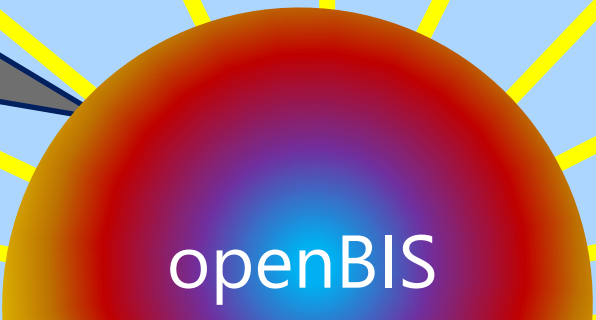


RESEARCH DATA MANAGEMENT



OPEN SCIENCE

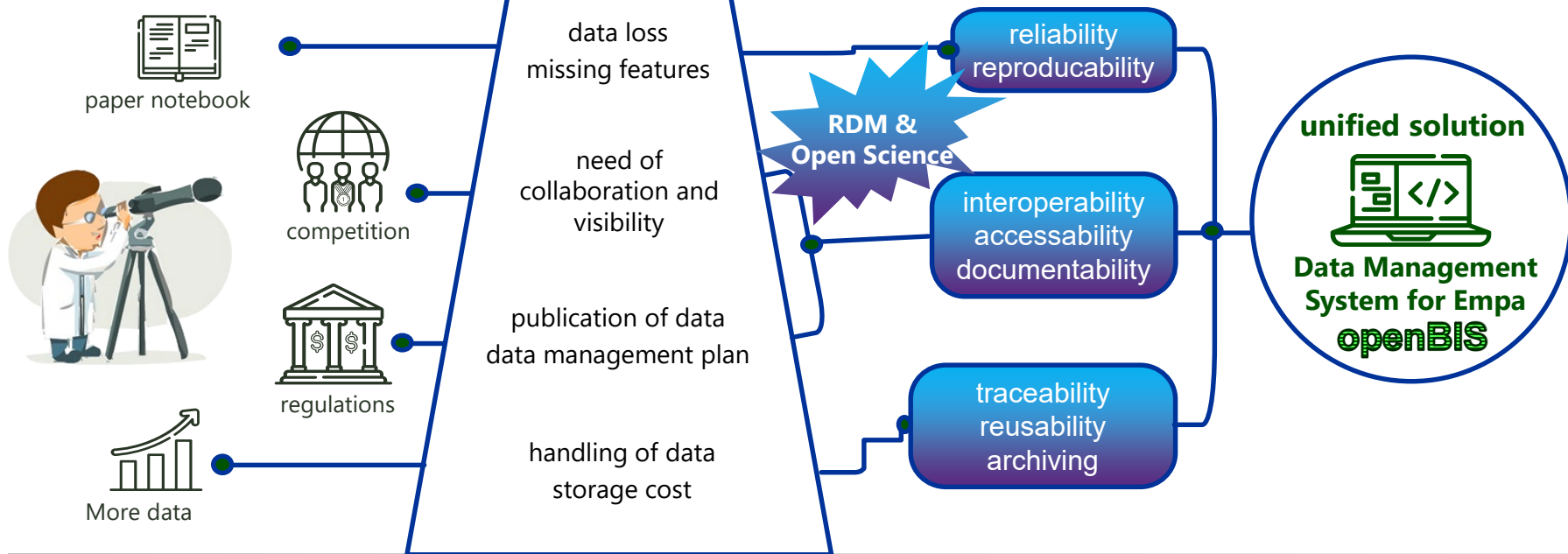
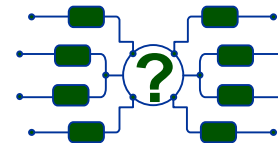
OPEN DATA

