

# Research Data Management in the life sciences

Elin Kronander

NBIS / SciLifeLab - ELIXIR Sweden

[data-management@scilifelab.se](mailto:data-management@scilifelab.se)

What comes to mind when you hear data management?

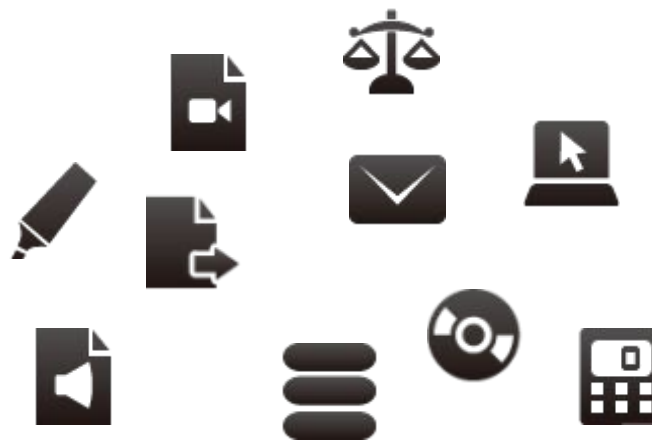


Data Life Cycle by [RDMkit](#) used under [CC-BY](#)



Go to [www.menti.com](http://www.menti.com) and use the code **7339 2159** or use this link:

<https://www.menti.com/alygf31121kr>



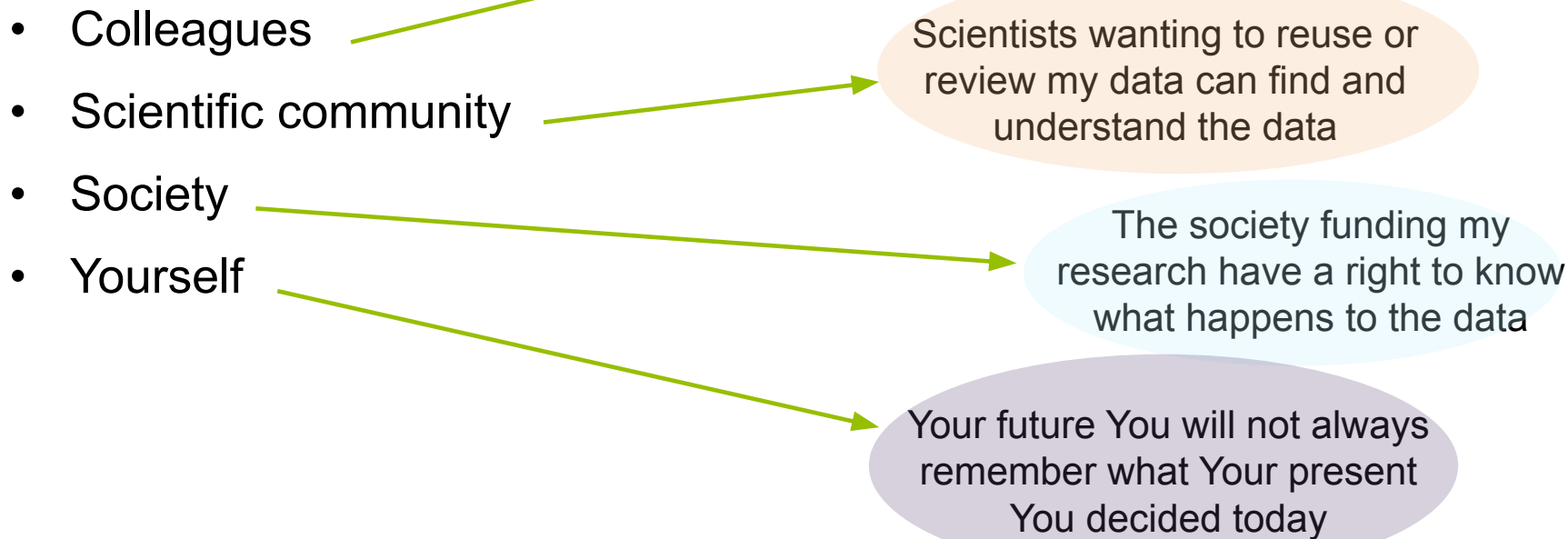
- raw data
- processed data
- data about data (metadata)
- ...

