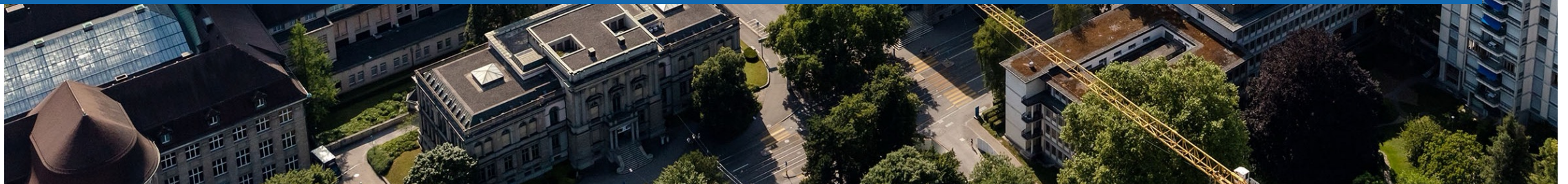# Supporting the reproducible research value chain in bioinformatics: from raw data to final publication

Caterina Barillari, Rostyslav Kuzyakiv, Michal Okoniewski
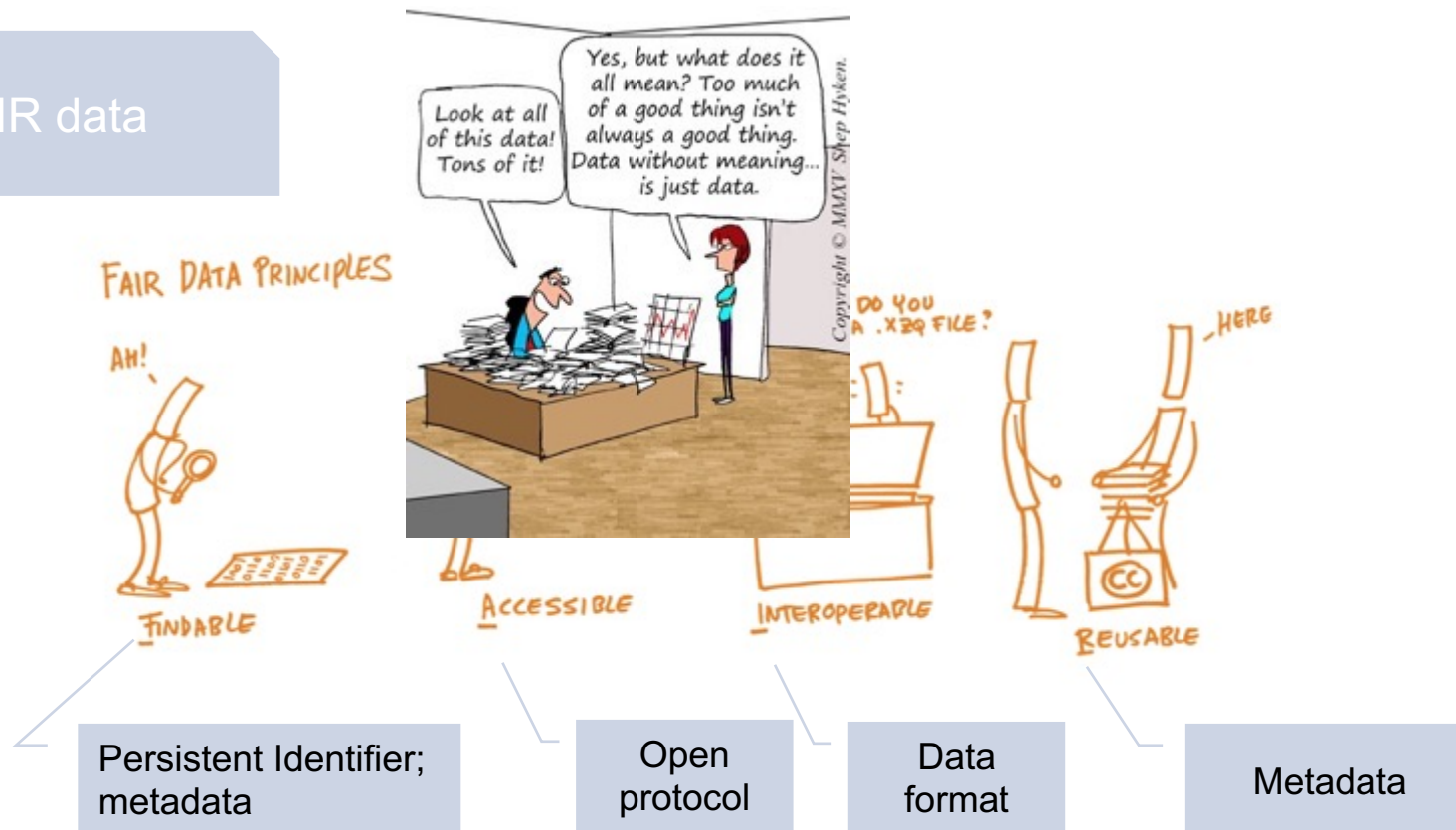*Scientific IT Services, ETH Zurich*