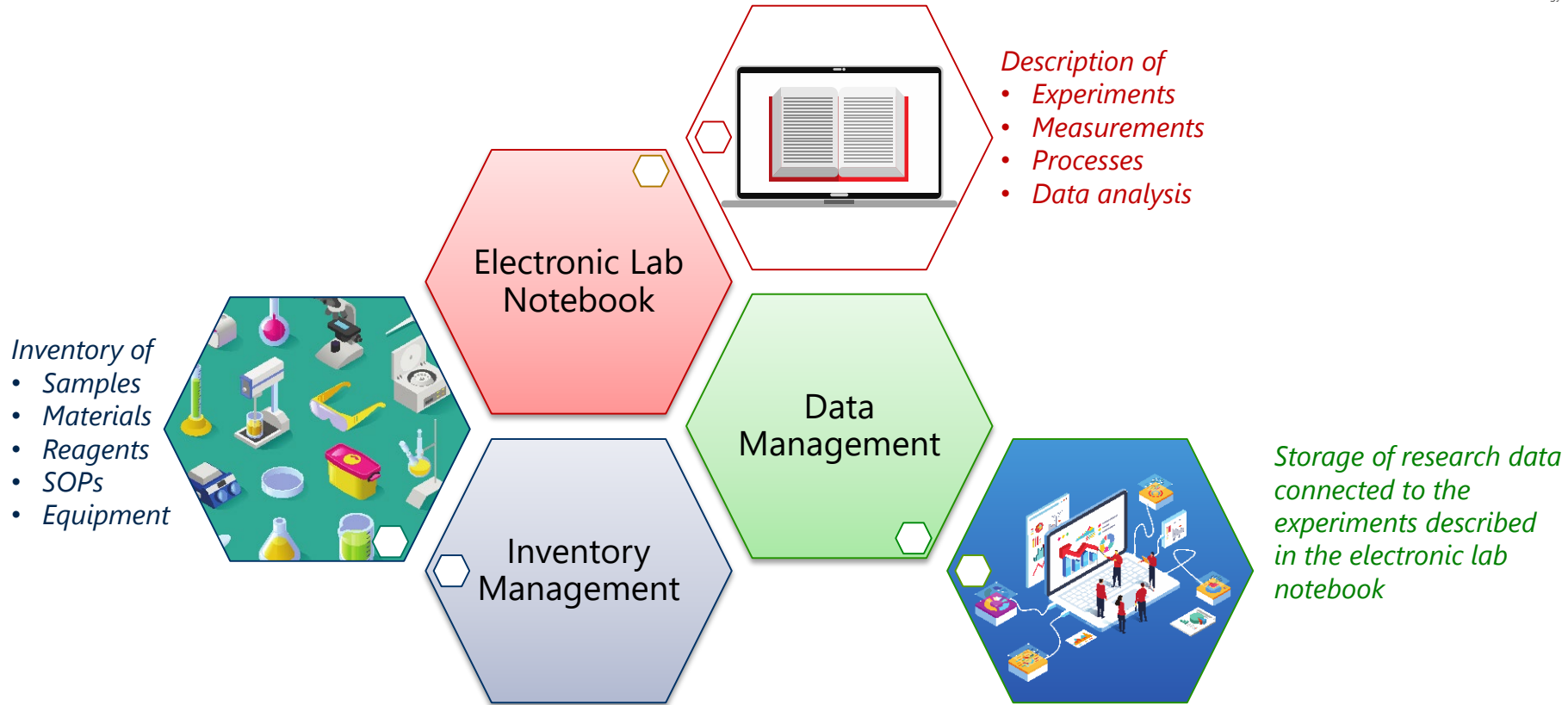


openBIS helps to solve the challenges

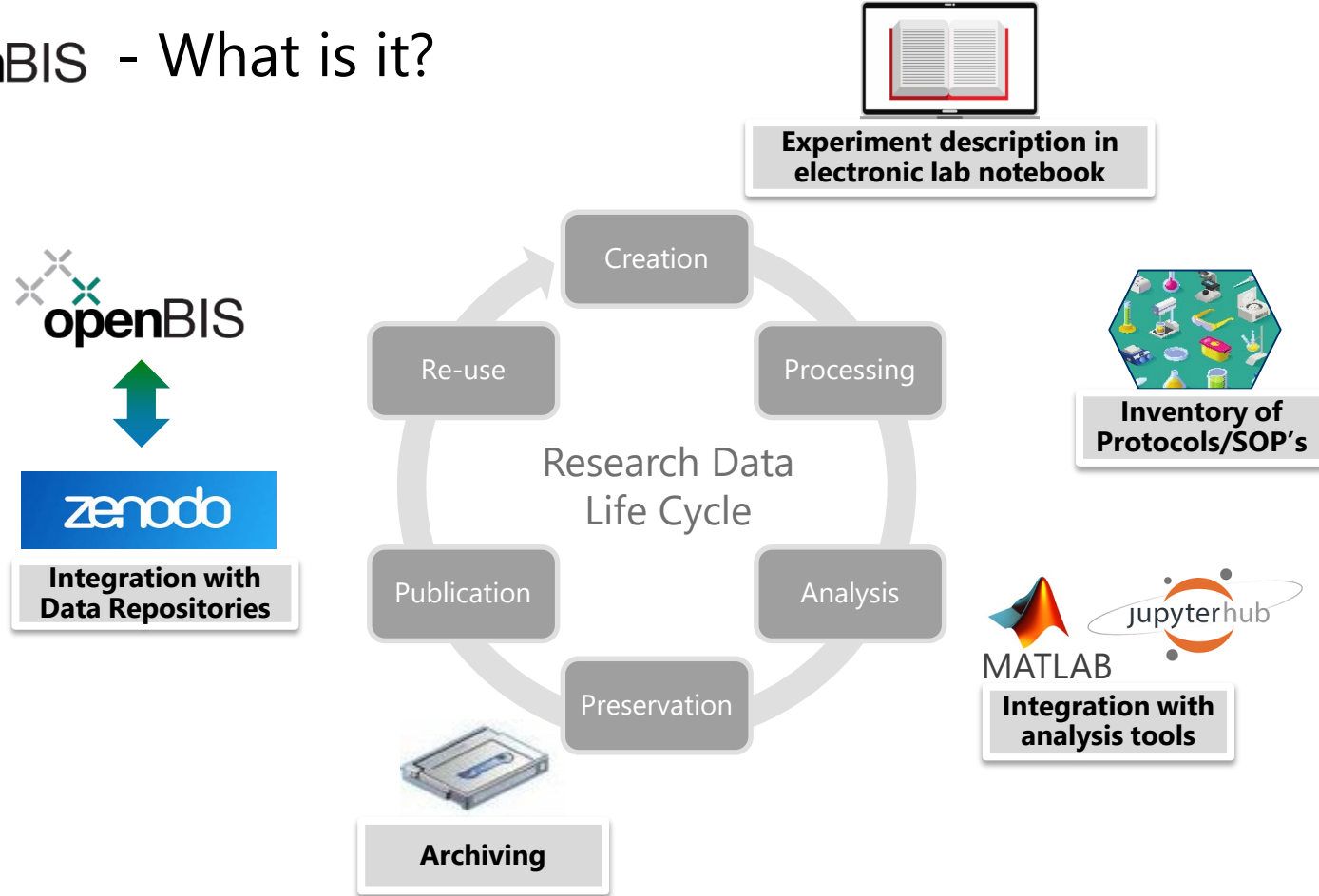


Landscape – Challenges – Solution

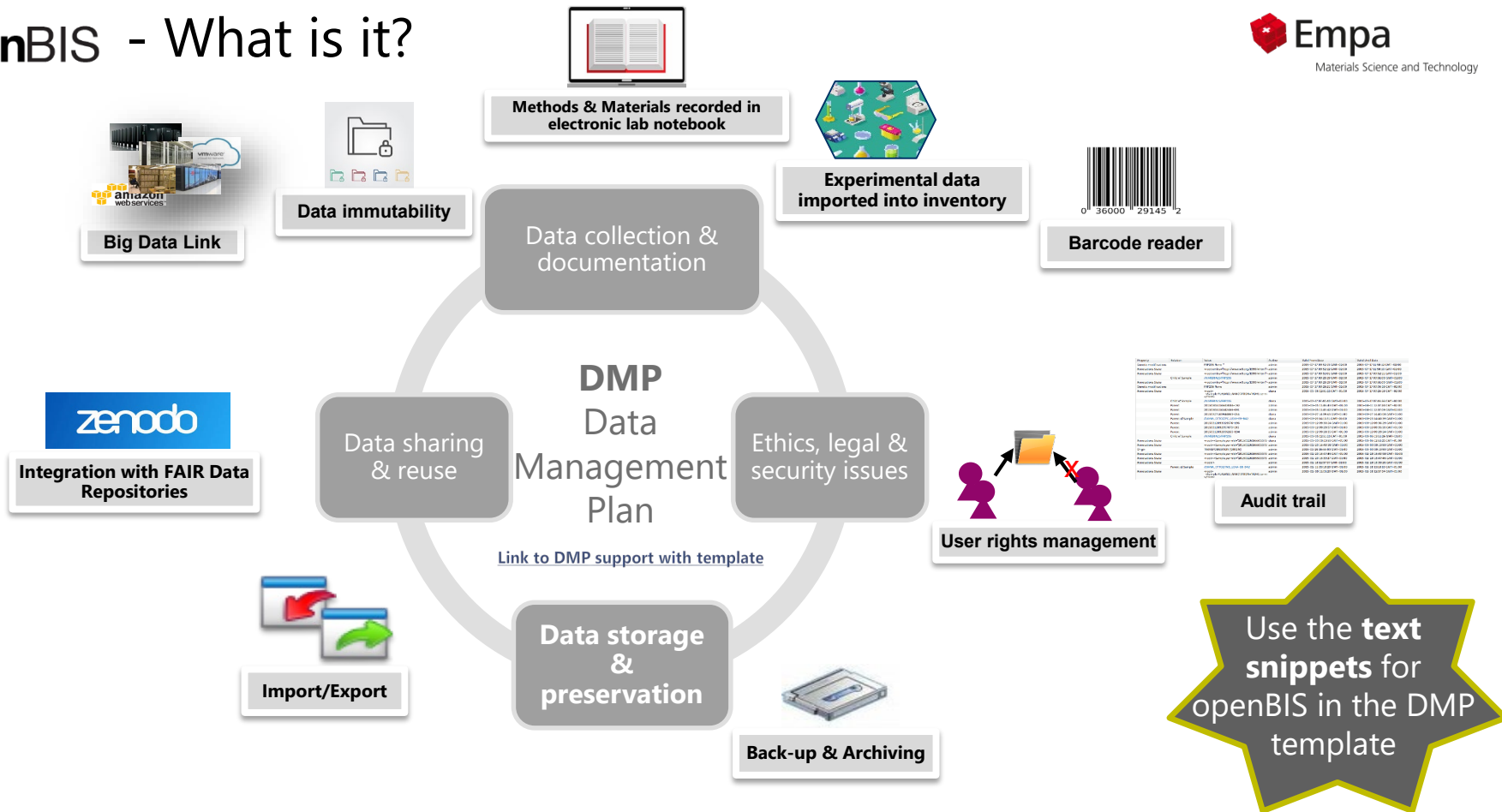




A data management system – Connecting lab inventory, research data & lab notebook in 1 tool

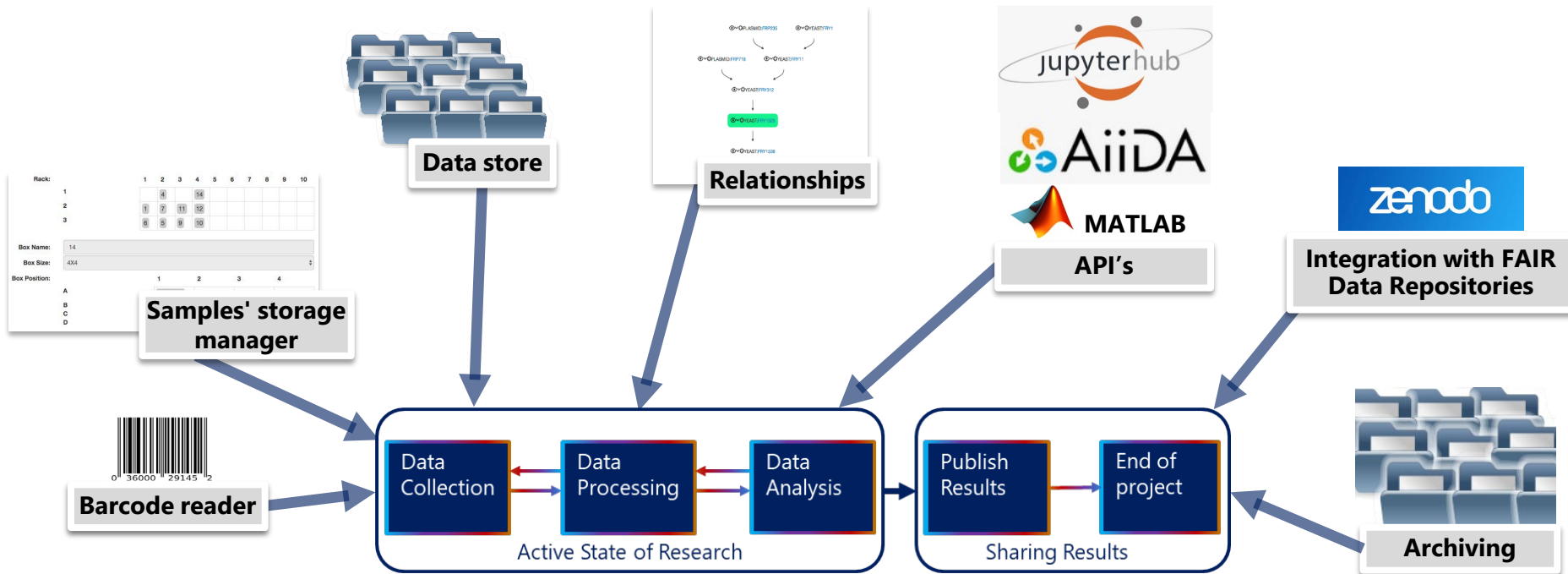


A data management system – Covering most of the data life cycle in 1 tool



A data management system – Solving DMPlan requirements of funding agencies with 1 tool