NB! not all data is digital



## Good data management practices in all phases of research

- Research documentation
- Data organisation
- Information security
- Ethics and legislation



---

*“Your primary  
collaborator is  
yourself six months  
from now, and your  
past self doesn’t  
answer e-mails,”*

-Rachael Ainsworth

---

# How do you know how an old result was generated?

# First step - Organization

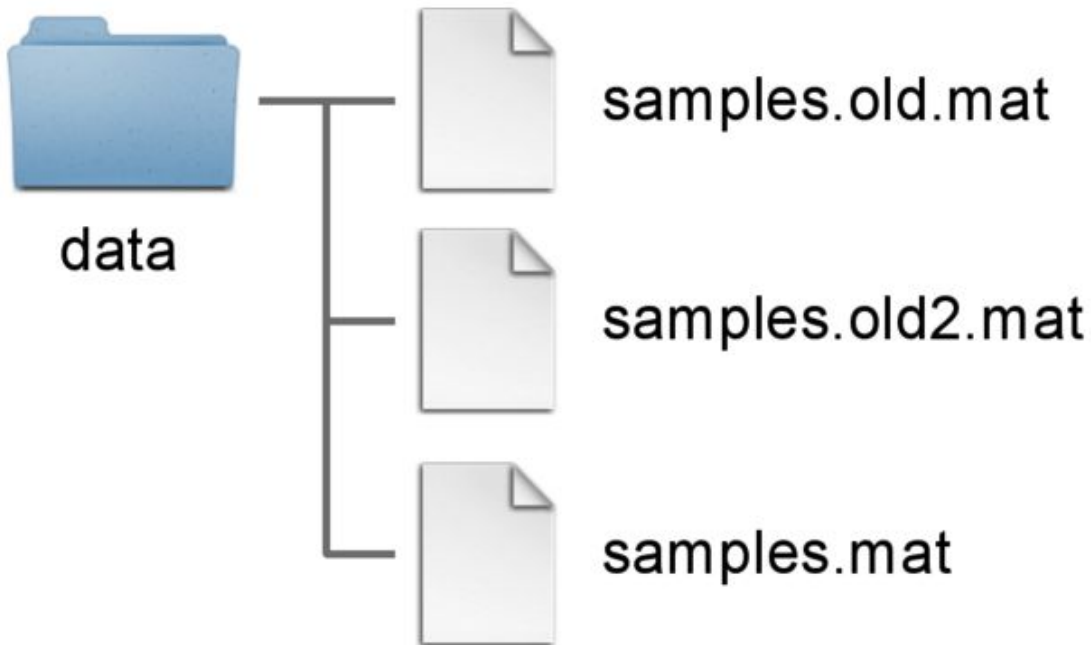




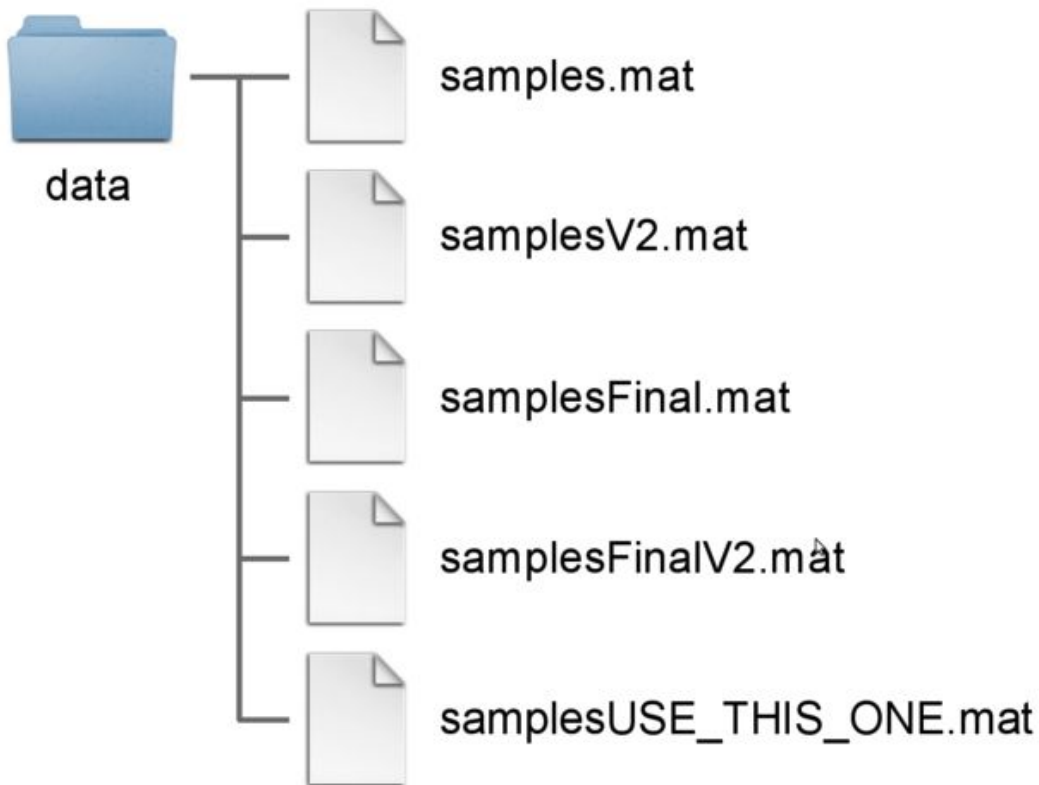
# I guess this is alright



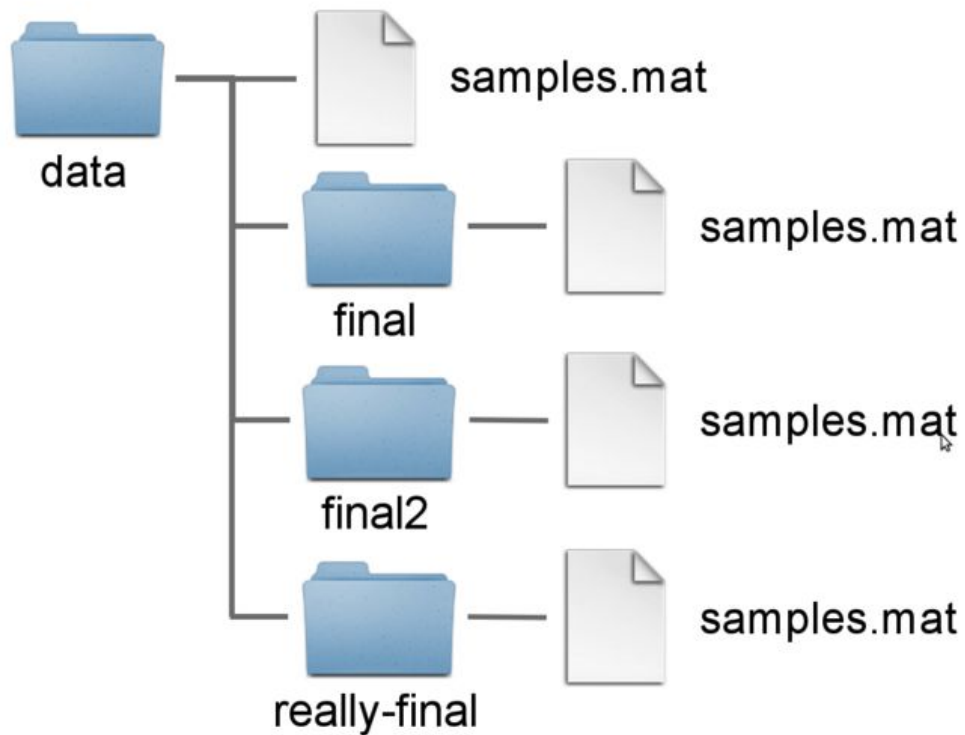
# I guess this is alright



# Which one is the most recent?

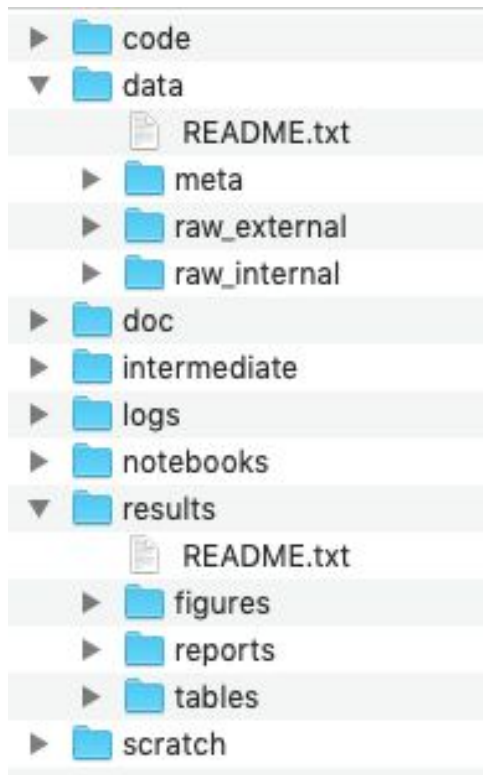


# Another (bad) common approach



# A possible solution





all code needed to go from input files to final results  
 raw and primary data, essentially all input files, **never** edit!

documentation for the study  
 output files from different analysis steps, *can be deleted*  
 logs from the different analysis steps

output from workflows and analyses

temporary files that can be safely *deleted or lost*

- 
- There is a **folder for the raw data**, which do not get altered, or intermixed with data that is the result of manual or programmatic manipulation. I.e., derived data is kept separate from raw data, and **raw data are not duplicated**.
  - **Code is kept separate from data.**
  - Use a **version control system** (at least for code) – e.g. **git**
  - There is a **scratch directory for experimentation**. Everything in the scratch directory can be deleted at any time without negative impact.
  - There should be a **README in every directory**, describing the purpose of the directory and its contents.
  - Use **file naming conventions** that makes it easy to find files and understand what they are (for humans and machines) and **document them**
  - Use **non-proprietary formats** – .csv rather than .xlsx
  - Etc...