SIB Days, 24.06.2024

# Open Science and Open Research Data

Open Research Data

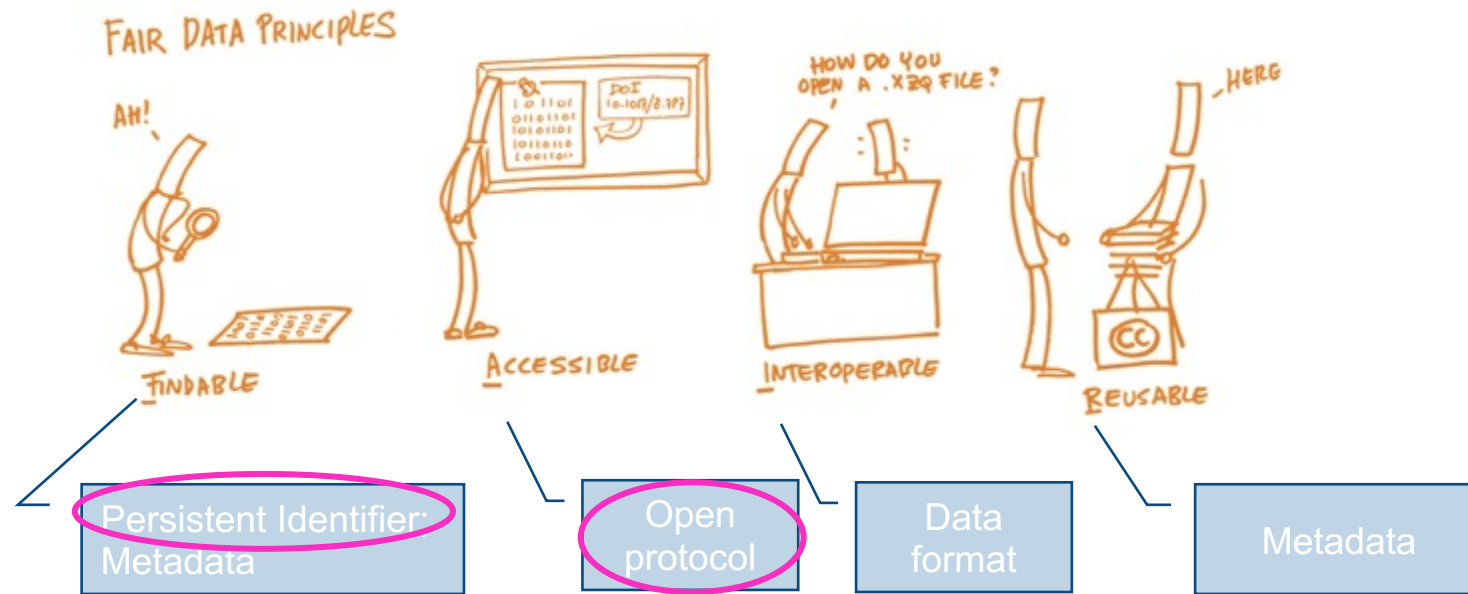**(!) Requirement from funding agencies, journals, academic institutions**

FAIR data



FAIR DATA PRINCIPLES

AH!