For managing research data from "bench" to publication in a central storage

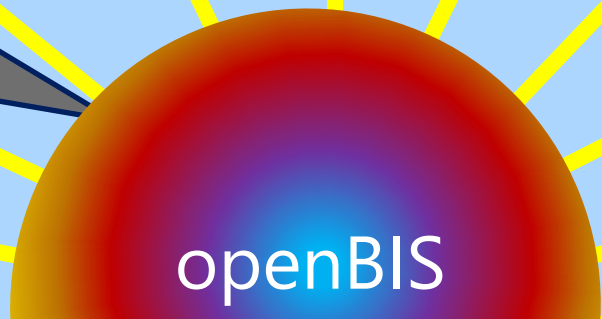


RESEARCH DATA MANAGEMENT



OPEN SCIENCE

OPEN DATA

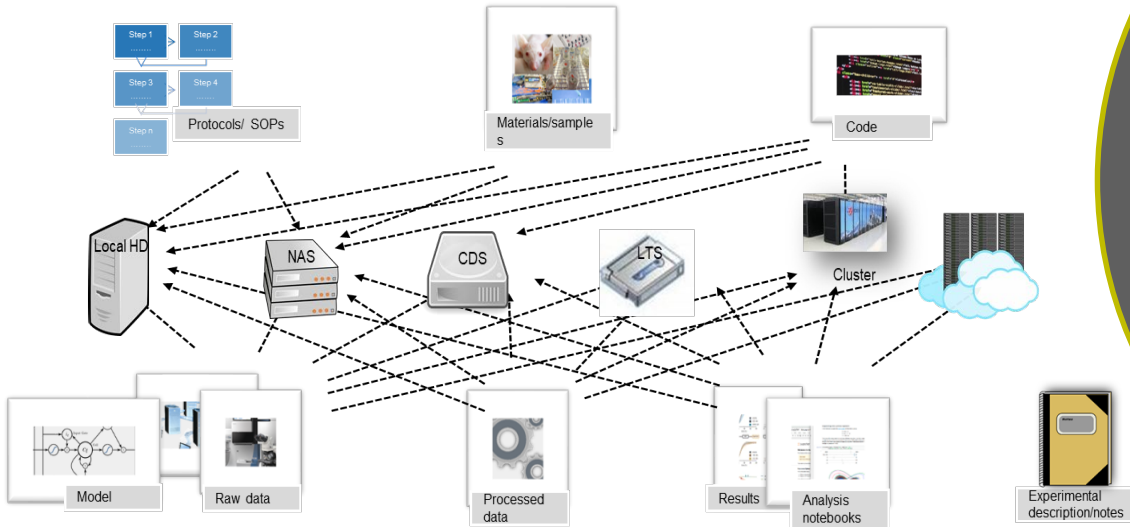


openBIS helps to solve the challenges

What are the main benefits?



A common scenario The "data spread"



The ideal solution Having everything connected in 1 place



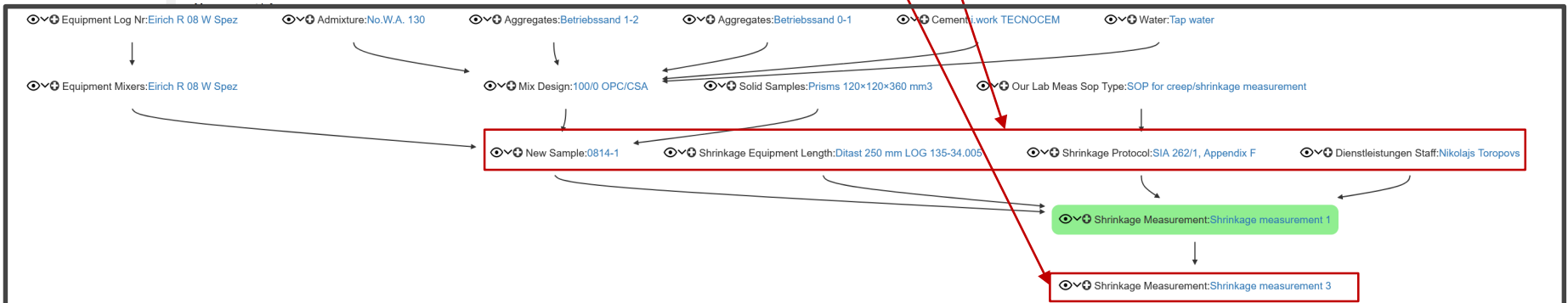
Shrinkage Measurement: Shrinkage measurement 1

Parents

Code	Name	Last calibration	Measurement accuracy (in mm)	Gauge reference length (in mm)	Shrinkage dimensionality	SOP ID	Short name of staff	Type
SAMPLE_1	0814-1							New Sample
SHRINKAGE_EQUIPMENT_LENGTH1	Ditast 250 mm LOG 135-54.005	2018-08-07 00:00:00 +0200	0.001	250				Shrinkage Equipment Length
SHRINKAGE_PROTOCOL_1	SIA 262/1, Appendix F				Linear shrinkage	4003		Shrinkage Protocol
STAFF_8	Nikolajs Toropovs						ton	Dienstleistungen Staff

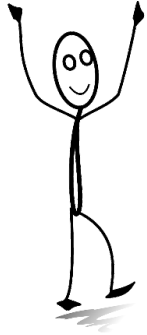
Children

Code	Name	Identifier	Shrinkage type	Sample's age at the measurement time point (days)	Shrinkage side A (in microns)	Shrinkage side B (in microns)	Notes
SHRINKAGE_MEASUREMENT_3	Shrinkage measurement 3	/GROUP_MICHELE.GRIFFA_AT_EMPA.CH /PROJECT_1/SHRINKAGE_MEASUREMENT_3		3	164.0	146.0	





openBIS benefits



- 1** Prevention of loss of research data & knowhow via structured documentation & storage
- 2** Easy & automatic data archiving over a long period
- 3** Easy & fast connection to repository Zenodo



Challenges at Empa


- 1. Data loss** due to change of personnel & no proper documentation of data
- 2. No archiving solution** easy accessible & reliable
- 3. No repository** available



- ❖ **Without openBIS** you need to solve these challenges on your own
- ❖ **Without proper data management** you risk losing funding money

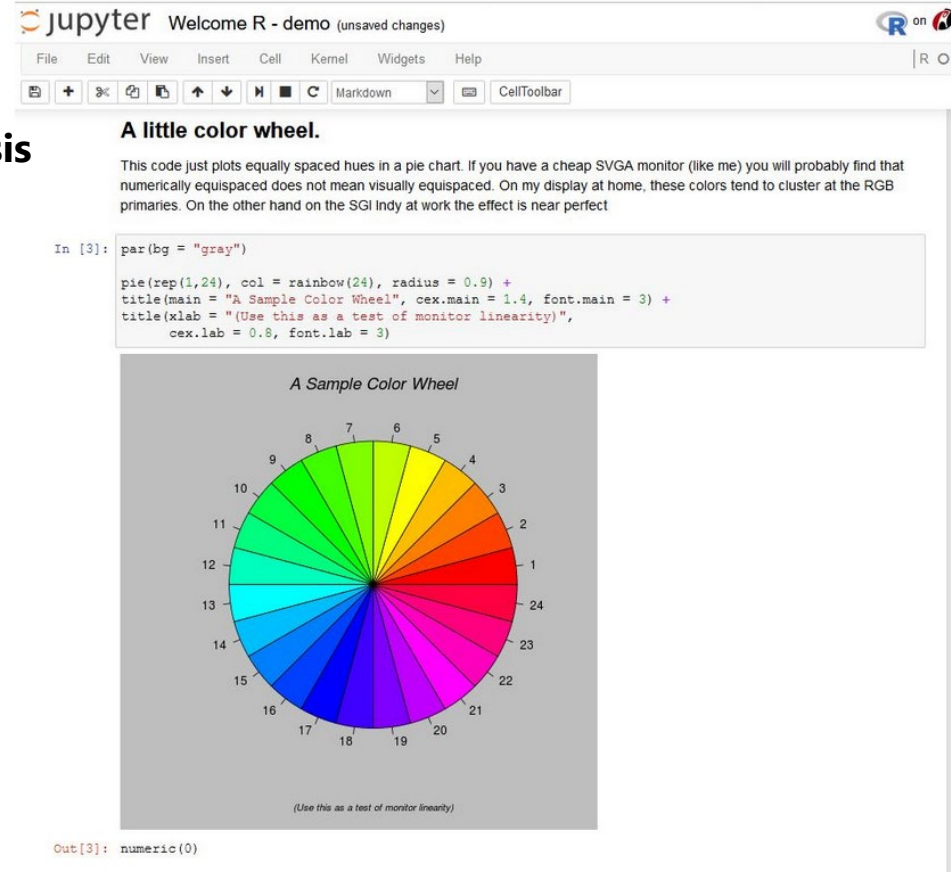
Data Science with openBIS

- ❖ **Jupyter notebooks combine code, documentation & outputs** like plots, images, videos etc.
- ❖ Useful for **interactive/exploratory data analysis** and **reproducibility**
- ❖ **Easy sharing** of code with documentation and results
- ❖ Like a **modern lab notebook** for reproducible coding