- There is a **folder for the raw data**, which do not get altered, or intermixed with data that is the result of manual or programmatic manipulation. I.e., derived data is kept separate from raw data, and **raw data are not duplicated**.
- **Code is kept separate from data.**
- Use a **version control system** (at least for code) – e.g. **git**
- There is a **scratch directory for experimentation**. Everything in the scratch directory can be deleted at any time without negative impact.
- There should be a **README in every directory**, describing the purpose of the directory and its contents.
- Use **file naming conventions** that makes it easy to find files and understand what they are (for humans and machines) and **document them**
- Use **non-proprietary formats** – .csv rather than .xlsx
- Etc...



Two starting points for your file naming strategy are:

- A file name is a principal identifier of a file
- File naming strategy should be consistent in time and among different people

Principles for naming files:

1. Consider file name lengths – beware of OS limitations and full path names!
2. Make names human readable – name describes content of file
3. Make names machine readable – Avoid spaces, punctations, accented characters etc.
4. Explain file naming strategy in associated README files (stored in the same location)

Examples of a **poor** file name:

”Honeybee project, experiment 2 done in Helsinki, data file created on the second of December 2020”

File name - Runnew\_again\_2NDTRY.xls

*Explanation* - N/A

Examples of a **good** file name:

”Honeybee project, experiment 2 done in Helsinki, data file created on the second of December 2020”

File name - 20201202\_HB\_EXP2\_HEL\_DATA\_V03.csv

*Explanation* - Time\_ProjectAbbreviation\_ExperimentNumber\_  
Location\_TypeOfData\_VersionNumber

Names for files and folders should be *consistent* and *meaningful to yourself and collaborators*, allow for *easy tracking/searching*, and be *somewhat descriptive of content*.

Example:

LD\_phyA\_off\_t04\_2020-08-12\_norm.csv

Based on the name, the file could contain information about:

- LD - Long day sampling, of the
- phyA - Phytochrome A genotype, in a
- off - Medium without sucrose, at
- t04 - Time point 4,
- 2020-08-12 - Sampled on Aug 12th, 2020, with
- norm - Normalised data

But! Not obvious from the letters and words alone. Explanation is required - README

## Group discussion

The following examples contain files from an imaginary project

- *phyA/phyB* - genotypes
- *sXX* - sample number
- *LD/SD* - light conditions (Long Day, Short Day)
- *on/off* - different growth media (on sucrose, off sucrose)
- *date format* - sample date
- *tXX* - sample timepoint
- *raw, norm* - raw or normalised data

```

2020-07-14_s12_phyB_on_SD_t04.raw.xlsx
2020-07-14_s1_phyA_on_LD_t05.raw.xlsx
2020-07-14_s2_phyB_on_SD_t11.raw.xlsx
2020-08-12_s03_phyA_on_LD_t03.raw.xlsx
2020-08-12_s12_phyB_on_LD_t01.raw.xlsx
2020-08-13_s01_phyB_on_SD_t02.raw.xlsx
2020-7-12_s2_phyB_on_SD_t01.raw.xlsx
AUG-13_phyB_on_LD_s1_t11.raw.xlsx
JUL-31_phyB_on_LD_s1_t03.raw.xlsx
LD_phyA_off_t04_2020-08-12.norm.xlsx
LD_phyA_on_t04_2020-07-14.norm.xlsx
LD_phyB_off_t04_2020-08-12.norm.xlsx
LD_phyB_on_t04_2020-07-14.norm.xlsx
SD_phyB_off_t04_2020-08-13.norm.xlsx
SD_phyB_on_t04_2020-07-12.norm.xlsx
SD_phya_off_t04_2020-08-13.norm.xlsx
SD_phya_ons_t04_2020-07-12.norm.xlsx
ld_phyA_ons_t04_2020-08-12.norm.xlsx

```

1. Should dates be put first, and if not, why?
2. What is the difference between using leading 0 (zero) and not?
3. Is there a difference between using upper and lower case letters?
4. What are the effects if we, as in the example, mix naming conventions?

- *phyA/phyB* - genotypes
- *sXX* - sample number
- *LD/SD* - light conditions (Long Day, Short Day)
- *on/off* - different growth media (on sucrose, off sucrose)
- *date format* - sample date
- *tXX* - sample timepoint
- *raw, norm* - raw or normalised data

2020-07-14\_s12\_phyB\_on\_SD\_t04.raw.xlsx  
 2020-07-14\_s1\_phyA\_on\_LD\_t05.raw.xlsx  
 2020-07-14\_s2\_phyB\_on\_SD\_t11.raw.xlsx  
 2020-08-12\_s03\_phyA\_on\_LD\_t03.raw.xlsx  
 2020-08-12\_s12\_phyB\_on\_LD\_t01.raw.xlsx  
 2020-08-13\_s01\_phyB\_on\_SD\_t02.raw.xlsx  
 2020-7-12\_s2\_phyB\_on\_SD\_t01.raw.xlsx  
 AUG-13\_phyB\_on\_LD\_s1\_t11.raw.xlsx  
 JUL-31\_phyB\_on\_LD\_s1\_t03.raw.xlsx  
 LD\_phyA\_off\_t04\_2020-08-12.norm.xlsx  
 LD\_phyA\_on\_t04\_2020-07-14.norm.xlsx  
 LD\_phyB\_off\_t04\_2020-08-12.norm.xlsx  
 LD\_phyB\_on\_t04\_2020-07-14.norm.xlsx  
 SD\_phyB\_off\_t04\_2020-08-13.norm.xlsx  
 SD\_phyB\_on\_t04\_2020-07-12.norm.xlsx  
 SD\_phya\_off\_t04\_2020-08-13.norm.xlsx  
 SD\_phya\_ons\_t04\_2020-07-12.norm.xlsx  
 ld\_phyA\_ons\_t04\_2020-08-12.norm.xlsx