FINDABLE · ACCESSIBLE · INTEROPERABLE · REUSABLE

Persistent Identifier; metadata

Open protocol

Data format

Metadata

# The FAIR Data Principles

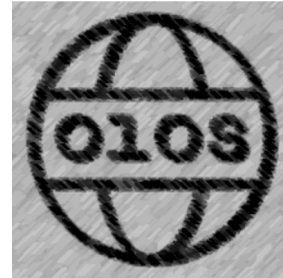| | |
|---|---|
| **F** | F1. (Meta)data are assigned a **globally unique** and **persistent identifie**r.<br>F2. Data are described with **rich metadata** (defined by R1 below).<br>F3. Metadata clearly and explicitly include the identifier of the data they describe.<br>F4. (Meta)data are registered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

https://www.go-fair.org/fair-principles/

# Prepare to meet the FAIR requirements when data are generated



FAIR DATA PRINCIPLES

AH!

FINDABLE

ACCESSIBLE

HOW DO YOU OPEN A .X3Q FILE?

INTEROPERABLE

HERE

REUSABLE

Persistent Identifier: Metadata

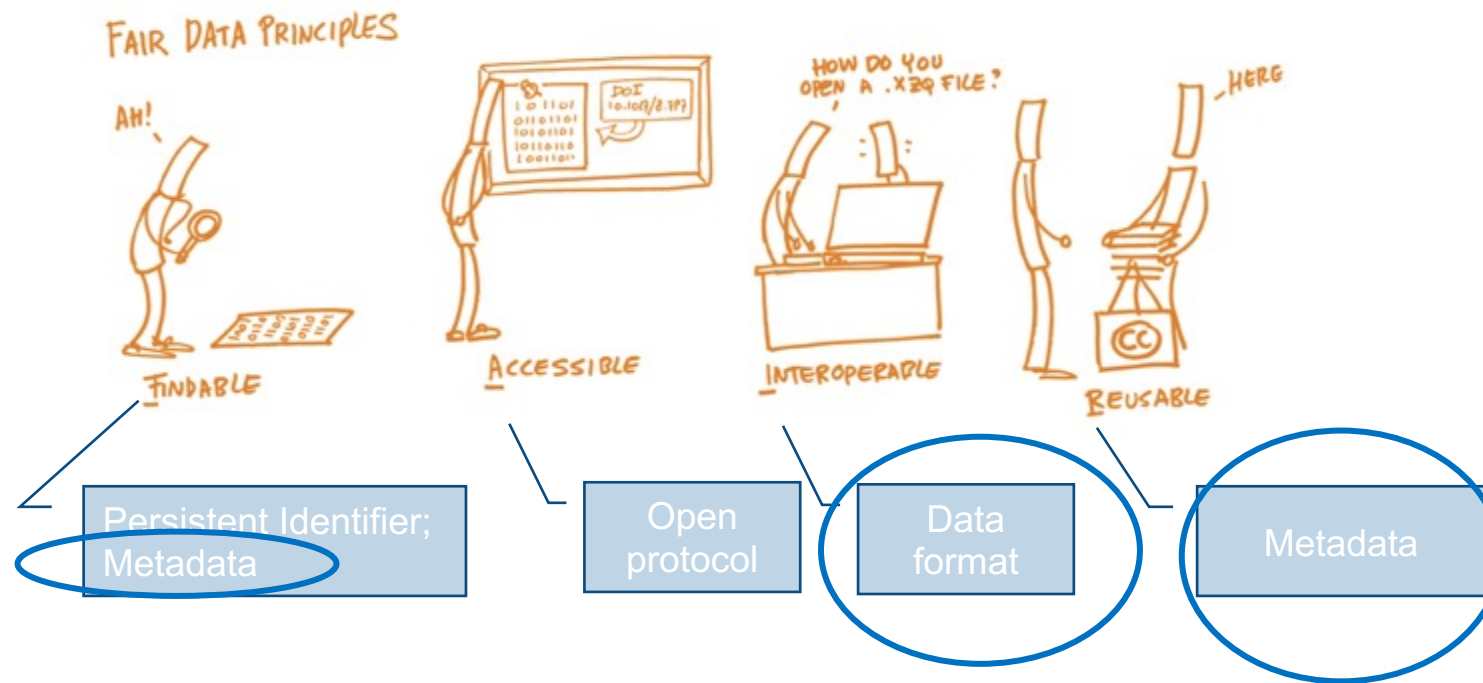Open protocol

Data format

Metadata

**i** When data are **published**

# How can we share data in a FAIR way?
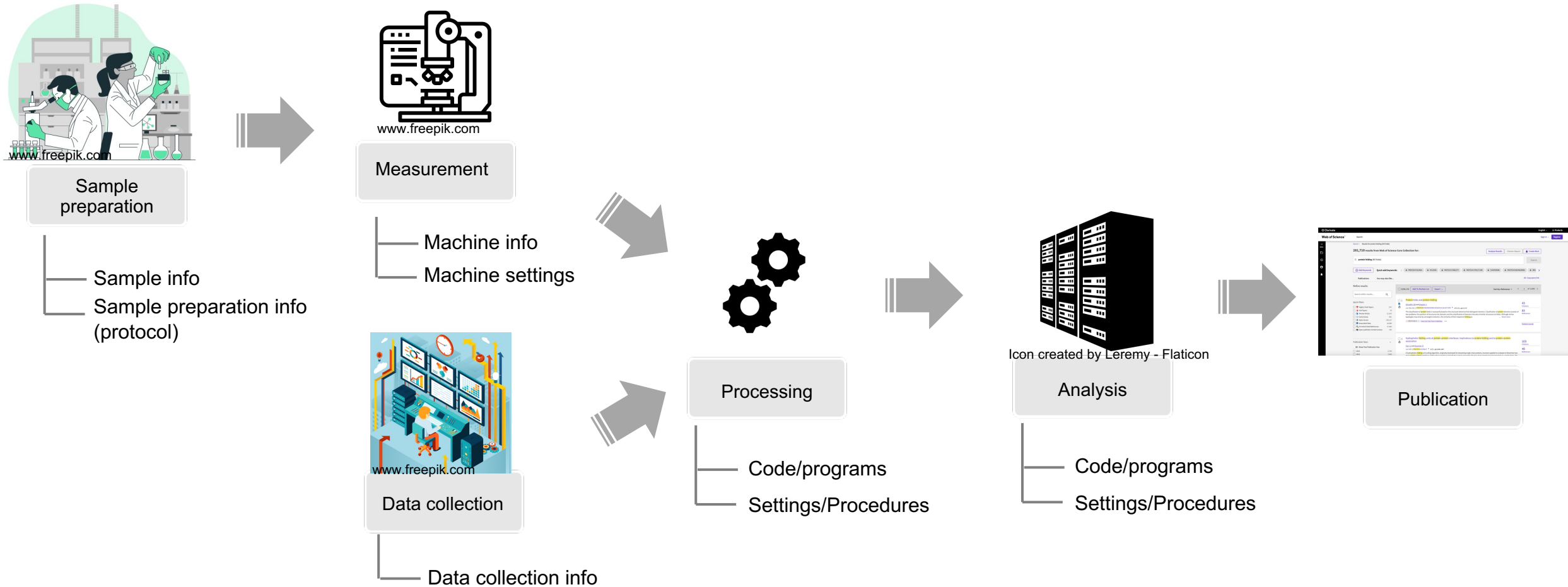
A few generic [data repositories recommended by SNSF](#)