openBIS comes with Jupyter Hubs for data analysis via **Jupyter Notebooks** that support R, Python, Octave

openBIS API available for Matlab and Python (Pybis)



The screenshot shows a Jupyter Notebook interface with the title "Welcome R - demo (unsaved changes)". The menu bar includes File, Edit, View, Insert, Cell, Kernel, Widgets, and Help. The toolbar contains icons for file operations, a dropdown menu set to "Markdown", and a "CellToolbar" button. The notebook content includes a section titled "A little color wheel." with a paragraph of text and a code cell. The code cell contains the following R code:

```
In [3]: par(bg = "gray")

pie(rep(1,24), col = rainbow(24), radius = 0.9) +
title(main = "A Sample Color Wheel", cex.main = 1.4, font.main = 3) +
title(xlab = "(Use this as a test of monitor linearity)",
      cex.lab = 0.8, font.lab = 3)
```

The output of the code is a pie chart titled "A Sample Color Wheel" with 24 segments of different colors, numbered 1 through 24. The chart is displayed on a gray background. Below the chart, the text "(Use this as a test of monitor linearity)" is visible. The output of the code cell is "Out[3]: numeric(0)".

OPEN ACCESS

OA POLICY OF EMPA

RESEARCH DATA MANAGEMENT



OA support



Golden way



Green way



Funding agencies

openBIS

OPEN SCIENCE

OPEN DATA

Repository

FAIR data principle



REPRODUCIBLE



RELIABLE



STORING BACK UP SECURING



DMP support

Scientific IT
RDM@empa.ch



DATA MANAGEMENT PLAN (DMP)

openBIS



COLLECTING ORGANIZING DOCUMENTING

F.A.I.R.



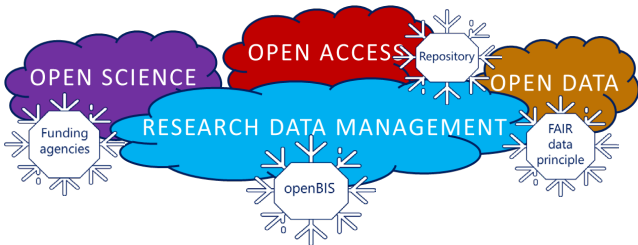
ARCHIVING SHARING

Repository Zenodo
Reproducible Data Analysis



REUSABLE

Solutions @ Empa for
RDM & Open Science



DigitalScience@Empa

Intranet Plattform <https://www.empa.ch/group/s909/overview>

Support

scientificit@empa.ch

<https://scientificit.empa.ch/>



Eleni Pratsini
Lead Scientific IT



Anusch Bachthaler
OpenBIS
Project Manager



Edoardo Baldi
MPC & CSCS
Software Developer



Simone Baffelli
Research Software
Engineer



Stefania Hauser
OpenBIS
(201, 312)



Matthias Flosslein
OpenBIS Lead S&E
Research Data
Management



Despina Adamopoulou
Data & Software Engineer
Full-Stack Developer



Aleksandr Yakusovich
Research Software
Engineer

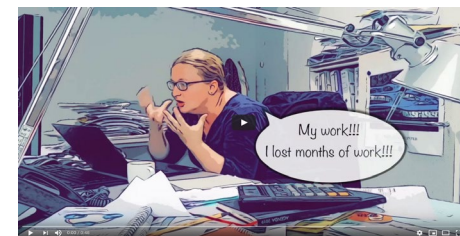


Videos

From Open Data to Open Knowledge
SIB - Swiss Institute of Bioinformatics • 507 views • 11 months ago

Discover SIB's vision on Open Data, one of the many facets of Open Science – the movement to make scientific research and its dissemination accessible to the society. Harmonizing licenses of datab...
CC

1:26



- https://youtu.be/t_rEXpfCTrg
- <https://youtu.be/tFWd2M2OXwQ>
- <https://youtu.be/6kHGbbdFuDE>
- <https://youtu.be/LCzjzP916o>
- <https://youtu.be/NdkIWkRi-ZQ>

RESEARCH DATA MANAGEMENT @ EAWAG

— SERVICES & SUPPORT —

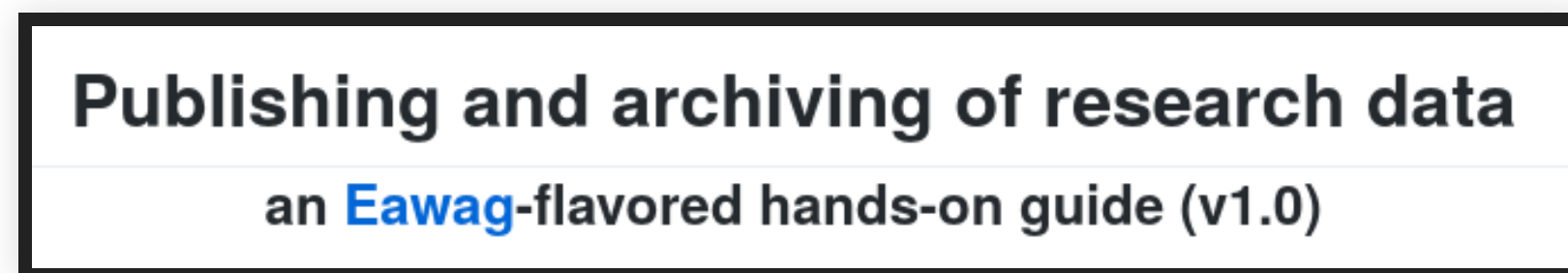
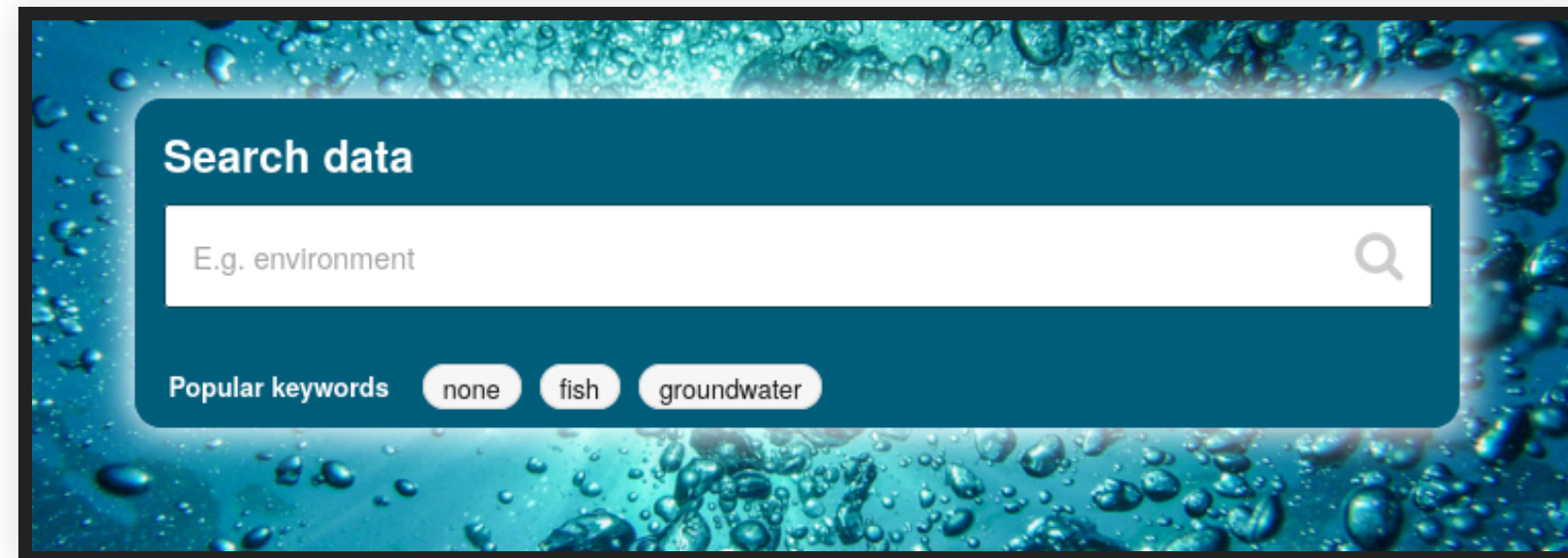


WHAT I DO?