1. Should dates be put first, and if not, why?
2. What is the difference between using leading 0 (zero) and not?
3. Is there a difference between using upper and lower case letters?
4. What are the effects if we, as in the example, mix naming conventions?

1. Using dates as leading information in file names makes finding data quickly harder as the more interesting information may be samples or timepoints (unless date is crucial to data).
2. Without leading zeros, sorting will make 12 appear before 1 and 2.
3. Upper and lower cases may sort differently
4. Mixed naming conventions can make it difficult to locate particular files, and/or sort a large number of files.

2020-07-14\_s12\_phyB\_on\_SD\_t04.raw.xlsx  
 2020-07-14\_s1\_phyA\_on\_LD\_t05.raw.xlsx  
 2020-07-14\_s2\_phyB\_on\_SD\_t11.raw.xlsx  
 2020-08-12\_s03\_phyA\_on\_LD\_t03.raw.xlsx  
 2020-08-12\_s12\_phyB\_on\_LD\_t01.raw.xlsx  
 2020-08-13\_s01\_phyB\_on\_SD\_t02.raw.xlsx  
 2020-7-12\_s2\_phyB\_on\_SD\_t01.raw.xlsx  
 AUG-13\_phyB\_on\_LD\_s1\_t11.raw.xlsx  
 JUL-31\_phyB\_on\_LD\_s1\_t03.raw.xlsx  
 LD\_phyA\_off\_t04\_2020-08-12.norm.xlsx  
 LD\_phyA\_on\_t04\_2020-07-14.norm.xlsx  
 LD\_phyB\_off\_t04\_2020-08-12.norm.xlsx  
 LD\_phyB\_on\_t04\_2020-07-14.norm.xlsx  
 SD\_phyB\_off\_t04\_2020-08-13.norm.xlsx  
 SD\_phyB\_on\_t04\_2020-07-12.norm.xlsx  
 SD\_phya\_off\_t04\_2020-08-13.norm.xlsx  
 SD\_phya\_ons\_t04\_2020-07-12.norm.xlsx  
 ld\_phyA\_ons\_t04\_2020-08-12.norm.xlsx

- Using spaces (use \_ or - instead)
- Dots, commas and special characters (e.g. ~ ! @ # \$ % ^ & \* ( ) ` ; < > ? , [ ] { } ' " )
- Using language specific characters (e.g. óężé), unfortunately they still cause problems with most software or between operating systems (OS)
- Long names
- Consider repetitions in file names, e.g if directory name is Electron\_Microscopy\_Images, and file ELN\_MI\_IMG\_20200101.img then ELN\_MI\_IMG is redundant
- Deep paths with long names (i.e. deeply nested folders with long names), as archiving or moving between operating systems may fail



- 
- For dates use the YYYY-MM-DD standard and place at the end of the file UNLESS you need to organize your files chronologically
  - Include version number (if applicable), use leading zeroes (i.e.: v005 instead of v5).
  - make sure the end-letter file format extension is present at the end of the name (e.g. .txt, .md, .csv, .FASTQ)
  - Add a README file in your top directory which details your naming convention, directory structure and abbreviations

## Group discussion

What are examples of potential benefits of agreeing on a File Naming Convention for a project?

- Easier to process - Team members will not have to over think the file naming process
- Easier to facilitate access, retrieval and storage of files
- Easier to browse through files, saving time and effort
- Harder to lose!
- Having logical and known naming conventions in place can also help you with version control.
- Check for obsolete or duplicate records

Names for files and folders should be *consistent* and *meaningful to yourself and collaborators*, allow for *easy tracking/searching*, and be *somewhat descriptive of content*.

Examples of a **good** file name:

"Honeybee project, experiment 2 done in Helsinki, data file created on the second of December 2020"

File name - 20201202\_HB\_EXP2\_HEL\_DATA\_V03.csv

*Explanation* -

Explanation is required - README

---

A file usually defined as the starting point of information about something (attracts attention!)