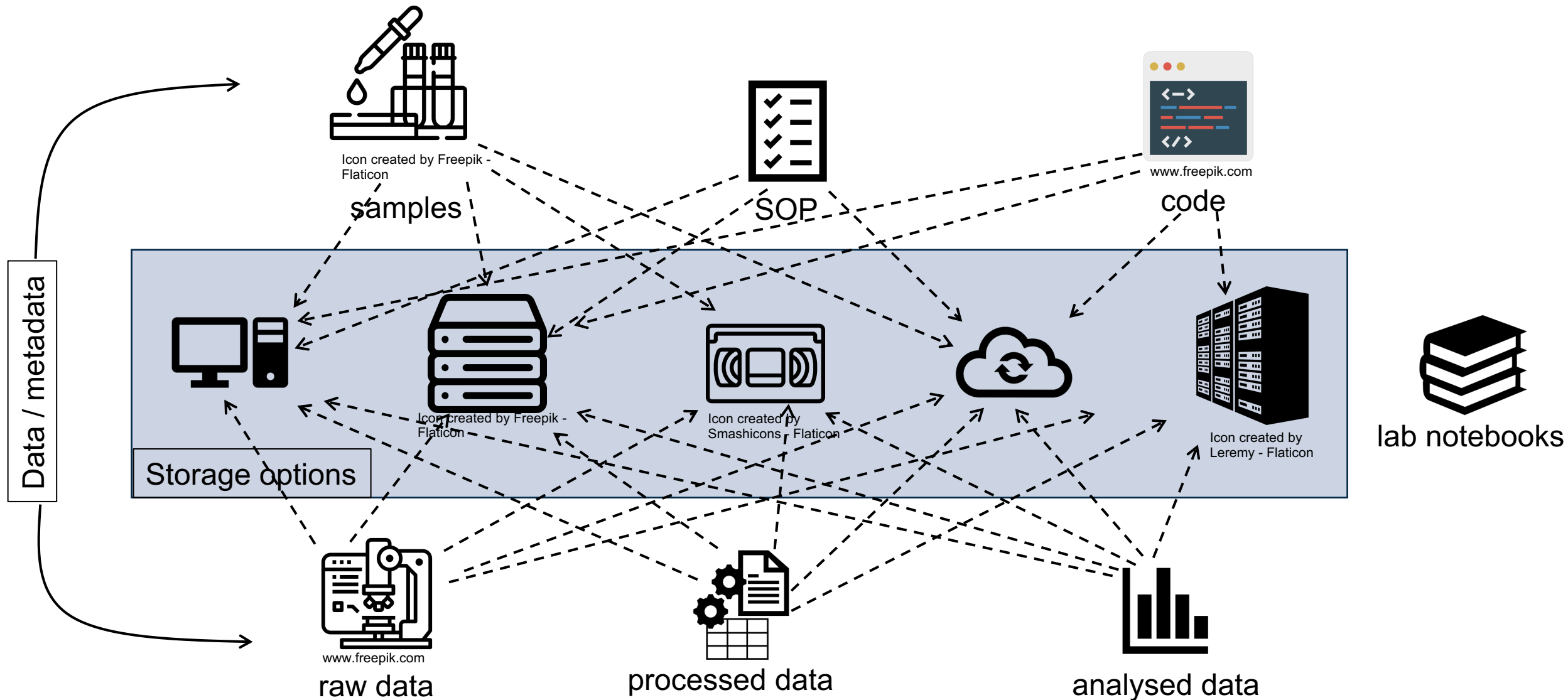# Prepare to meet the FAIR requirements when data are generated
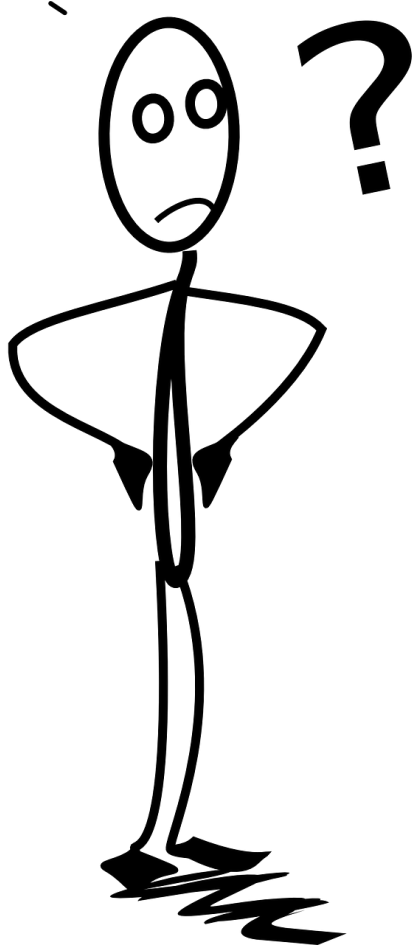


When data are **generated**

# Data and information generation during a research project



Sample preparation
- Sample info
- Sample preparation info (protocol)

