

Data organisation practices

Introduction to Data Management Practices course

NBIS DM Team

data@nbis.se

<https://nbisweden.github.io/module-organising-data-dm-practices/>



- What to consider for maintaining data organization strategies in a project
- What to consider when settling for a file structure
- Understanding good practices for data storage, processing and documentation (**FAIR-ification**)



Credit: This image was created by Scriberia for The Turing Way community and is used under a CC-BY licence.

What measures do you take in order to avoid file chaos in your data organisation?



Digitalbevaring.dk

- Research **data is a core component** of any research project or publication.
- Good data management practices are **important in all phases** of research
 - Ethics and legislation
 - Information security
 - Research documentation
 - Project organisation
- Research data needs to stay authentic and be secured **beyond the project's** time frame



Digitalbevaring.dk

The Data Lifecycle

Data *content* is subject to changes at all phases of its life cycle

- Creation
- Error corrections
- New variables
- Changing file formats
- Etc.

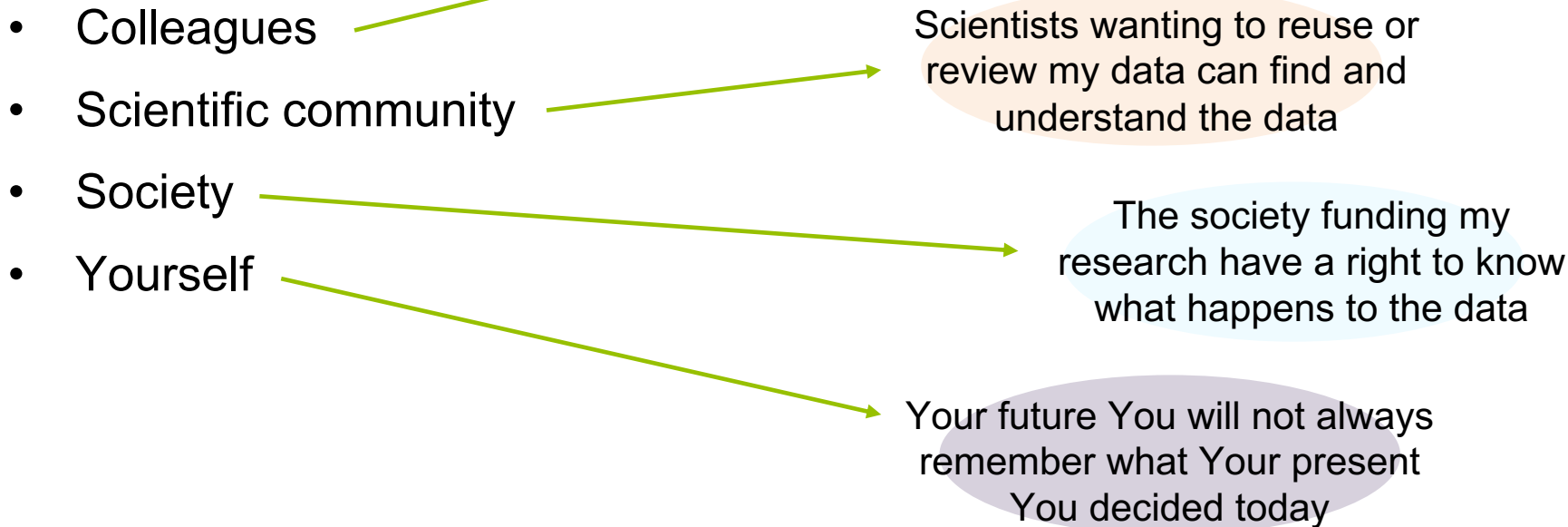


Data *structure* may change at all phases of data life cycle

- Data splits
- Re-organization
- Storage change
- Etc.

- Large impact in the planning phase
 - Intended folder structure
 - Clustering of files
 - File naming convention
 - Standards for dates/measures
 - Documentation procedures
- Prepare your project for receiving data by making a **Data Management Plan**



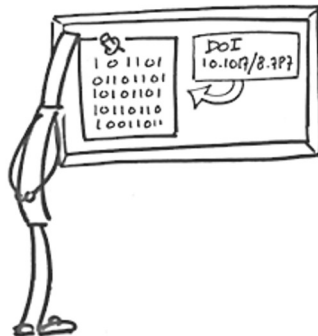


FAIR DATA PRINCIPLES

AH!



FINDABLE

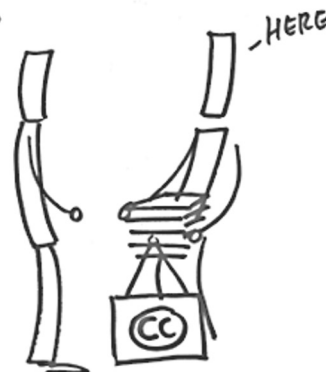


ACCESSIBLE

HOW DO YOU
OPEN A .XZQ FILE?



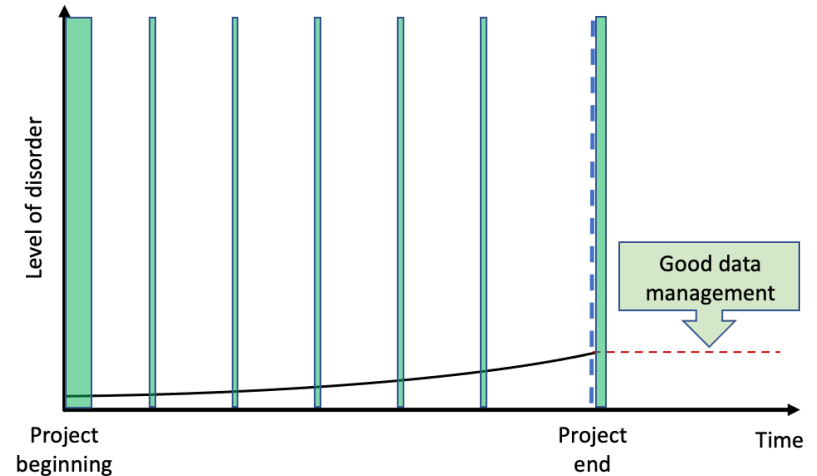
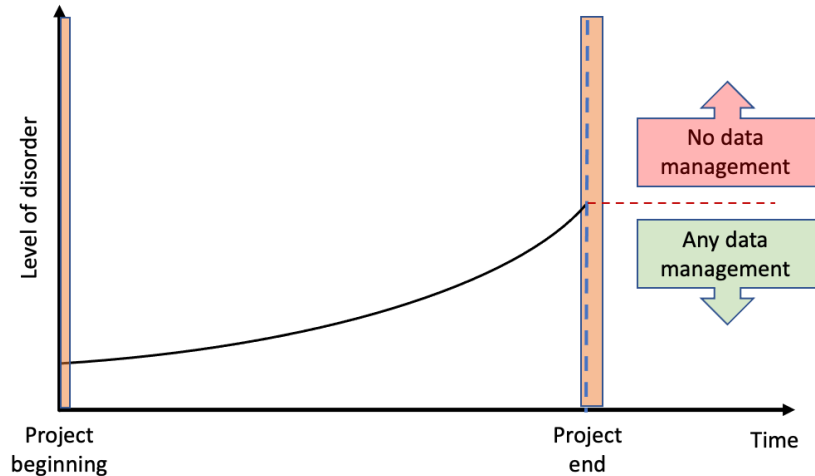
INTEROPERABLE



