# KU LEUVEN

# Organising & standardising research data that underpin your publication

Author: Veerle Van den Eynden, KU Leuven RDM Competence Centre

Trainer: Bruna Piereck

VIB RDM training, Gent

October 2024

# Overview

- *Organise and structure data and documentation files*

- *Logical, structured and descriptive file names*

- *Open / standard file formats*

- *File versioning in a project*

- *Data standards*
    - *make data interoperable and reusable*
    - *commonly understandable*

KU LEUVEN

# Organise / structure files

- Different options exist

- Here examples of good practices that researchers at KU Leuven use

- Find what works for you, in line with technical knowledge / support available and your data collection methods and active data storage system

- Centre organisation around:
  - Research project
  - Research paper

**KU LEUVEN**

# Folder structure

- File Explorer, OneDrive, MS Teams, ...

- Develop a structure organised by:
  - Paper, Project, Researcher, Experiment, Instrument

- Folders should:
  - follow a structure with folders and subfolders that correspond to the project design and workflow
  - have a self-explanatory name that is only as long as is necessary
  - have a unique name

- Good practice: ReadMe file in top folder

- Consider read / write access to folders for colleagues / collaborators

- When paper is published (or end of project): package structure and files into zip bundle and move to archival storage

An example:

```
project/
    code/               code needed to go from input files to final results
    data/               raw and primary data (never edit!)
        raw_external/
        raw_internal/
        meta/
    doc/                documentation of the study
    intermediate/       output files from intermediate analysis steps
    logs/               logs from the different analysis steps
    notebooks/          notebooks that document your day-to-day work
    results/            output from workflows and analyses
        figures/
        reports/
        tables/
    scratch/            temporary files that can safely be deleted or lost
    README.txt          file and folder description
```

Source: https://rdmkit.elixir-europe.org/data_organisation

*In Research Coordination Office at KULeuven, each project has a designated folder.*
*When a new project is started, a new folder is made. There raw data, syntax files,*
*questionnaires, ethical approval, etc are kept.*
*All researchers have access to the shared drive and to all folders.*

RDM Competence Centre

KU LEUVEN

# Record file

- Record file
  - A textual or tabular
  - List all data and documentation files of a project, paper, etc.
  - Specifies standard information for each dataset:
    - Unique ID
    - Dataset name
    - Description
    - Origin
    - Owner
    - Person responsible
    - Purpose, e.g. project name
    - Storage location, e.g. where on server, OneDrive, etc.
    - Contains personal data Y/N
    - Size / volume
    - Access: who has / needs access to the data

*In the research group they mostly develop algorithms for simulations. Every researcher has to keep a register (Word file) that lists which code repositories (on GitLab or GitHub) are used with the URL, and where data files are stored.*
*These registers are available with read access for colleagues.*

KU LEUVEN

# eLab Notebook

# Data management plan

**KU LEUVEN**

# File naming

- Develop a **logical** structure for **meaningful** file names

- Order 4-7 elements from generic to specific

- Suggested elements:
    - Project / experiment name, acronym or number
    - Creator name or initials
    - Date of creation: use ISO8601 format YYYYMMDD (and if needed time HHMMSS)
    - Type of data: sample ID
    - Version number: v01, v02, 00.01, 01.01 (leading zeros ensure correct sorting of files)
    - Location

- No spaces: use underscore (_), hyphen (- ) or Capitalized letters to separate elements

- Avoid special characters such as "/ \ : * ? " < > [ ] & $

- Independent of the location of the file on a computer

- Include a txt-file that explains your naming convention in your documentation

Dataset Challenges and Opportunities for Academic Parents during COVID-19

Eva Lantsoght

Anonimized dataset of the survey on the impact of COVID-19 on academic parents. Participants who did not give consent were filtered out as well.

Files (210.1 kB)

| Name | Size |
| --- | --- |
| dataset for Zenodo.xlsx | 210.1 kB |

md5:a2c2da0f619e5366f57314929ac7fa3f

Citations

Show only:  Literature (0)  Dataset (0)  Software (0)  Unknown (0)  Citations to this version

# File naming examples

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

- File name: 20201202_HB_EXP2_HEL_DATA_V03.xls
- Explanation: Date_ProjectAbbreviation_ExperimentNumber_Location_TypeOfData_VersionNumber

Source: https://rdmkit.elixir-europe.org/data_organisation

RDM Competence Centre

KU LEUVEN

# File naming examples

Cropped image of an ant head taken on the third of December 2020 by Meg Megson

- File name: 20201203_MM_HEAD_CROPPED_V1.psd
- Explanation: Date_CreatorData_Type_Modification_Version

Source: https://rdmkit.elixir-europe.org/data_organisation

RDM Competence Centre

KU LEUVEN

# File naming examples

Version 4 of the survey procedures for the British Dental Health Survey.