Using them as documentation files for:

Folder level – Explaining folder contents, naming, file history, organisation/structure etc

Data – Explaining file names and contents

README.txt — Edited

README for the Honey Bee project field measurements data folder

This folder contains raw data collected manually from field measurements over several time points.

file naming convention:

Time\_ProjectAbbreviation\_ExperimentNumber\_Location\_TypeOfData\_VersionNumber

For example

20201202\_HB\_EXP2\_HEL\_DATA\_V01.csv

20201202\_HB\_EXP2\_HEL\_DATA\_V03.csv

20201202\_HB\_EXP2\_HEL\_DESCR\_V03.csv

Time - is the date at the start of experiment YYYY-MM-DD

ProjectAbbreviation - is HB for Honey Bee

ExperimentNumber - EXP1, EXP2, EXP3 or EXP4

Location - refers to a city, HEL for Helsinki, STO for Stockholm or OSL for Oslo

TypeOfData - DATA for numeric measurements, DESCR for qualitative values

VersionNumber - Version number is increased each time point of data collection as V01, V02 and so on.

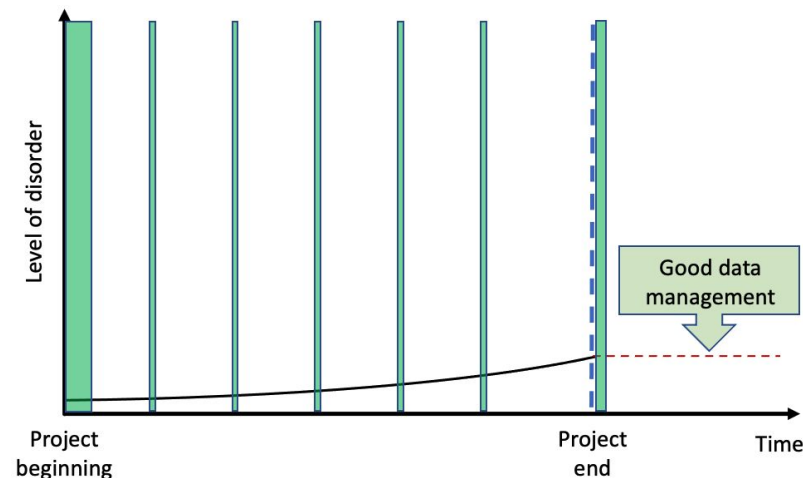
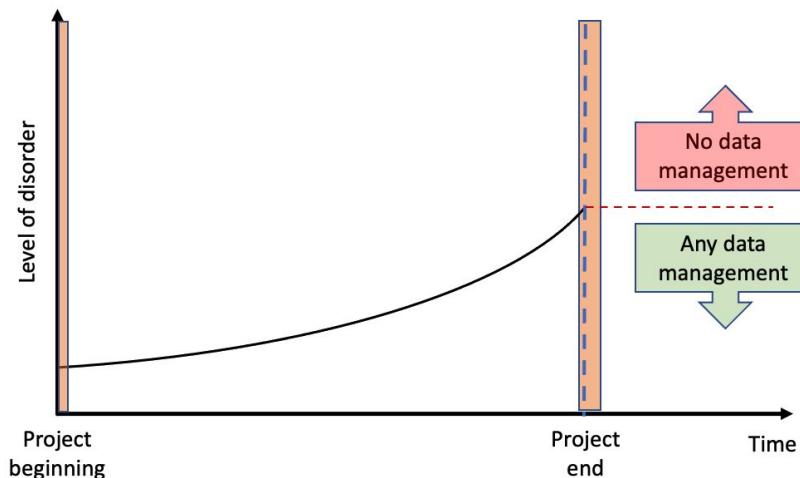
## Discussion

Think of an example where you would have benefited from having access to a README-file when working with data. Describe to your neighbor what you would have wanted such a file to contain.

Files will become unorganised over time (particularly downloads and/or desktop folders)