REUSABLE

Adopting good practices for data organization, makes research data more **FAIR**

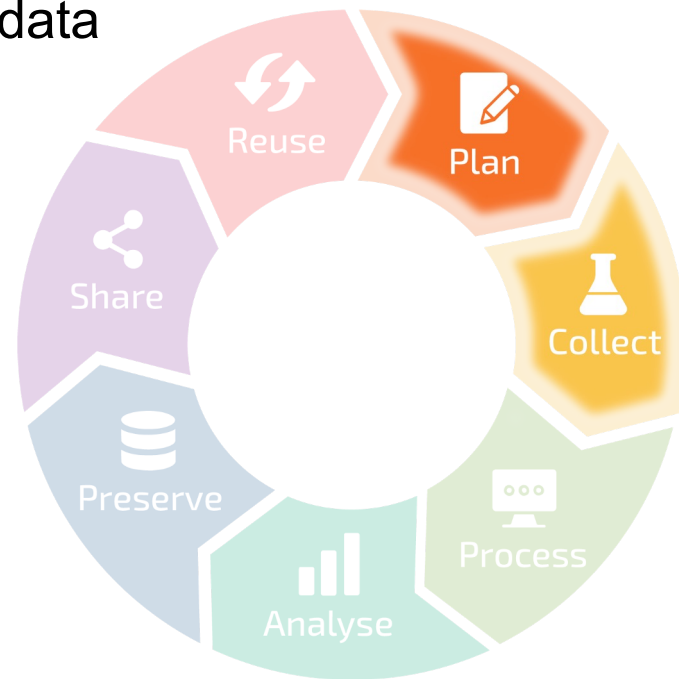
Planned data organization reduces disorder and optimizes time vs. effort



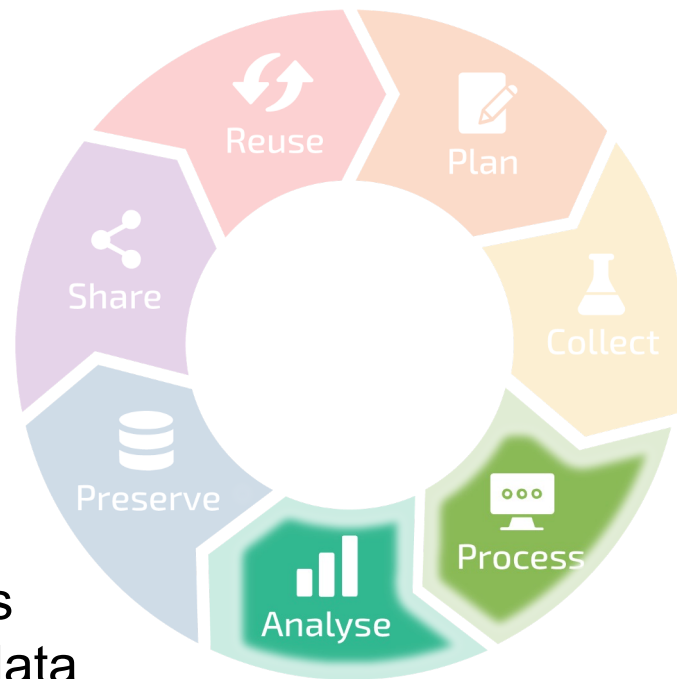
Rank the following data organization steps from 1-5 (1 being the one you believe you think is most important, and 5 the least). Also mark with an "X" the steps you have implemented in your own research.

- File naming convention
- Folder naming convention
- File versioning system
- File organisation documentation (README.txt)
- File and folder maintenance (moving, deleting)

- Raw data is the purest form of scientific data
- NOT for analysis
- Preservation and access restriction
- Extense of raw data package(s)
- Versioning
- Documentation
- Status of sub-selections
- Associated data
- Long term / Publication package

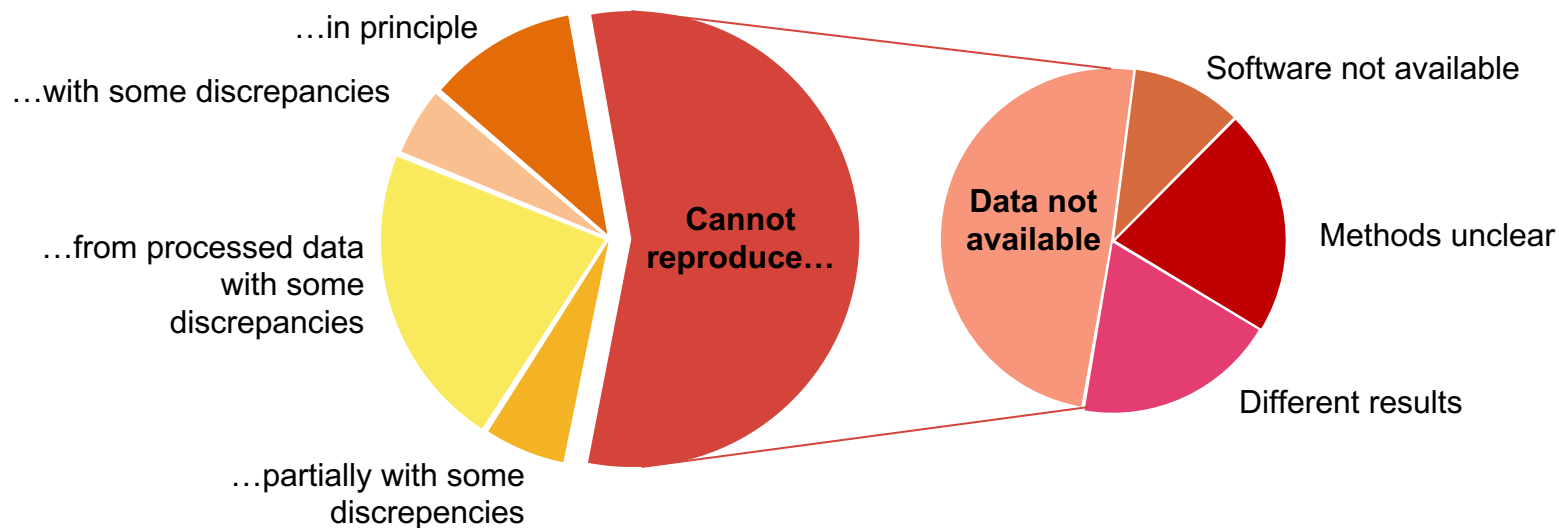


- Data (subset) + Analyses = Results
 - Document for reproducibility
 - Describe analyses, avoid "default"
 - Self-sufficient descriptions
-
- Rich metadata
 - Store code close to data
 - Document and publish code for analyses
 - Cross-reference code and data in metadata



Reproduction of data analyses in 18 articles on microarray-based gene expression profiling published in Nature Genetics in 2005–2006:

Can reproduce...