- Data management
- Software development
- DevOps
- Improve ERIC
- Teaching best practices

ERIC internal



Eawag Research Data Institutional Collection

- Development, Maintenance
- No review, quality control or curation
- Support
 - data organization, formats and annotation
 - process automation

ERIC **open**

- **D**evelopment, **M**aintenance
- **DOI** reservation & registration
- Support reg. **workflow** in sync with article review & publishing
- (Meta)data **dissemination**, interlinking (ORCID, article-DOIs, ...)
- **Review, quality control and curation**

Data Management Plans

- **Guide:** Eawag SNF / SNSF Data Management Plan Guide
- DMP reviews

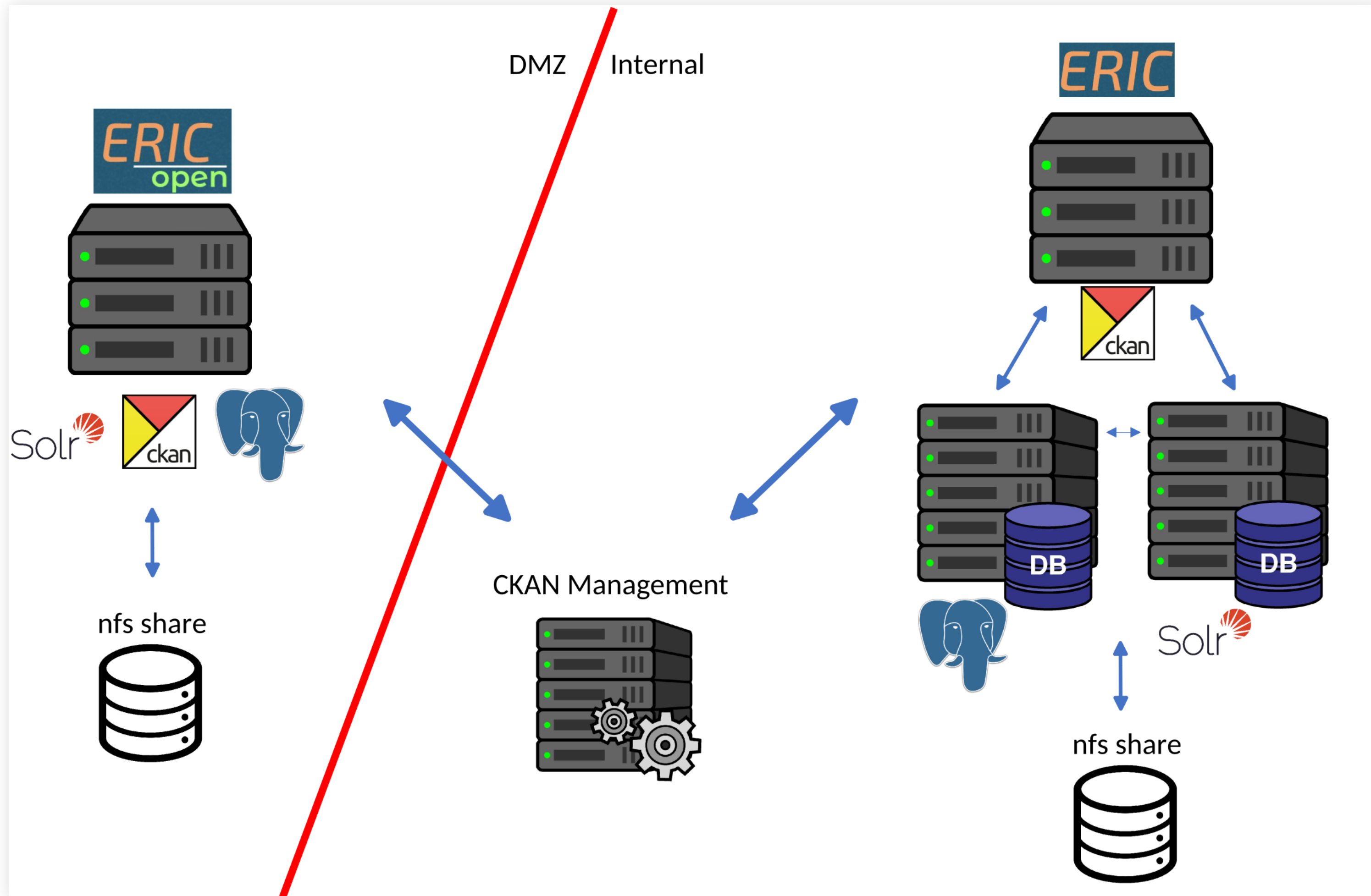
Data management planning

- Sampling, surveys
- Data transmission/transport
- Backup strategy, data safety
- Data security, encryption
- Legal issues, licensing, personal data, anonymization

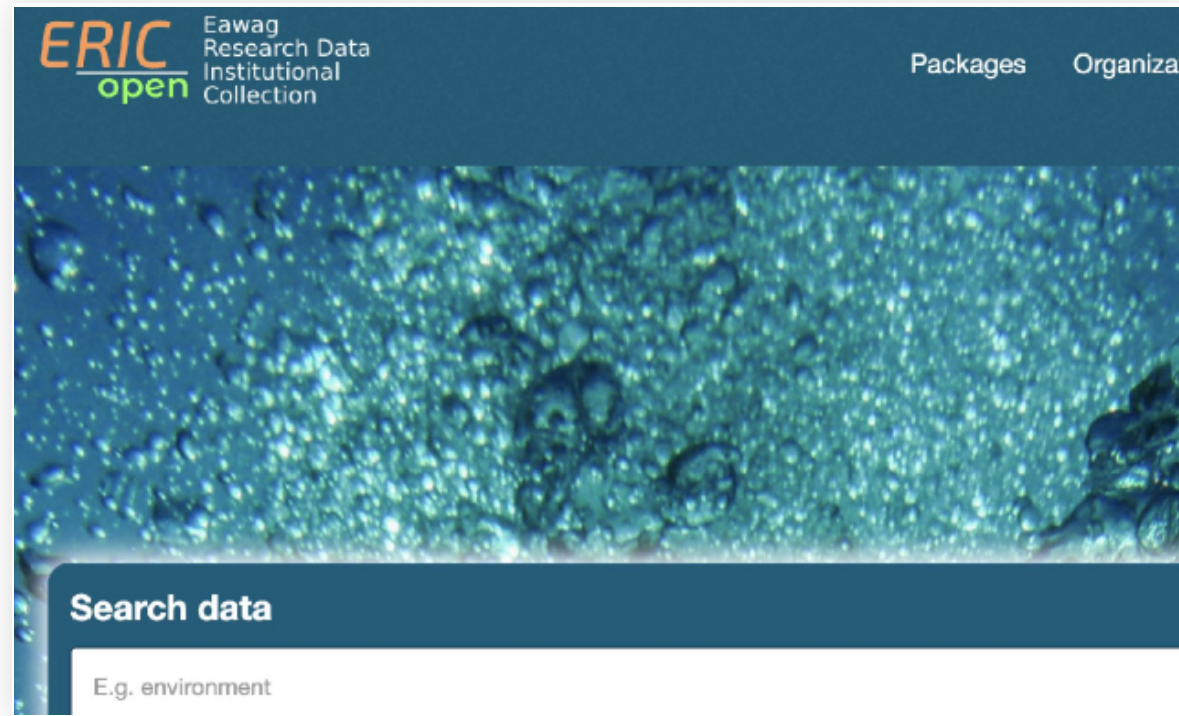
- Comprehensive Knowledge Archive Network (**CKAN**)
- Findable Accessible Interoperable Reusable (**FAIR**)
- **SNF approved**
- **Secure** infrastructure
- **RDM knowledge base**
- **FAQ** ([here](#))
- **Eawag** branding
- **Meta data** schema from **DataCite**
- All extensions on **GitHub** in the **eawag-rdm** organization