Files can multiply across folders and versions, decreasing findability

Organising will reduce clutter and maintenance requirements over time

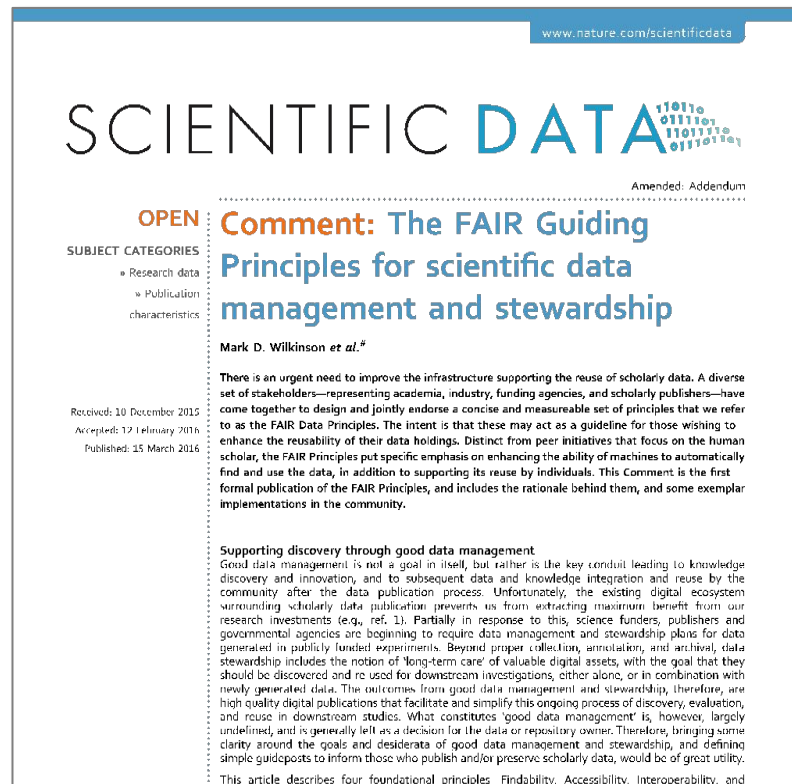


- ❑ **Secure/organise data & analyses**, by using folder structures, file naming conventions and README files, managing back-ups, access restrictions, versioning, docs, scripts and transcripts
- ❑ **Deposit and share data** using restricted or public access data repositories that promote FAIR data principles
- ❑ **Adhere to community standards**, such as file formats, data dictionaries, controlled vocabularies and metadata
- ❑ **Maintain a Data Management Plan**, outlining the project's data management practices





- Promote **efficient data discovery and reuse** by providing guidelines to make digital resources
  - Findable
  - Accessible
  - Interoperable
  - Reusable
- Address aspects **enabling software and infrastructure** to automatically find and use research data



www.nature.com/scientificdata

## SCIENTIFIC DATA

Amended: Addendum

**OPEN** Comment: The FAIR Guiding Principles for scientific data management and stewardship

**SUBJECT CATEGORIES**

- Research data
- Publication characteristics

Received: 10 December 2015  
Accepted: 17 February 2016  
Published: 15 March 2016

Mark D. Wilkinson *et al.*<sup>§</sup>

There is an urgent need to improve the infrastructure supporting the reuse of scholarly data. A diverse set of stakeholders—representing academia, industry, funding agencies, and scholarly publishers—have come together to design and jointly endorse a concise and measurable set of principles that we refer to as the FAIR Data Principles. The intent is that these may act as a guideline for those wishing to enhance the reusability of their data holdings. Distinct from peer initiatives that focus on the human scholar, the FAIR Principles put specific emphasis on enhancing the ability of machines to automatically find and use the data, in addition to supporting its reuse by individuals. This Comment is the first formal publication of the FAIR Principles, and includes the rationale behind them, and some exemplar implementations in the community.

**Supporting discovery through good data management**

Good data management is not a goal in itself, but rather is the key conduit leading to knowledge discovery and innovation, and to subsequent data and knowledge integration and reuse by the community after the data publication process. Unfortunately, the existing digital ecosystem surrounding scholarly data publication prevents us from extracting maximum benefit from our research investments (e.g., ref. 1). Partially in response to this, science funders, publishers and governmental agencies are beginning to require data management and stewardship plans for data generated in publicly funded experiments. Beyond proper collection, annotation, and archival, data stewardship includes the notion of ‘long-term care’ of valuable digital assets, with the goal that they should be discovered and re-used for downstream investigations, either alone, or in combination with newly generated data. The outcomes from good data management and stewardship, therefore, are high quality digital publications that facilitate and simplify this ongoing process of discovery, evaluation, and reuse in downstream studies. What constitutes ‘good data management’ is, however, largely undefined, and is generally left as a decision for the data or repository owner. Therefore, bringing some clarity around the goals and desiderata of good data management and stewardship, and defining simple guideposts to inform those who publish and/or preserve scholarly data, would be of great utility.

This article describes four foundational principles: Findability, Accessibility, Interoperability, and

Wilkinson et al. (2016). The FAIR Guiding Principles for scientific data management and stewardship. *Scientific Data*, 3(1), 160018. doi:10.1038/sdata.2016.18

## Vision: As open as possible, as closed as necessary by 2026

research process. Already existing data that have only been used in their original form and that are already managed and made accessible by another actor are not covered by this recommendation.

### Metadata should also be published with open access

Both research data and data describing research data (known as metadata) should be published with open access. If there are obstacles to publishing research data, the focus should in the first instance be on making metadata openly accessible on the internet. In this way, users can find information on what research data exists, even when there are obstacles to open publication, for example lack of a suitable publication platform or technical limitations that prevent all data from being published.

### Publication according to the FAIR principles

Publication of research data can be done using various digital platforms, for example via the higher education institution where the research is conducted or via other relevant national and/or international portals, infrastructures and similar organisations and platforms. The publication of research data shall always be based on the FAIR principles.

### The Swedish Research Council's recommendation on data management according to FAIR

The Swedish Research Council recommends that the research data produced through research are managed according to the FAIR principles, clarified via the criteria developed by the Swedish Research Council to achieve FAIR data.