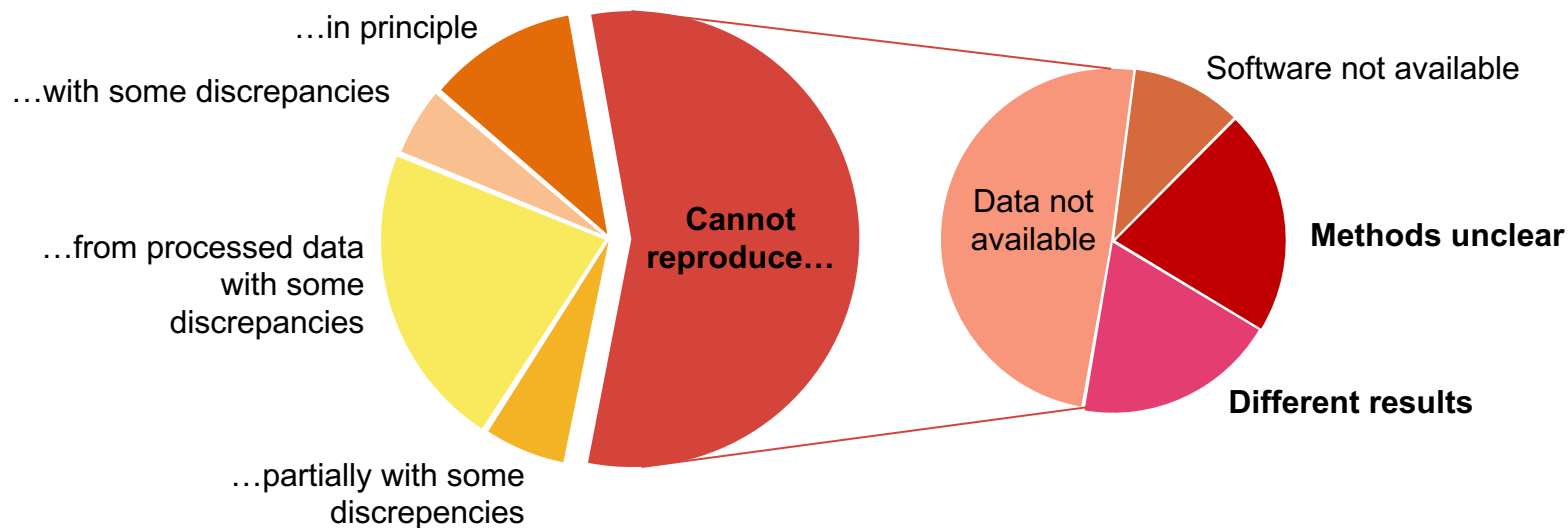
Summary of the efforts to replicate the published analyses.

Adopted from: Ioannidis et al. Repeatability of published microarray gene expression analyses.

Nature Genetics **41** (2009) doi:10.1038/ng.295

Reproduction of data analyses in 18 articles on microarray-based gene expression profiling published in Nature Genetics in 2005–2006:

Can reproduce...



Summary of the efforts to replicate the published analyses.

Adopted from: Ioannidis et al. Repeatability of published microarray gene expression analyses.

Nature Genetics **41** (2009) doi:10.1038/ng.295

- Tabular data is not a data type, but a mode to organise data
 - Do not mix with notes/code/analyses
 - Data visualization \neq data use
 - Missing values
-
- Interoperable as .csv or .tsv



Tabular data

- Column = Variable
- Row = Observation
- Cell = Value

Open Access training					
Date	Length (hours)	Registered	Attended	Delivered by	Canceled
16/01/17	1	26	23	JM	N
05/02/17	1	38	26	JM	N
17/02/17	1	19	25	PG	N
07/03/17	1	27	17	JM	N
29/03/17	1	32	15	PG	N
02/04/17	1	41		PG	Y
24/04/17	2	44	44	JM	N
25/05/17	1	43	37	PG	N
16/06/17	1	15	15	JM	N

Do not:

- Spatially distribute data
- Combine values in cells
- Split compatible data in tables
- Use colors

[illegible]

✓ Raw means raw!

✓ Tidy data tables

One cell—one value

One column—one variable

One row—one observation

✓ Beware of Excel “features”

Misguided “auto-corrections” of
dates, casing, numbers etc.