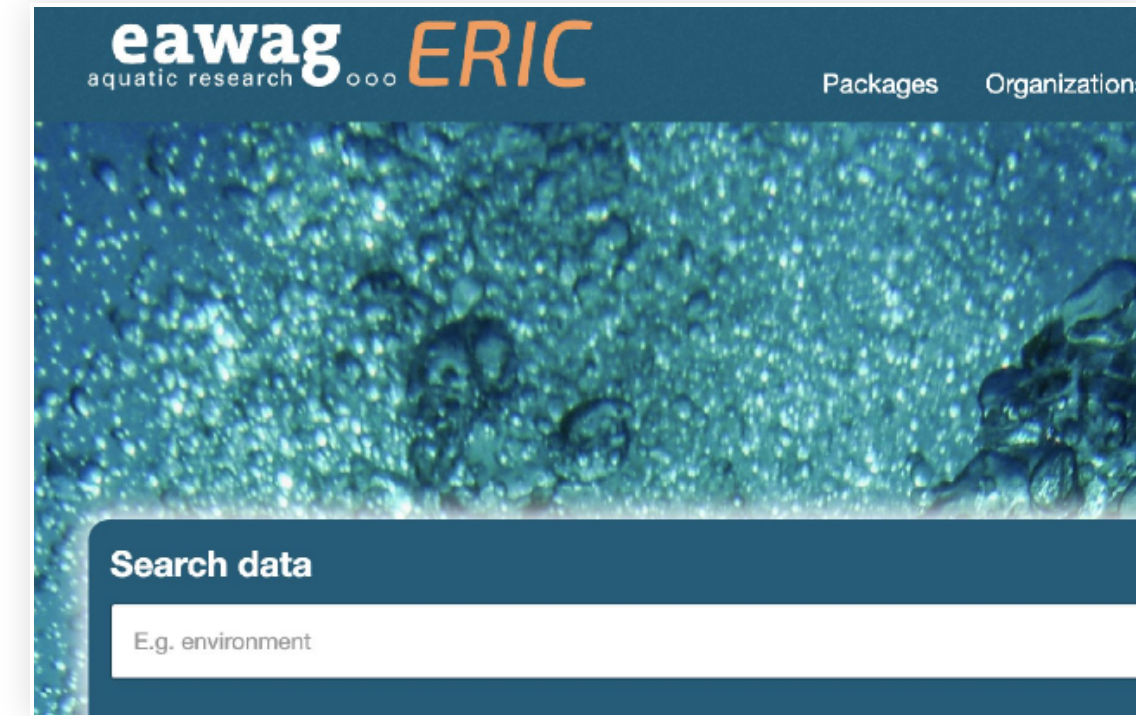
INFRASTRUCTURE



OPEN VS INTERNAL



- opendata.eawag.ch
- Openly accessible
- Strict checking
- Pure **read** data portal
- **Immutable** datasets



- data.eawag.ch
- Eawag internal
- Unchecked data
- **Read/write** platform
- **Mutable** datasets

PUBLICATION WORKFLOW



Good practices throughout

- Keep your data **save** and **clean** (**minimal guide**)
- **Version** your data (**renku**)
- Use **version control** (**git**, ...)

First time **ERIC** use

- Log in with your **Eawag** credentials
- Contact the **data manager** of your department to get **editor** privileges

Upload to **ERIC internal**

- Create a **package/group**
- Upload your **neat** and **organized** data
- Write **rdm@eawag.ch** (naming your package)

RDM initial feedback

- Reserved DOI for your publication
- Data publication (checklist)

Iterative package improvement

- Review of checklist
- Rigorous checks (data protection, checksums, data usage, basic code review, check links, ...)
- Feedback

We are aiming for an iteration of 1! ;-)

You **should**

- Add your current **code**
- Add **links** to repositories, related packages
- Provide accurate and plentiful **meta data**

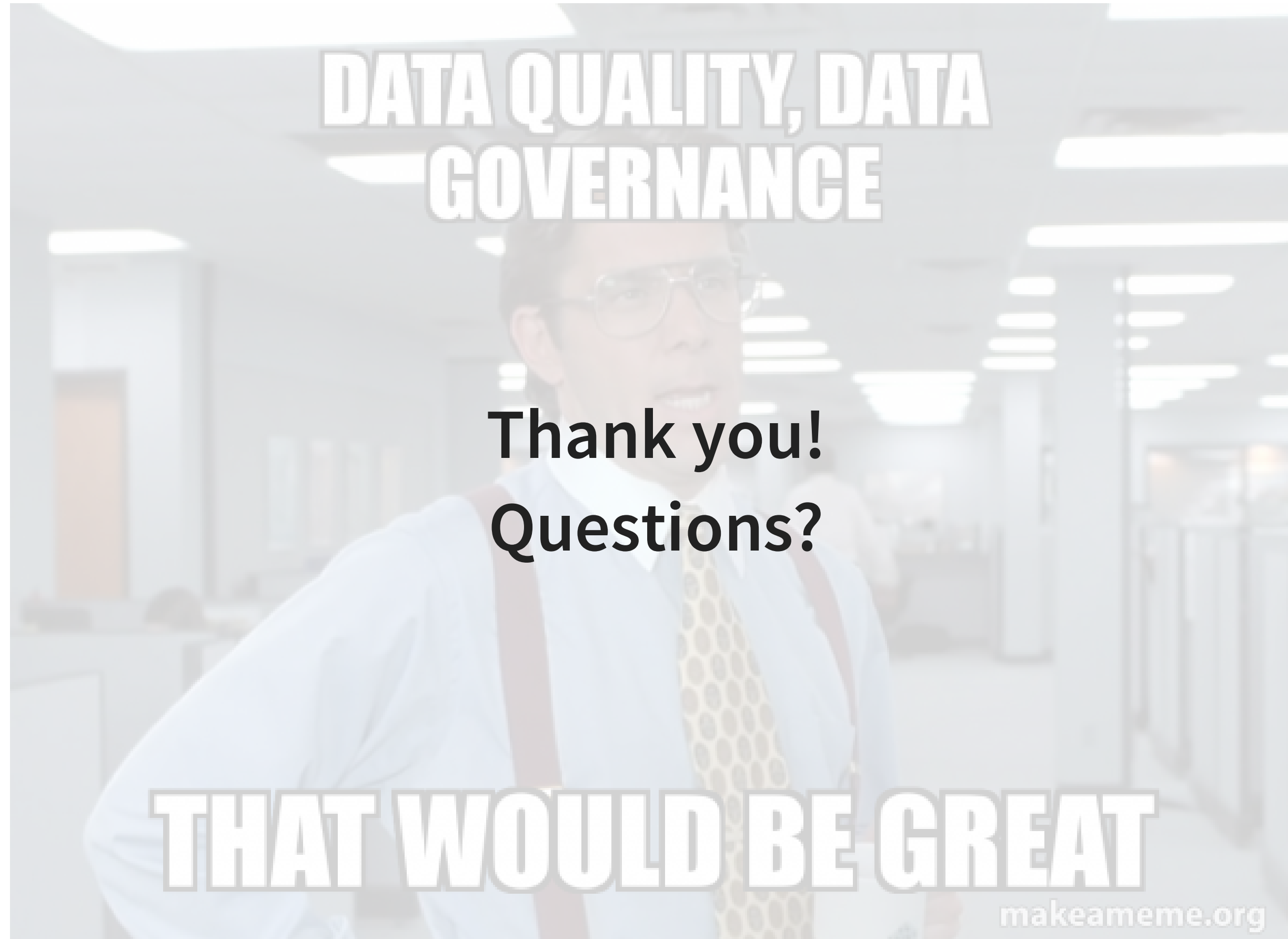
You **should not**

- Add your **paper** to the data package
- Incorporate **Eawag specific information** like paths on a network share
- Violate data protection **guidelines**

Now the data can be made public!