The FAIR principles should be implemented taking into account applicable legislation, and, as far as is possible and applicable, based on the technical, organisational and/or discipline-specific preconditions that apply.

The recommendations relates in the first instance to research data (and metadata) financed by public funds that can be published with open access, but the application of the FAIR principles can be made broader than this, and be used also for research data that cannot be published entirely openly. The recommendation on data management according to FAIR is overarching, and aims to create a common starting point for the implementation of FAIR data management.

*[...] The publication of research data shall always be based on the FAIR principles.[...]*

*The Swedish Research Council's recommendation on data management according to FAIR*

*The Swedish Research Council recommends that the research data produced through research are managed according to the FAIR principles, clarified via the criteria developed by the Swedish Research Council to achieve FAIR data. [...]*



Data Life Cycle by [RDMkit](#) used under [CC-BY](#)

- **FAIR data  $\neq$  Open data**  
Data can be Open without being FAIR  
Data can be FAIR without being open  
*“As open as possible, as closed as necessary”*
- **FAIR software/FAIR training materials**
- **Data can be more or less FAIR**

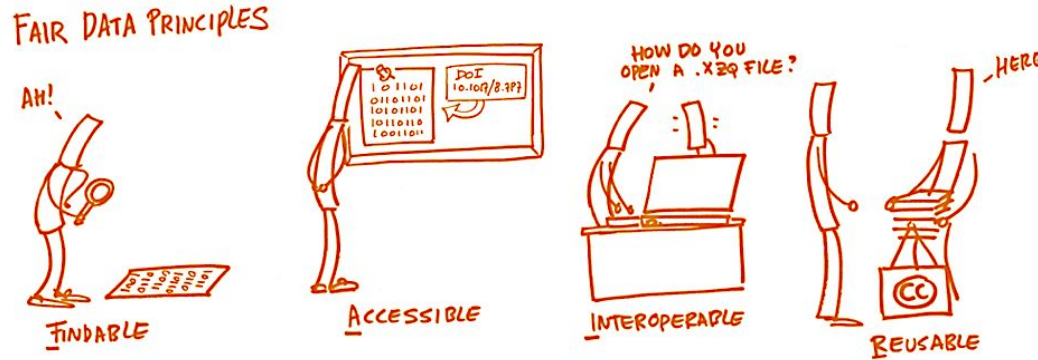


Image: <https://book.fosteropenscience.eu/>


- **FAIRify by README**, adopting good practices for data organization, makes research data more FAIR
- **FAIRify by planning**, thinking ahead and continuously document your strategies in a Data Management Plan using a guiding tool <https://dsw.scilifelab.se/>
- **Deposit data** in a repository
- **Get support by data stewards**

- Guide writing a data management plan
- Identify a suitable repository for publishing your data
- Assist during the submission process when publishing your data and code
- Advice on what needs to be done when working with sensitive human data
- Advice on describing data with proper metadata for documentation and publishing
- Data transfers, data organisation, backup, and security procedures



## Contact us

- [data-guidelines.scilifelab.se](https://data-guidelines.scilifelab.se)
- [data-management@scilifelab.se](mailto:data-management@scilifelab.se)



**SciLifeLab RDM Guidelines**

Get support About Contact

Knowledge hub for the management of life science research data in Sweden

The purpose of these guidelines is to serve as an information resource to life science researchers in Sweden regarding research data management.


Home Research data life cycle Resources Topics

## Research Data Management

Research Data Management (RDM) concerns the organisation, storage, preservation, and sharing of data that is collected and analysed during a research project. Proper planning and management of research data will make project management easier and more efficient while projects are being performed. It also facilitates sharing and allows others to validate as well as reuse the data.

### Research data life cycle

The research data life cycle can be divided into several phases as seen in the wheel below; **plan, collect, process, analyse, preserve, share and reuse**. Click on a section of the wheel below to get an introduction to that phase of the research data life cycle, including information on relevant resources and training material.



RDMkit

### Get Support

Do you need support with research data management?

We offer support to anyone involved in life science research that is affiliated with a Swedish university or research institute.

[Click here to get support](#)

Seminar series:

### SciLifeLab Data Platform- how we support data-driven life science research

Join the next event in our SciLifeLab Data Management seminar series.

Next date: Nov 21, 2023

[More information on the event page](#)

### Events & Training

Upcoming conferences, webinars, workshops, and training opportunities in Sweden related to data-driven life science can be found on the SciLifeLab Data Platform.

RDM life cycle from RDMkit licensed under Creative Commons Attribution 4.0 International License.

<https://data-guidelines.scilifelab.se/>

[data-management@scilifelab.se](mailto:data-management@scilifelab.se)

**21 November** - SciLifeLab Data Platform- how we support data-driven life science research

**22 November** - Data Management Plans in practice - research funder perspectives and practical demos

**30 November** - SciLifeLab FAIR storage – how to apply for storage resources

**23-25 April 2024** - Introduction to Data Management Practices course  
Stockholm