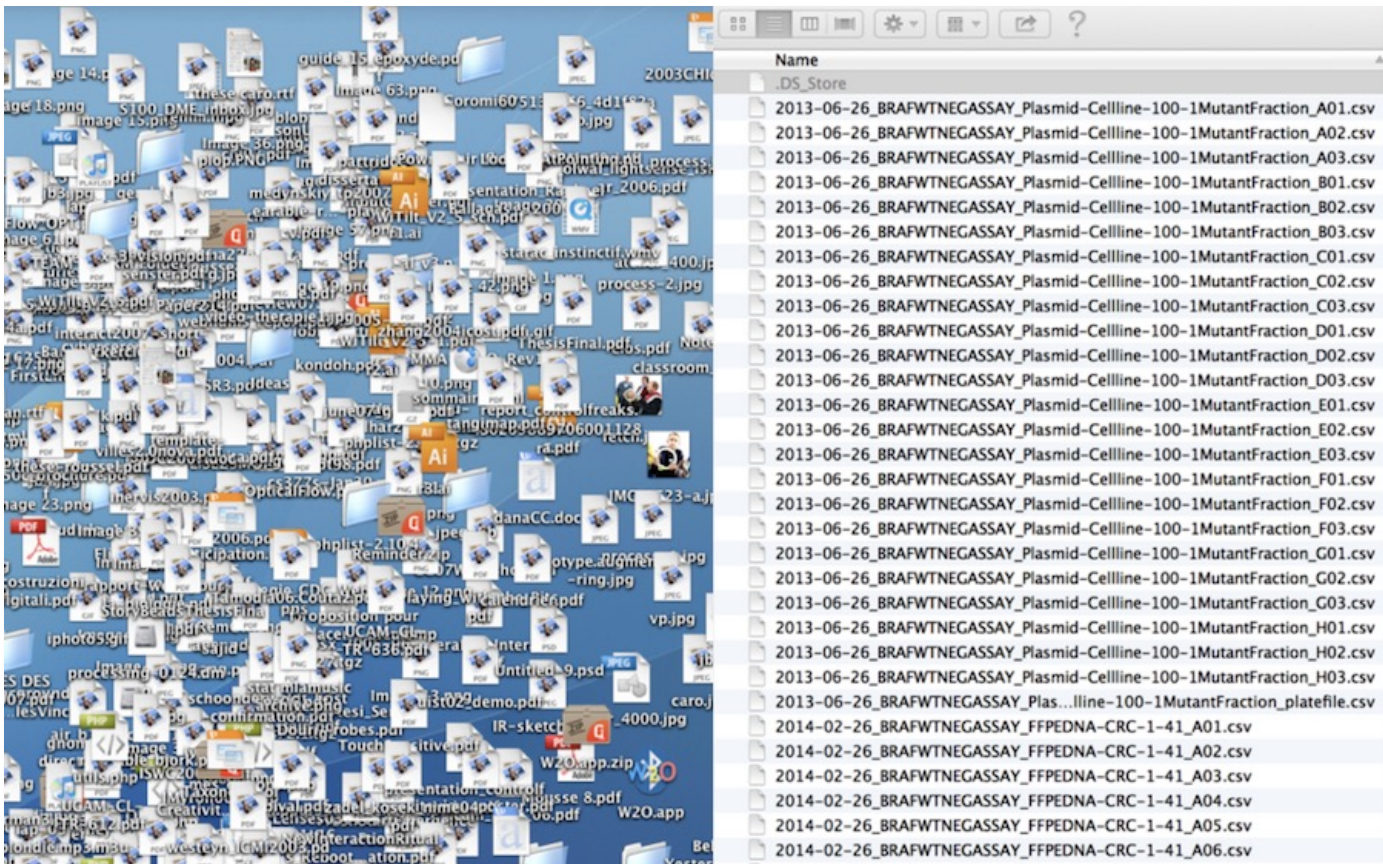
Misaligned formulas

Limited numerical precision

Limited number of rows/columns

	A	B	C	D	E	F	G	H	I	J	K
1	data							analysis			
2	id	biomarker1	biomarker2	biomarker3	biomarker4			variation	ave	problem	
3	81	0.08502	0.07002	0.07735	0.07746			0.008	0.0775		
4	82	0.0658	0.06859	0.06958	0.06799			0.002	0.068	no	
5	83	0.07757	0.07497	0.0801	0.07755			0.003	0.0775		
6	84	0.07185	0.06957	0.07474	0.07205			0.003	0.0721	yes	
7	85	0.06959	0.07361	0.07113	0.07145			0.002	0.0714	maybe	
8	86	0.09291	0.10439	0.09425	0.09718			0.006	0.0972		
9	87	0.07878	0.08143	0.07203	0.07742			0.005	0.0774		
10	88	0.07907	0.077	0.08227	0.07944			0.003	0.0794		
11	89	0.07299	0.07616	0.08131	0.07682			0.004	0.0768		
12	90	0.07487	0.0664	0.0671	0.06946			0.005	0.0695		
13											
14	mean	0.076845	0.076214	0.076986	0.076682						
15								biomarker QC			
16	notes							b1	b2	b3	b4
17	* patient id86 may need removing due to missing notes							0.46336967	0.875281336	0.918250702	0.14953926

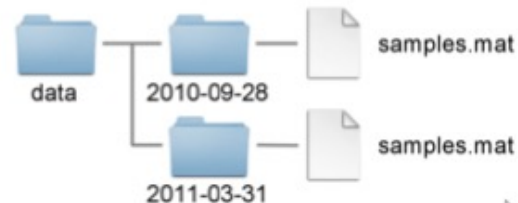
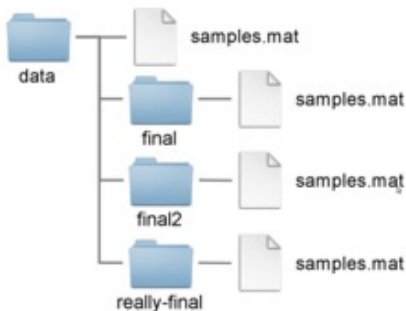
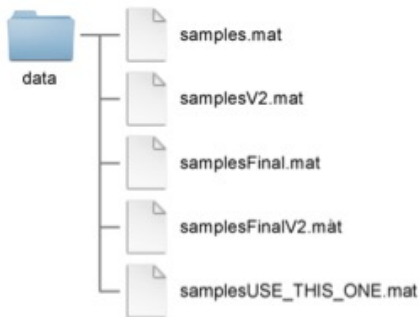
Organising files and folders



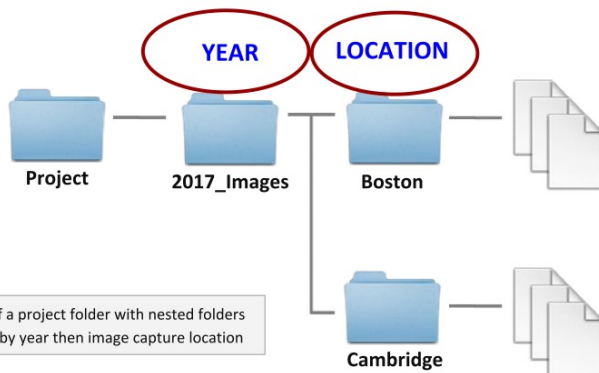
Your primary collaborator is yourself from 6 months ago, and that person is really difficult to communicate with!