Measurement
- Machine info
- Machine settings

Data collection
- Data collection info

Processing
- Code/programs
- Settings/Procedures

Analysis
- Code/programs
- Settings/Procedures

Publication

www.freepik.com

Icon created by Leremy - Flaticon

# The "Data Spread": a Common Scenario in Academic Research



samples

SOP

code

Icon created by Freepik - Flaticon

www.freepik.com

Data / metadata

Storage options

Icon created by Freepik - Flaticon

Icon created by Smashicons - Flaticon

Icon created by Leremy - Flaticon

lab notebooks

www.freepik.com

raw data

processed data

analysed data

Which tools can we use to manage all these data and information?

# Management of materials and samples



*Not scalable*

*No sharing*

*No efficient search*

*Easy to use*

**Spreadsheets / tabular files**

+*Scalable*

+ *Sharing*

+ *Search functionality*

 - *Require time for set up and maintenance*

**Database/ LIMS**

# Management of protocols

| Step 1 ……… | Step 2 ……… |
|---|---|
| Step 3 ……… | Step 4 ……. |
| Step n ……… | |

**Text files**

*Not scalable*
*No sharing*
*No efficient search*
*Easy to use*

WIKI

*Scalable*
*Sharing*
*Search functionality*
*Versioning*

**Database/ LIMS**

*Scalable*
*Sharing*
*Search functionality*
*Require time for set up and maintenance*

# Laboratory Information Management System (LIMS)



- ❑ LIMS are software for managing laboratory operations:
  - ▪ **sample tracking**
  - ▪ **sample data tracking**
  - ▪ **protocol management**

- ❑ Nowadays **LIMS are often combined with ELNs** in one platform.

# Management of research data files

Metadata
Raw data
- Creator
- Date
- Description

Metadata
Processed data
- Creator
- Date
- Description

Metadata
Results
- Creator
- Date
- Description

Metadata
Model
- Creator
- Date
- Description

Lab X
- Projects
  - P1_Title
    - Exp1
      - Raw
      - Processed
    - Exp2
      - Raw
      - Processed
    - Analysis
      - Code
      - Results
  - P2_Title
    - Exp1
    - Exp2
    - Analysis
- Protocols
  - Cloning
  - Imaging
  - Biochemistry
- Grants
  - SNF
  - H2020
- Presentations
  - 2016
  - 2017
  - 2018
- Photos
  - 2017-06_Retreat
  - 2016-12_XmasDinner

**Files / folders hierarchy**

SQLite
OME
openBIS

**Data management platform**

# Data management platforms

**Generic**

PostgreSQL

FileMaker
An Apple Subsidiary

SQLite

A
Microsoft Access

**Scientific**

PostGIS
Spatial PostgreSQL

openBIS

SLIMS
Agilent Technologies

rasdaman
raster data manager

OME

RSpace

- System that allows **structured organization** of data
- Data are described by **metadata**
- Usually more FAIR-compliant than Files / Folders
- Searchable, scalable, flexible
- Allows user rights management

# Metadata


www.digitalbevaring.dk

- ❑ **Metadata** is the *data about your data*

- ❑ (Machine-readable) **Metadata** is a key element of the **FAIR data** principles

- ❑ Use of structured metadata **facilitates data organization** and searches

- ❑ Existing **metadata schemas** are preferred (can be extended, if necessary)



https://rdamsc.bath.ac.uk/

# Metadata schema

- ❑ Defines the structure for the metadata.
- ❑ Schema defined by a scientific community to enable the best description of a resource type for their needs.
- ❑ **Generic metadata** schema examples:

## Table 1: DataCite Mandatory Properties

| ID | Property | Obligation |
|----|----------|------------|
| 1 | Identifier (with mandatory type sub-property) | M |
| 2 | Creator (with optional given name, family name, name identifier and affiliation sub-properties) | M |
| 3 | Title (with optional type sub-properties) | M |
| 4 | Publisher | M |
| 5 | PublicationYear | M |
| 10 | ResourceType (with mandatory general type description sub-property) | M |

https://schema.datacite.org/

## Table 2: DataCite Recommended and Optional Properties

| ID | Property | Obligation |
|----|----------|------------|
| 6 | Subject (with scheme sub-property) | R |
| 7 | Contributor (with optional given name, family name, name identifier and affiliation sub-properties) | R |
| 8 | Date (with type sub-property) | R |
| 9 | Language | O |
| 11 | AlternateIdentifier (with type sub-property) | O |
| 12 | RelatedIdentifier (with type and relation type sub-properties) | R |
| 13 | Size | O |
| 14 | Format | O |
| 15 | Version | O |
| 16 | Rights | O |
| 17 | Description (with type sub-property) | R |
| 18 | GeoLocation (with point, box and polygon sub-properties) | R |
| 19 | FundingReference (with name, identifier, and award related sub-properties) | O |

# Metadata schema

❑ Defines the structure for the metadata.
❑ Schema defined by a scientific community to enable the best description of a resource type for their needs.
❑ **Discipline-specific** metadata schema examples:



Microscopy images

Social Sciences

Ecology

Biology

SPHN Meta-Data Catalogue

Health data

# Electronic Laboratory Notebooks

❑ Scientists have always documented their findings in paper notebooks.



C. Darwin

Beagle voyage notebooks

*1831-1836*



ELN example

*2023*

❑ An **electronic laboratory notebook** (also known as **electronic lab notebook** or **ELN**) is a software program or package designed to replace more traditional paper laboratory notebook (*https://www.limswiki.org/index.php/Electronic_laboratory_notebook*).

# ELNs vs. paper notebooks

**+**

- Sharing
- Rights management
- Search functionality
- Easier to link digital data
- Can be backed up
- No issues with handwriting
- Track changes

**−**

- Learning curve
- Change in working mode required
- Time needed for introduction in lab

# Wiki and note-keeping applications





❑ Wikis and note-keeping applications often considered ELN solutions.

❑ Popular in academia for ease of use.

❑ These are a straight replacement of paper notebooks, with some added functionalities, but do not provide a solution for data management.

# Structured ELNs



- ❑ Additional functionalities compared to note-keeping applications (e.g. workflow management, chemical structures drawing, etc).

- ❑ Can be **discipline-specific** or **cross-disciplines**.

- ❑ Can have **LIMS** functionalities.

- ❑ Some systems offer an all-in-one solution for **RDM**.

- ❑ Can be integrated with third party applications.

# Where to start when choosing which ELN to use?



- ❑ Is it for personal use or group use?
- ❑ Can I/we use a cloud-based solution?
- ❑ Do I/we need specific features?
- ❑ What do I/we want to do with the ELN? (e.g. only write experimental descriptions, manage samples, manage data – how big?, etc.)
- ❑ Commercial v. open-source.
- ❑ Budget?
- ❑ Can I export my data?

# Some useful references on how to choose an ELN



https://www.nature.com/articles/d41586-018-05895-3



https://datamanagement.hms.harvard.edu/electronic-lab-notebooks



https://zenodo.org/records/4723753



https://eln-finder.ulb.tu-darmstadt.de/home

http://phdcomics.com/comics/archive_print.php?comicid=1689

**Reproducible Data Analysis**

# What do we mean by *Reproducibility*?

Computing Environment
& Infrastructure

Data

Workflow

Results

Code

« ***Reproducibility*** is **obtaining consistent results** using the same <u>input data</u>; <u>computational steps</u>, <u>methods</u>, and <u>code</u>; and <u>conditions of analysis</u>. This definition is synonymous with "computational reproducibility"… »

- All components need to be reproducible!

*National Academies of Sciences, Engineering, and Medicine. 2019. Reproducibility and Replicability in Science. https://doi.org/10.17226/25303.*

# Code management


Code



**Version control systems**

> Software tools specialized on managing and documenting changes to source code over time

> Necessary for managing large code bases

> Standard in professional software development



**Interactive notebooks**

> Applications that combine documentation, code, input and output generated by the code, e.g. graphs, plots (*Nature 515, 151–152*)

> Useful for exploratory data analysis and reproducibility

# Workflow Management



**Workflow management systems**

An incomplete list of **286** Computational Data Analysis Workflow Systems

https://github.com/common-workflow-language/common-workflow-language/wiki/Existing-Workflow-systems

A curated list of **109** Awesome Pipeline frameworks & libraries + **30** Workflow platforms

https://github.com/pditommaso/awesome-pipeline

# Reproducible Environment

**Problem**

Full reproducibility requires the possibility to recreate the system that was originally used to generate the results.