- Update **ERIC**
- Extend **functionality**
- Improve and standardize **documentation**
- Better support for other parts of the **DLC** like "life" database clusters



Plan & Design: Data Management Plan (DMP)

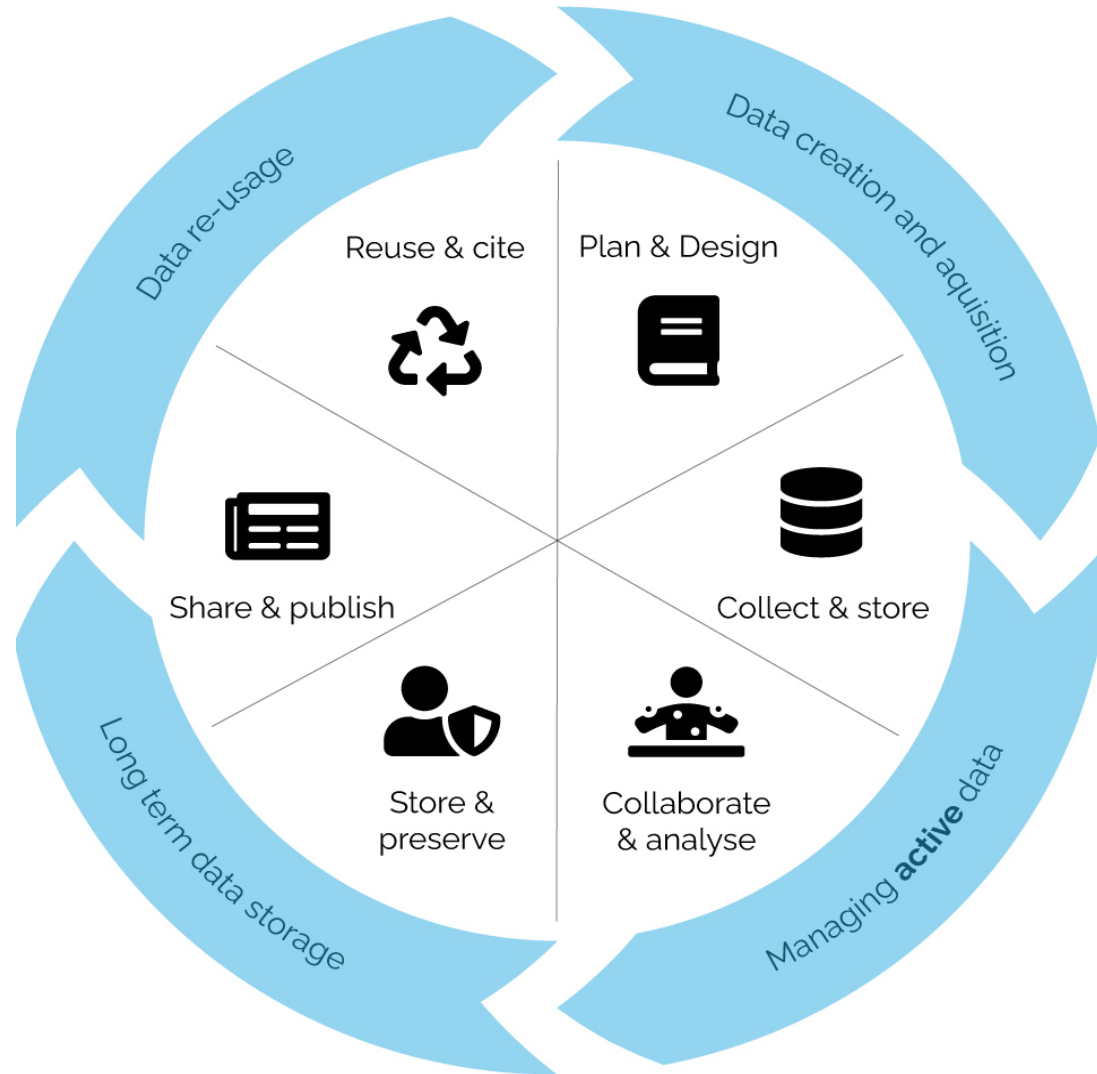
Plan & Design: Why?



Bibliothèque de l'EPFL, «RDM Horror stories | Episode 2 – Stranger Data Things», 11th February 2020.

<https://bit.ly/3qPWMIS>

DMP



Covers the whole Research Data Life Cycle

Plan & Design: DMP

- What types of data will be collected and which code (incl. software) will be created or used?
- How will you document the data used and code programmed?
- Where will data and code be stored?
- Who owns the data and code is responsible for security and backup?
- Which data and code will be shared and preserved?
- How will data be shared and with whom?

Plan & Design: DMP

Applications and Projects

Grant application 1

1. Personal data

- Responsible applicant
- Other applicants
- Applicants' employment
- Project partners




2. Application data

- Basic data I
- Basic data II
- Use-inspired project
- Re-submission
- Continuation of
- Link to other SNSF projects
- Further requested and available funds (not from the SNSF)
- University or research institution
- Requested funding
- Data management plan (DMP)**
- Research requiring authorisation or notification
- Exclusion of external reviewers
- General remarks on the project




3. Annexed documents (upload)

- Research plan
- CV and major achievements
- Quotes
- Cover letter
- Official certificates
- Weave/Lead Agency and International Co-Investigator Scheme
- Other annexes



1. Data collection and documentation

-  1.1 What data will you collect, observe, generate or reuse?
-  1.2 How will the data be collected, observed or generated?
-  1.3 What documentation and metadata will you provide with the data?





2. Ethics, legal and security issues

-  2.1 How will ethical issues be addressed and handled?
-  2.2 How will data access and security be managed?
-  2.3 How will you handle copyright and Intellectual Property Rights issues?

3. Data storage and preservation

-  3.1 How will your data be stored and backed-up during the research?
-  3.2 What is your data preservation plan?

4. Data sharing and reuse

-  4.1 How and where will the data be shared?
-  4.2 Are there any necessary limitations to protect sensitive data?
-  4.3 All digital repositories I will choose are conform to the FAIR Data Principles.
-  4.4 I will choose digital repositories maintained by a non-profit organisation.

Plan & Design: DMP



- Keep it short and simple
- Be stingy with words
- Have one idea per sentence
- Use the active form
- Use positive phrases
- Use concrete terms

«we used the method» not «the method was used»
 «the results are different» not «the results are not the same»
 «it will be published in Nature» not «it will be published in a reputable journal»



- Don't write in «sophisticated style»
- Save on adjectives and adverbs
- Avoid unnecessary constructions
- Don't nominalise
- Don't use empty modifiers
- Don't use tautologous modifiers

e.g. «It is clear that», «the fact is that», «in an attempt to», «in order to»
 «reduce» not «achieve a reduction in length»
 e.g. «basically», «indeed», «quite», «actually»
 e.g. «completely finish», «may potentially», «ultimate result», «blue in colour»

Plan & Design: DMP

- 1. Organize yourselves in groups of two (5 minutes)**
- 2. Each group will engage with the first section of the SNSF DMP (20 minutes)**
 - Read requirements
 - Write answers and questions
 - Discuss with other group members
 - Designate presenter
- 3. Presentation and discussion of findings (20 minutes)**

Plan & Design: DMP - Data Collection and Documentation

1.1 What data will you collect, observe, generate or reuse?

- Type, format (NEAD), content, volume of data, reference to data (if reused)

1.2 How will the data be collected, observed, generated?

- Standards methodology, quality assurance
- File organisation and versioning (folder structures, git, ELN/LIMS, etc.)

1.3 What documentation and metadata will you provide?

- Scientific Metadata (README, metadata standards)
- General Metadata (Depending on choice of data repository)

Plan & Design: DMP - Ethics, Legal and security issues

2.1 How will ethical issues be addressed and handled?

- Information and consent to using personal data, location of critical infrastructure as well as rare and protected species
- Requirements for assessments by ethical review boards, permission by third parties
- Description of Pseudonymisation or Anonymisation Methods

2.2 How will the data access and security be managed?

- Distinguish datasets according to the level of risk (cf. §2.1) and use an adverb to describe the level of risk («high», «medium», «low»)
- State Storage Location, secure transmission, access restriction, IT infrastructure

2.3 How will you handle copyright and Intellectual Property Rights Issues?

- Consider non-disclosure agreements, potential patents, research collaborations across institutions
- Recommendation to use CC0 where possible

Plan & Design: DMP - Data Storage and Preservation

3.1 How will your data be stored and backed-up during the research?

- Backup strategy for work at all stages of research (amount of storage needed, frequency of updates, responsibilities, security measures)

3.2 What is your data preservation plan?

- Data formats
- Selection mode for data to be preserved (all relevant data related to reported results, long term preservation of unique datasets)

Plan & Design: DMP - Data Sharing and Reuse

4.1 How and where will the data be shared?

- Repository of choice (non-commercial preferred and required for contribution of up to 10'000 CHF for storage)
- Metadata Policy of said repository

4.2 Are there necessary limitations to protect sensitive data?

- Reasons data cannot be published at certain times (Section §2.1)

4.3 All Digital Repositories I will choose conform to FAIR Data?

- Check box

4.4 All Digital Repository I will choos are mainained by a non-profit oranisation?

- If no, provide justification (costs will not be covered)

Thank you for your attention!

Feedback!

Please give us a short feedback

Questions?