Good practices

- Organise files hierarchically
- Use folders to divide files into categories
- Choose a file naming strategy
- Create documentation files

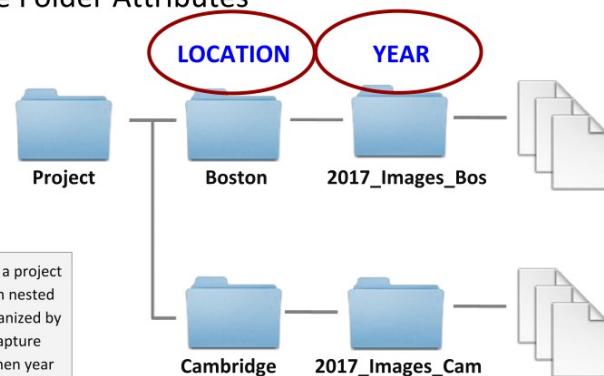


Create Folder Attributes



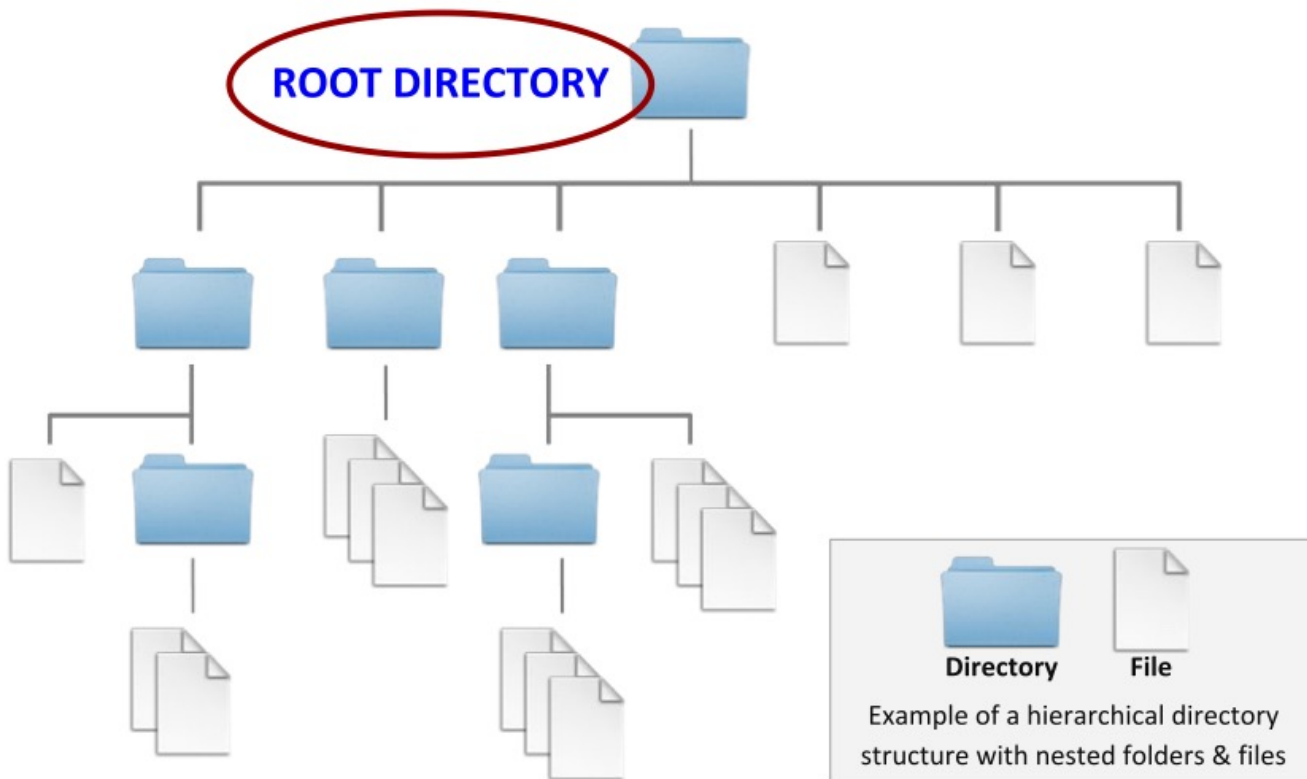
Example of a project folder with nested folders organized by year then image capture location

Create Folder Attributes

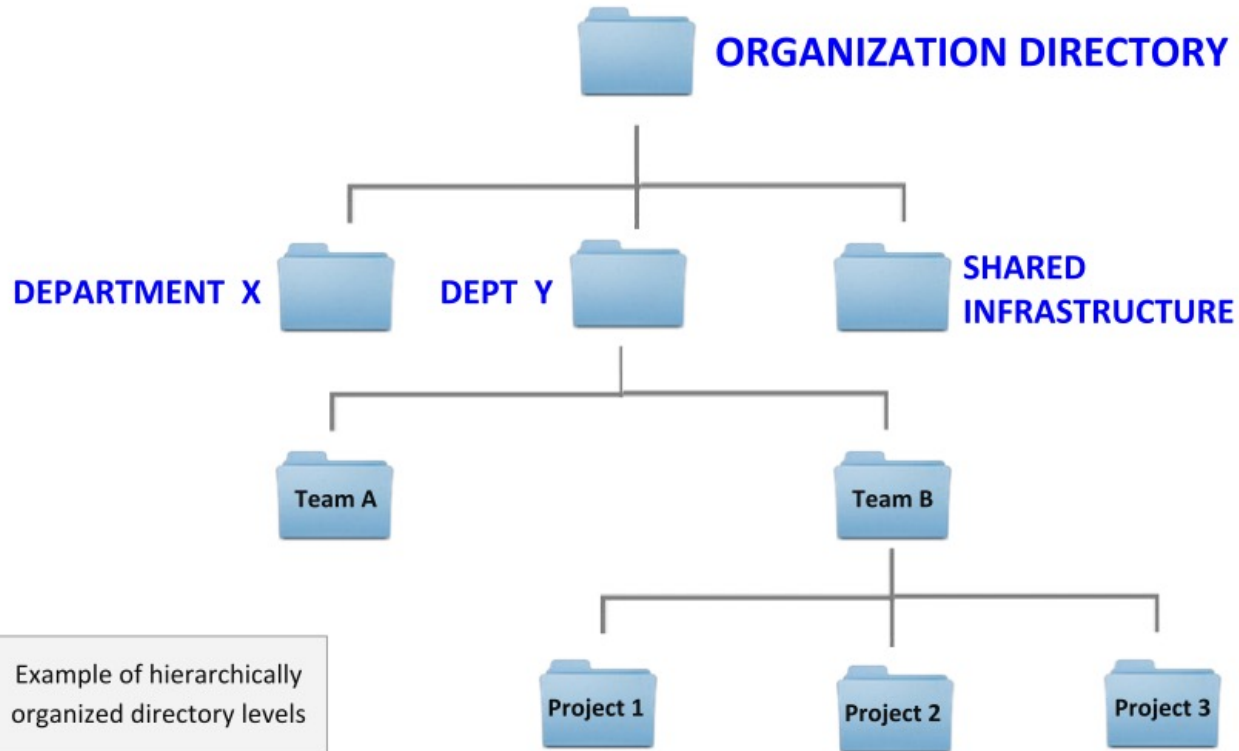


Example of a project folder with nested folders organized by image capture location then year