**Solution**

Bundle your application and all dependencies = Environment Isolation + Dependency Management

Environment and Package Management

**Containers**

**Virtual Machines**

**Environment & package management**

# Platforms for running code reproducibly



RRP

# Take home messages

❑ Efficient RDM during the lifetime of a project is necessary to meet FAIR data requirements.

❑ RDM should be an integral part of every researcher's daily work.

❑ Several tools are available for RDM. There is no „*one-fits-all*" solution, but every use-case should find the most appropriate solution(s) for them.

❑ Talk to the RDM experts in your institution!

# Scientific IT Services: bridging the gap between research and IT

# Who is Scientific IT Services of ETH Zurich?





High Performance Computing

Software Development

Scientific Computing & Data Co-Analysis

Scientific Visualization

Data Science & Machine Learning

Research Data Management

Confidential Research Data

Consulting & Training

- ❑ A section of ETHZ IT Services

- ❑ Around 50 experts in various areas of scientific computing

- ❑ With a background in different areas of science

ETH *zürich*

# Scientific IT Services: 4 groups

Research Platforms

High Performance Computing

Software development

Computational and Data Science Support

# Scientific IT Services

Research Platforms

- ❏ Research data management services (openBIS)
- ❏ Confidential research data
- ❏ LeonhardMed
- ❏ GFB sequencing core facility data support

openBIS

NV NGSvivo

# Scientific IT Services

- ❑ Scientific computing
- ❑ EULER cluster



High Performance
Computing

# Scientific IT Services







Software
development

❑ openBIS development
❑ Custom software development
❑ Scientific visualizations
❑ Scientific code support

# Scientific IT Services



**Gap**





MACHINE LEARNING

☐ Scientific IT consulting
☐ "Glue" between the SIS units
☐ Data science
☐ Data co-analysis
☐ Machine Learning/ AI

Computational and Data Science Support

# Bionformatics consulting



sequencing facilities

sequencing data

*researcher*

genomes and databases

results

publications

Research

data analysis

co-analysis

*SIS bioinformatics consultant*

High performance computing clusters

bioinformatic software

# Scientific IT Services



Coordinates between partners

Continuity of service

Single point of contact

**SIS Subscription**

Customized to your needs

Supported by a team of experts

Work on-site

F indable   A ccessible   I nteroperable   R eusable

added value

openBIS

NGSvivo

git

docker

The ETHZ Scientific IT Services data management solution for research groups

# openBIS: a complete solution towards FAIR data management



❑ Developed at **ETHZ** since 2007.

- *Description of experiments*
- *Description of processes and data analysis*

Electronic Lab Notebook

- *Store data connected to experimental description*

Data Management

Inventory Management

- *Samples*
- *Materials*
- *Reagents*
- *Equipment*
- *SOPs*

**https://openbis.ch**

# Inventory management



Lab equipment

Lab samples & materials

Lab procedures

Samples' storage manager

Barcode reader

User rights management

# Electronic Lab Notebook



Personal folder

User rights management

Entities relations

# Data management

❑ Data are always connected to experimental descriptions

# Data analysis: JupyterHub & MATLAB

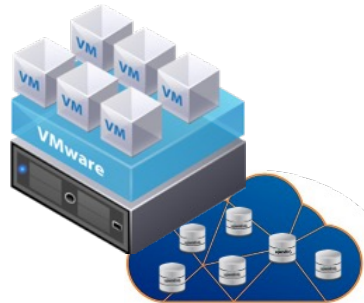# Data publication: export to ETH Research Collection & Zenodo

# RDM services offered BY ETHZ SIS

❑ Services for ETHZ researchers

❑ Services for Swiss academic scientists (openrdm.swiss)

openBIS installation on ETHZ infrastructure (ETH customers) /cloud (Swiss academic customers).

openBIS + OS maintenance & upgrades

Consulting

Tailored data modelling

Training

User support