Presentation slides: lib4ri.ch > Learn
> Trainings

Appendix

Appendix: Eawag

- **Four links under data.eawag.ch:**
 - <https://opendata.eawag.ch/eawagrmd/help/quickstart.html>
 - <https://opendata.eawag.ch/eawagrmd/help/opendata.html>
 - <https://doi.org/10.25678/000066>
 - https://www.internal.eawag.ch/fileadmin/intranet/informatik/datenman/rdm/directive_archiving_of_researchdata.pdf
- **Difference between ERIC/internal (data.eawag.ch) and ERIC/open (opendata.eawag.ch)**
- **Services are in the form of guides and consulting. Most notable guides in addition to the one mentioned above are**
 - <https://doi.org/10.25678/000033>
 - <https://opendata.eawag.ch/eawagrmd/software-licensing.html>
- **Finally the list of resources can be helpful:**
 - <https://opendata.eawag.ch/eawagrmd/resources.html>

Appendix: Empa

- o General overview of topics:

<https://www.empa.ch/web/s909/overview>

- o Support topics like DMP template of Empa:

<https://www.empa.ch/web/s909/support1>

OpenBIS

- o General overview: <https://www.empa.ch/group/s909/openbis>

- o Documentation & trainings info:

<https://www.empa.ch/group/s909/documentation-tutorials>

Appendix: File Formats EPFL

Bibliothèque de l'EPFL, Research Data, fast guide #4», 2019,
<https://bit.ly/3NFloYx>

TYPE OF DATA	APPROPRIATE	ACCEPTABLE	DEPRECATED
Tabular (extensive metadata)	CSV – HDF5	TXT – HTML – TEX – FASTQ ^[3] – POR	
Tabular (minimal metadata)	CSV – TAB – ODS – SQL – TSV	XML (if appropriate DTD) – XLSX	XLS – XLSB
Textual / Presentation	TXT – PDF – ODT – ODM – TEX – MD – HTM – XML – EXTXYZ ^[4] – ODF	PPTX – RTF – DOCX – PDF (with embedded forms) – EPS – IPF	DOC – PPT – DVI – PS
Code / Computation	M – R – PY – IYPNB – RSTUDIO – RMD – NETCDF – AIML	SDD	MAT – RDATA
Image & Spectroscopy	TIF – PNG – SVG – JPEG – FITS	JCAMP – JPG – JP2 – TIF – TIFF – PDF – GIF – BMP – DM3 – OIR – LSM ^[5]	INDD – AIT – PSD – SPC
Audio	FLAC – WAV – OGG – MXL – MIDI – MEI – HUMDRUM	MP3 – AIF	
Video	MP4 – MJ2 – AVI – MKV	OGM – MP4 – WEBM	WMV – MOV – QT
Geospatial	NETCDF – tabular GIS attribute data – SHP – SHX – DBF – PRJ – SBX – SBN – POSTGIS – TIF – TFW – GEOJSON	MDB – MIF	
3D structures & images	X3D – X3DV – X3DB – PDF3D – POV – PDBML	DWG – DXF – PDB	PXP
Generic	XML – JSON – RDF		

Appendix: File Formats ETH Zürich

Assessment of various file formats

Table 1: Our assessment of future readability of some common file formats. (For more detailed information we refer to the recommendations of the Bundesarchiv (German), the KOST (German or French), the Memoriav, the Forschungsdatenzentrums Archäologie & Altertumswissenschaften IANUS (Germany), the Library of Congress and the Harvard Library)

File type	Recommended	Suitable to only a limited extent	Not suitable for archiving
Text	<ul style="list-style-type: none"> PDF/A (*.pdf, preferred subtypes 2b and 2u) Plain Text (*.txt, *.asc, *.c, *.h, *.cpp, *.m, *.py, *.r etc.) coded as ASCII, UTF-8, or UTF-16 using byte order mark XML (inclusive XSD/XSL/XHTML etc.; with included or accessible schema and character encode explicitly specified) 	<ul style="list-style-type: none"> PDF (*.pdf) with embedded fonts Plain text (*.txt, *.asc, *.c, *.h, *.cpp, *.m, *.py, *.r etc.) (ISO 8859-1 coded) Rich Text Format (*.rtf) HTML and XML (The ASCII text is readable over long term; try to avoid external links.) <p>Not accepted for publication, OK for supplementary materials:</p> <ul style="list-style-type: none"> Word *.docx PowerPoint *.pptx LaTeX, TeX (The ASCII text is readable over long term; open source software required for formatting and the resulting PDF should be included.) OpenDocument formats (*.odm, *.odt, *.odg, *.odc, *.odf) 	<ul style="list-style-type: none"> Word *.doc PowerPoint *.ppt
Spreadsheet or table	<ul style="list-style-type: none"> Comma- or tab delimited text files (*.csv) 	<ul style="list-style-type: none"> Excel *.xlsx (container format) OpenDocument spreadsheets (*.ods) 	<ul style="list-style-type: none"> Excel *.xls, *.xlsb (binary formats)
Raw data and workspace		<ul style="list-style-type: none"> ASCII Text is suitable for long-term use, but the data import may be time-consuming. S-Plus files (*.sdd) may be saved as text files. Matlab *.mat files may be saved in HDF Format. Saving nontrivial ASCII Matlab *.mat files should be avoided because they are not readable with the Matlab load command (see table 2). Network Common Data Format or NetCDF (*.nc, *.cdf) Hierarchical Data Format (HDF5) (*.h5, *.hdf5, *.he5) 	<ul style="list-style-type: none"> Binary files such as the standard Matlab files *.mat or the R files *.RData
Raster image (bitmap)	<ul style="list-style-type: none"> TIFF (*.tif) (uncompressed, preferentially TIFF 6.0, Part 1: baseline TIFF). TIFF is preferred as compared to PNG or JPEG2000. Portable Network Graphics (*.png, uncompressed) JPEG2000 (*.jp2, lossless compression) Digital-Negative-Format (*.dng) to keep raw data of digital fotos in addition to an second copy in TIFF format 	<ul style="list-style-type: none"> TIFF (*.tif) (compressed) GIF (*.gif) BMP (*.bmp) JPEG/JFIF (*.jpg) JPEG2000 (lossy compression) (*.jp2) 	
Vector graphics	<ul style="list-style-type: none"> SVG without JavaScript binding (*.svg) 		<ul style="list-style-type: none"> Graphics InDesign (*.indd), Illustrator (*.ait) Encapsulated Postscript (*.eps) Photoshop (*.psd)
CAD	<ul style="list-style-type: none"> AutoCAD Drawing (*.dwg) Drawing Interchange Format, AutoCAD (*.dxf) Extensible 3D, X3D (*.x3d, *.x3dv, *.x3db) 		
Audio	<ul style="list-style-type: none"> WAV (*.wav) (uncompressed, pulse-code modulated) 	<ul style="list-style-type: none"> Advanced Audio Coding (*.mp4) MP3 (*.mp3) 	
Video ¹	<ul style="list-style-type: none"> FFV1 codec (version 3 or later) in Matroska container (*.mkv) 	<ul style="list-style-type: none"> MPEG-2 (*.mpg, *.mpeg) MP4, which is also called MPEG-4 Part 14 (*.mp4) QuickTime Movie (*.mov) ² Audio Video Interleave (*.avi) Motion JPEG 2000 (*.mj2, *.mjp2) 	<ul style="list-style-type: none"> Windows Media Video (*.wmv)

Footnotes

¹ In addition to the file format (or container format), also the codec and the compression method are important. See Ianus, Memoriav and KOST for further information.

² In the Version of Nov 21, 2016 of the current document, the format QuickTime Movie was downgraded from „Recommended“ to „Suitable to only a limited extent“. Apple discontinued the support of Windows QuickTime Player in the year 2016. Windows Media Player thus only supports file format versions 2.0, or earlier, of QuickTime Movie files.

Appendix: References (Slide 18)

¹ SPARC Europe, «The Open Data Citation Advantage», 2017, <https://sparceurope.org/open-data-citation-advantage/>.

² Digital Science, «The state of Open Data Report», 2019, https://digitalscience.figshare.com/articles/report/The_State_of_Open_Data_Report_2019/9980783/2

³ European Commission and PwC, «Cost-Benefit analysis fro FAIR research Data», 2019.

<https://op.europa.eu/en/publication-detail/-/publication/d375368c-1a0a-11e9-8d04-01aa75ed71a1>

⁴ Baker, M., “1,500 scientists lift the lid on reproducibility”. *Nature* 533, 452–454 (2016).

<https://doi.org/10.1038/533452a>

Appendix: Icon References

Slide 4:

- Le Moign, Vincent, «Lab Scientist Icon», <https://icon-icons.com/icon/lab-scientist/101049>, free for commercial use.
- Flaticon, «Checkliste», https://www.flaticon.com/de/kostenloses-icon/checkliste_2666469, free for personal and commercial use.
- PLoS, «Open Access logo», https://de.wikipedia.org/wiki/Datei:Open_Access_logo_PLoS_white.svg, CC-0.
- «Databases and People», <https://freesvg.org/databases-and-people>, CC-0.

Slide 8

- Felixmh, «Krischen-Früchte-Natur-Symbol», free commercial use.