Top down

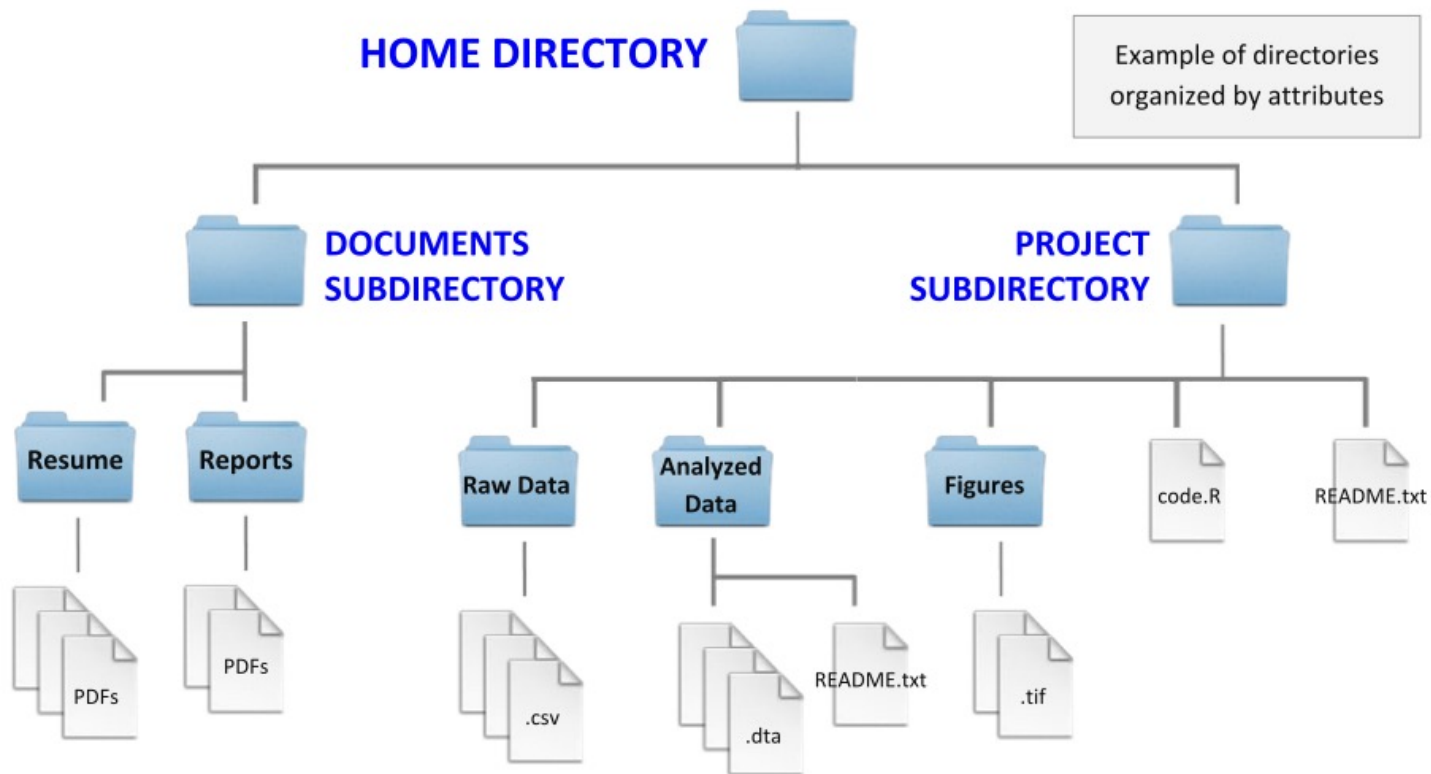


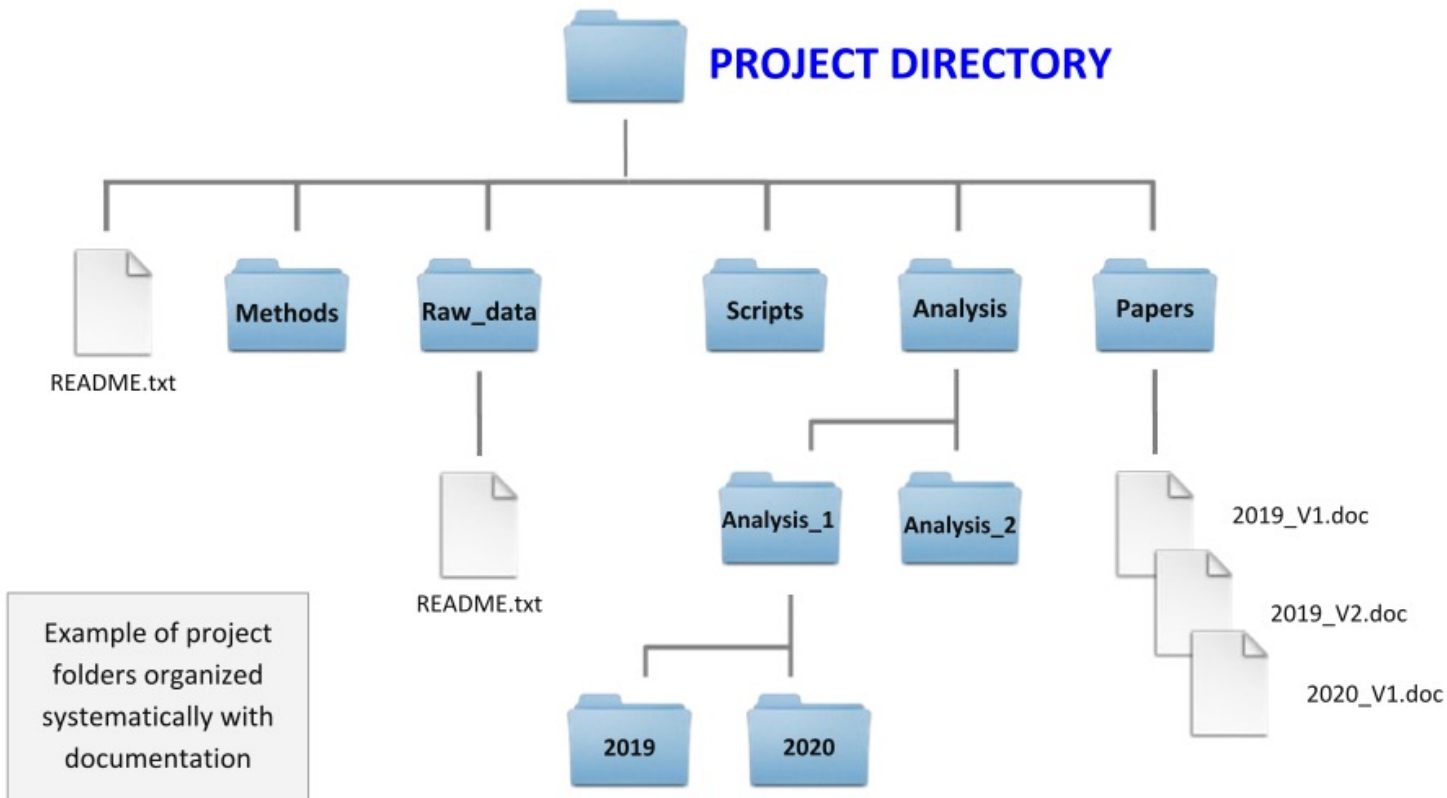
Bottom up



Example of hierarchically
organized directory levels

Hierarchy?





Exercise 3

Considering your own file structure, or the file structure used in your research group, write in the shared document:

1. If your current file structure is top-down or bottom-up, or a mix of the two.
2. Who the "inventor" of your file structure is (title, not name!), and to what extent you can influence how it is organized?
3. One thing you believe you could improve in your current file structure.



Digitalbevaring.dk

- Year or other date
- Type of data, document or file
- Project stages
- Analysis version or revision
- Experiments
- Instruments
- Time periods
- Geographic location
- Storage requirements
- Team member, institution or project site

Helpful characteristics!



Digitalbevaring.dk

- A file name is a principal identifier for the file
- Consistent in time and among different people
- Practically useful when accessing files, such as sorting and filtering

Chronologically

(ISO 8601 date standard)

```
20171028_001.tiff
20171028_002.tiff
20171028_003.tiff
20171029_001.tiff
20171029_002.tiff
```

Classification or code

(standardized)

```
USNM_379221_01.tiff
USNM_379221_02.tiff
USNM_379221_03.tiff
USNM_379222_01.tiff
USNM_379222_02.tiff
```

Alphabetically

(depending on type of files)

```
bos_20171028_001.tiff
bos_20171028_002.tiff
bos_20171029_001.tiff
cam_20170922_001.tiff
cam_20170922_002.tiff
```

- Human readable
- Machine readable

A well suited file naming protocol should be:

1. Human readable

A name describes the content of the file, connects to concept of a *slug* from semantic URLs (e.g. www.scilifelab.se/this-is-a-slug).









2. Machine readable





Avoid spaces, deliberate punctuation, accented or odd characters, inconsistent letter casing

3. Default ordered

Put something numeric first, use the ISO 8601 standard for dates (YYYYMMDD, or YYYY-MM-DD), left pad single digits with zeros (01, 02, 03... 10)

01_marshall-data.md	01.md
01_marshall-data.r	01.r
02_pre-dea-filtering.md	02.md
02_pre-dea-filtering.r	02.r
03_dea-with-limma-voom.md	03.md
03_dea-with-limma-voom.r	03.r
04_explore-dea-results.md	04.md
04_explore-dea-results.r	04.r
90_limma-model-term-name-fiasco.md	90.md
90_limma-model-term-name-fiasco.r	90.r
Makefile	Makefile
figure	figure
helper01_load-counts.r	helper01.r
helper02_load-exp-des.r	helper02.r
helper03_load-focus-statinf.r	helper03.r
helper04_extract-and-tidy.r	helper04.r
tmp.txt	tmp.txt

 2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_H01.csv
 2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_H02.csv
 2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_H03.csv
 2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_platefile.csv
 2014-02-26_BRAFWTNEGASSAY_FFPEDNA-CRC-1-41_A01.csv
 2014-02-26_BRAFWTNEGASSAY_FFPEDNA-CRC-1-41_A02.csv
 2014-02-26_BRAFWTNEGASSAY_FFPEDNA-CRC-1-41_A03.csv
 2014-02-26_BRAFWTNEGASSAY_FFPEDNA-CRC-1-41_A04.csv

 2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_H01.csv
 2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_H02.csv
 2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_H03.csv
 2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_platefile.csv

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

Explanation - Time_ProjectAbbreviation_ExperimentNumber_
Location_TypeOfData_VersionNumber

- Ensures orientation and longevity of data
 - Add to top level folder in separate plain text file
 - Collect all information to be documented and cross-reference where necessary
-
- Shared projects where files are handles by several individuals with different areas of responsibility
 - Project ending to guarantee remembrance of file structure over longer time (past your present You)

- Consider upgrading your file structure with an explicit **Conditions and use** agreement, and apply it to file permissions, at least for raw data.
- Edits and changes to raw data should be restricted to only a few trusted individuals with particular responsibilities.



Digitalbevaring.dk

Credit: Illustration from Digitalbevaring.dk / Jørgen Stamp (CC BY 2.5 Denmark license).

When we collect data and organise it for research purposes, findability is not necessarily our primary motivation.

- Increase the number of **file copies**

Pros - Someone (always?) have the latest version of a file!

Cons - Who has the latest version of a file?

- Increase the number of **shortcuts**

Pros - Easy to create, file name of shortcut can be changed and may even increase findability when named differently.

Cons - Shortcut may break if original file location or name is changed. Easy to lose orientation if not maintained regularly.

When we collect data and organise it for research purposes, findability is not necessarily our primary motivation.

- Increase the amount of **metadata**

Pros - Possible to enrich any file with unlimited amounts of metadata

Cons - Can be cumbersome to keep updated as number of files increase and file names are changed

Keyword tagging

(Metadata.txt content)

20220115_MyFile_Project1_Location_Dataiteration1_V1.xml

First version of X data from Y, with additions of Z made by A and B on 20220110 including suggestions by C.

Keywords HumptyDumpty Genome_Assembly

20220115_MyFile_Project1_Location_Dataiteration2.xml

Contains X data from Y, with additions of Z made only by A on 20220111 not including suggestions by C.

Keywords Published

Associated metadata to increase findability of files over e.g. multiple projects

For the following filename, construct a metadata explanation:

20220310_GenAnn_AssemblyProj2_GOT_Dataiteration_2nd_try.xml

Using the downloaded compressed file directory

Example_project_begin.zip :

1. Create a hierarchical folder system based on the file names and contents.
2. Rename files in a consistent manner if required, such that it reflects both contents and file version. Consider number and date formats as well.
3. Optionally, create a file for tagging files with metadata and keywords in accordance with file contents.

➤ How to store during the project?

- **Storage/processing locations**

For data collection, analysis, reporting, code, transfers etc.

- **Back-up and data recovery**

Strategies to mitigate risks of data-loss and data corruption?

(Beware of laptops and external storage)

- **Technical requirements**

Software and systems required to access / process the data?

➤ How to protect data?

- **Information classification**

Suitable storage based on the characteristics of the data?

- **Access control**

Who will have access to what data and how will it be enforced?

- **Data protection procedures**

Other strategies to mitigate risks of unwanted data disclosure or sabotage.

- Data has a life cycle

Raw (experiment) data – produce, collect, license, get access, ...

Processed – generate, clean, aggregate, label, transform, analyse, ...

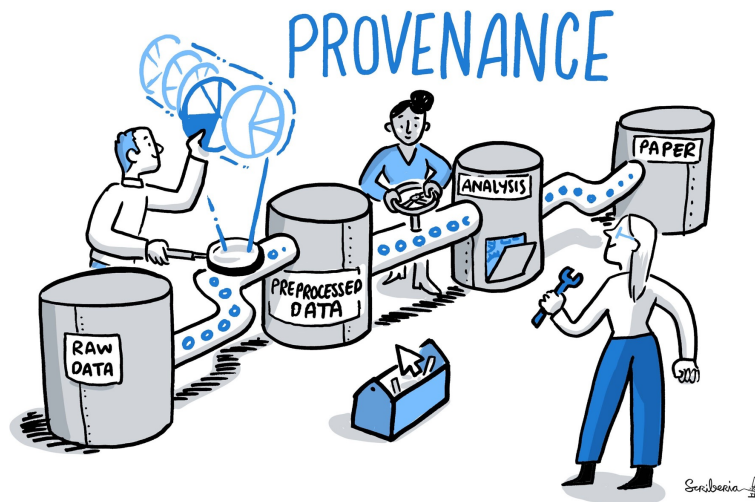
Long-term storage – document, select, convert, package, submit, ...

Published – FAIRify, promote reuse, ...

- Maintain data integrity and authenticity

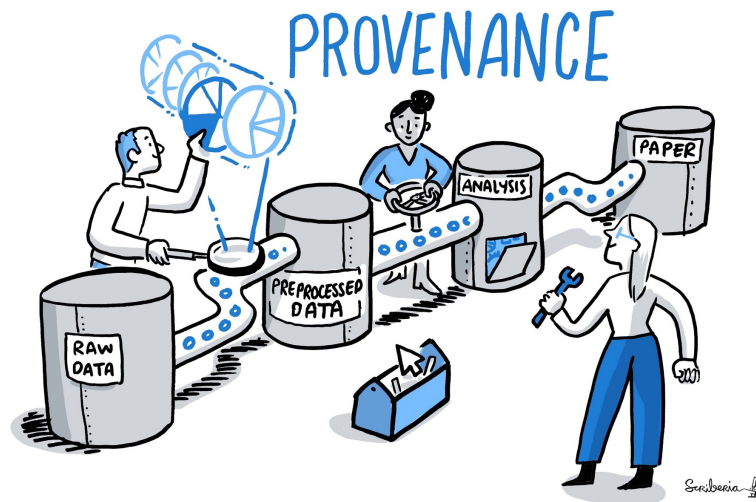
- Plan a storage strategy

- Plan a backup and disaster recovery strategy



What storage categories do you use and what factors do you consider when selecting which category to use or not to use?

- **Portable devices** – Laptops, tablets, external hard-drives, flash drives and Compact Discs
- **Cloud storage** - E.g. Google Drive, OneDrive, Dropbox, a University's OwnCloud, Open Science Framework and Tresorit
- **Local storage** – Desktop computers and personal laptops
- **Networked drives** – Shared drives on university servers, NAS servers (Network Attached Storage) or infrastructures (such as SNIC)



- **Temporary, short-term storage for non-sensitive data**, e.g. in the field or to transport data and files when online transmission is not possible.
 - In combination with **encryption and strong password protection**, especially if working with sensitive information.
- ✓ Conduct **regular checks** to ensure your device is working and that files are accessible.
 - ❖ **Not for long-term storage or master copies** of your data

- Granting **shared, remote and easy access to data and other files** to all involved in the project
- **Read the terms of service.**
Especially focus on rights to use content given to the service provider
- **Opt for European, national, or institutional** cloud services which store data in Europe if possible
- ❖ **Not your only storage and backup solution**
- ❖ **Not for unencrypted (sensitive) personal data**

- When working on different (local) workstations, e.g. laptop at home and the desktop in the office:
 - **always make sure that you are working on the most current version**, for example with the help of versioning software or guidelines
 - make sure that the most **current version is always backed up somewhere else**
- ✓ Suitable as a primary storage for projects involving only very few people
- ❖ Avoid if data will be moved back and forth between personal computers frequently

- **Use in projects involving many people** who need access to data and files
- Use a suitable security strategy to **protect data and files against unauthorised access**
- **Agree on rules for versioning files and data** to ensure that everyone can locate and access what they need
- **Long-term store complete data** that has been analysed, which can be cost efficient and offer increased security
- Restrict access where possible using rights and permissions, e.g. **write protect a master copy** and only grant access to specific files/folders when necessary

A minimalist strategy

- ❑ There are at least

- Three copies of the data, of which...

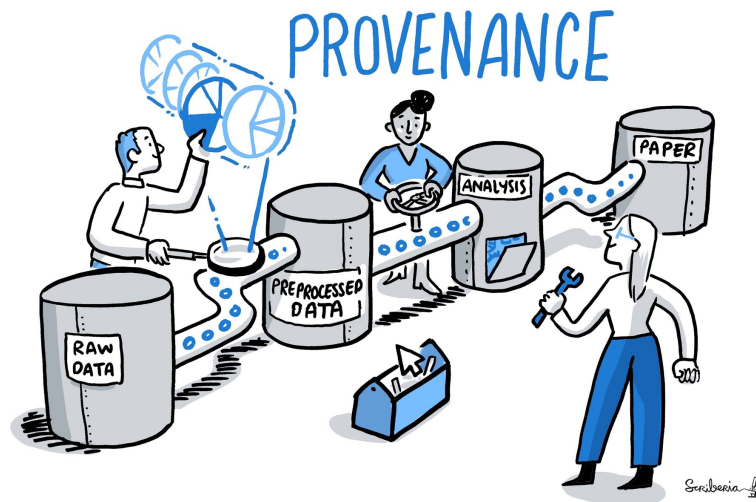
- ✓ Two are kept on different types of storage media
 - ✓ Two are at different locations
 - ✓ One is located off-site

- ❑ All copies are checked regularly to make sure that they work

- ❑ The process is known and applied in the project (automated)

... also determine what you want to back up, and find out whether your institution already has a backup strategy.

What are examples of potential causes for data loss in a research project?



Credit: This image was created by Scriberia for The Turing Way community and is used under a CC-BY licence.