- BDHS_SurveyProcedures_00-04.pdf

- Explanation: Project acronym_Type_version number

Source: https://rdmkit.elixir-europe.org/data_organisation

RDM Competence Centre

**KU LEUVEN**

# Batch file renaming

- Need to rename large amounts of file names?
    - Images from digital cameras with automatically assigned files names/numbers
    - Default filenames generated by proprietary software or instruments
    - Removing spaces, odd characters, etc
    - Add meaningful elements to file name, e.g. project acronym, date, etc.

- Use a batch renaming tool for consistent, structured renaming
    - Bulk Rename Utility (Windows)
    - Renamer (Mac)
    - Gnome Commander (Linux)

- Demo: Renaming Files Using Bulk Rename Utility - YouTube

**KU LEUVEN**

# Batch rename example



Experiment measuring vertical dynamic running load with 13 treadmill users.

File renaming for 78 files:
- Find single digits in filename and add leading zero
- Add project name 'MALL' as prefix
- Add creation date as suffix, with underscore

# Exercise: folder structure & file naming

- The role of basal epithelial cells for small airway loss and epithelial injury in chronic lung disease.
    - Design a suitable folder structure for this research project
    - What would be useful elements for file names?

**KU LEUVEN**

# Exercise: folder structure & file naming

**Folders**

- Data
  - Scans
  - Processed data
  - Images
  - Measures
- Doc
  - SOP
  - …
  - Papers
- Code
- Results

**Elements**

- DONOR, COPD, IPF + number
- Whole, Part
- Mild, moderate, severe
- Lung
- Sample number

**KU LEUVEN**

# Exercise: open / standard file formats

- List all file formats you use in your daily work / research. (group activity)

- After a few minutes you will share with the group
  - Open formats
  - Proprietary formats (standard / not)

**KU LEUVEN**

# Open / standard file formats

- Use oppen/standard file formats

  - Long term access

  - Use of research data

- Good source: fairsharing.org

  - examples:

- Containers: TAR, ZIP

- Databases: XML, CSV, JSON

- Video: MPEG (mp4), AVI

- Sounds: WAVE, AIFF, MP3, FLAC

- Statistics: DTA, POR, SAS, SAV

- Images: TIFF, JPEG 2000, PNG, GIF

- Tabular data: CSV, TXT

- Text: XML, PDF/A, HTML, JSON, TXT, RTF

- 3D: X3D, C3D

- Neuroimaging: DICOM, Nifti

- Mass spectrometry: mzML

- Sequencing data: FASTA, FASTQ

- Microscopy: OME Next Generation File Format, Bio-formats conversions

KU LEUVEN

# File versioning

- Manage multiple versions

- Enable reverting to an earlier version

- Easy methods for small demands of versioning:
  - File naming
  - Cloud storage file versioning, e.g. OneDrive

- For automatic management of versioning
  - conflict resolution
  - back-tracing capabilities,
  - proper version control

- Git

- GitHub

- GitLab

- BitBucket

KU LEUVEN

# Data standards

- Make data interoperable
- Easier to understand
  - by multiple communities
- Reusable more widely

  - International, common standards

  - Community standards

**KU LEUVEN**

# Question

- Which standards do you already use in your research?

KU LEUVEN

# Community standard: biodiversity data

## Wolf observations Flanders 2022



Because biodiversity data are collected worldwide using the same data standards, collecting the same attributes and variables, they can be combined into large comparable datasets on the GBIF platform.

## GBIF platform wolf data 2022

KU LEUVEN

# GBIF & Darwin Core

## Appears in Datasets

APPEARS IN 69 CHECKLIST DATASETS:

**GBIF Backbone Taxonomy**
As *Canis lupus* Linnaeus, 1758

**Catalogue of Life Checklist**
As *Canis lupus* Linnaeus, 1758

**The European Nucleotide Archive (ENA) taxonomy**
As *Canis lupus*

**World Register of Marine Species**
As *Canis lupus* Linnaeus, 1758

**Integrated Taxonomic Information System (ITIS)**
As *Canis lupus* Linnaeus, 1758

**International Barcode of Life project (iBOL) Barcode Index Numbers (BINs)**
As *Canis lupus* Linnaeus, 1758

**Global Names Usage Bank**
As *Canis lupus* Linnaeus, 1758

**TAXREF**
As *Canis lupus* Linnaeus, 1758

**The Paleobiology Database**
As *Canis lupus* Linnaeus, 1758

APPEARS IN 545 OCCURRENCE DATASETS:

**Répartition historique du loup en France métropolitaine**
View occurrences

**NSW BioNet Atlas**
View occurrences

**Norwegian Biodiversity Information Centre - Other datasets**
View occurrences

**Swiss National Mammal Databank: Larger Carnivores Monitoring Program (KORA)**
View occurrences

**iNaturalist Research-grade Observations**
View occurrences

**Fauna Atlas N.T.**
View occurrences

**UAM Mammal Collection (Arctos)**
View occurrences

**SA Fauna (BDBSA)**
View occurrences

**Victorian Biodiversity Atlas**
View occurrences

## Darwin Core standard

| | | |
|---|---|---|
| Record-level Terms | Dublin Core terms, institutions, collections, nature of data record | Simple Darwin Core (flat) |
| Occurrence | evidence of species in nature, observers, behavior, associated media, references. | |
| Event | sampling protocols and methods, date, time, field notes | |
| Location | geography, locality descriptions, spatial data | |
| Identification | linkage between Taxon and Occurrence | |
| Taxon | scientific names, vernacular names, names usages, taxon concepts, and the relationships between them | |
| GeologicalContext | geologic time, chrono-stratigraphy, biostratigraphy, lithostratigraphy | |
| ResourceRelationship | explicit relationships between identified resources (e.g., one organism to another, taxon to location, etc.) | Generic Darwin Core (relational) |
| MeasurementOrFact | measurements, facts, characteristics, assertions, references | |

KU LEUVEN

# Standards

## International

- ISO 8601 standards for date / time
- ISO 3166 standard for country codes
- Getty Thesaurus for geographical names

## Community

- DICOM MRI data
- NACE code: Statistical classification of economic activities in European Community
- Standard International Age Classification, UNStat 1982

C. Learning and education services

1. Enrolment in regular and adult education — 2-4; 5 y.gr. 5-24; 10 y.gr. 25-64; 65+

2. Educational attainment — 5 y.gr. 15-24; 10 y.gr. 25-64; 65+

3. Illiteracy — 5 y.gr. 10-24; 10 y.gr. 25-64; 65+

G. Health, health services and nutrition

1. Morbitiy and handicaps (for mortality see I) — u 1; 1-4; 10 y.gr. 5-74; 75+

2. Usage of health services — u 1; 1-4; 10 y.gr. 5-74; 75+

3. Food consumption — u 1; 1-4; 10 y.gr. 5-74; 75+

4. Malnutrition — u 1; 1-4; 10 y.gr. 5-74; 75+

D. Earning activities and the inactive

1. Labour force participation — u 15; 5 y.gr. 15-24; 10 y.gr. 25-54; 5 y.gr. 55-74; 75-84; 85+

2. Employment/unemployment/ underemployment — u 15; 5 y.gr. 15-24; 10 y.gr. 25-54; 5 y.gr. 55-74; 75-84; 85+

*When age classification categories are applied consistently at an international level, datasets can be easily linked, combined and compared. But: different disciplines / purposes will need different categories !*

# Incompatible dates

## Terrorism attacks on buildings

| Date | Country | Target | Place |
|------|---------|--------|-------|
| 09/11/2001 | USA | WTC | New York |
| 13/11/2015 | France | Bataclan | Paris |
| 12/10/1984 | UK | Grand Hotel | Brighton |
| | | | |

## Tweets

| Tweet ID | Date-Time | Tweet text |
|----------|-----------|------------|
| 320217690004393984 | 2001-09-11T16:53:41Z | Lorem ipsum dolor sit amet |
| 320206007982755840 | 20001-09-11T16:07:16Z | Lorem ipsum dolor sit amet |
| 320205389780090880 | 2001-09-11T16:04:48Z | Lorem ipsum dolor sit amet |
| 320202492031930368 | 2001-09-11T15:53:17Z | Lorem ipsum dolor sit amet |
| 320197516371062784 | 2001-09-11T15:33:31Z | Lorem ipsum dolor sit amet |
| 320197511107211265 | 2015-11-13T15:33:30Z | Lorem ipsum dolor sit amet |
| 320195708835749889 | 2015-11-13T15:26:20Z | Lorem ipsum dolor sit ame |
| 320193833260425216 | 2015-11-13T15:18:53Z | tLorem ipsum dolor sit amet |
| | | |

KU LEUVEN

# Compatible dates: Linking 5 minute weather data with time of sunrise / sunset

TimeStamp in both datasets facilitates interoperability

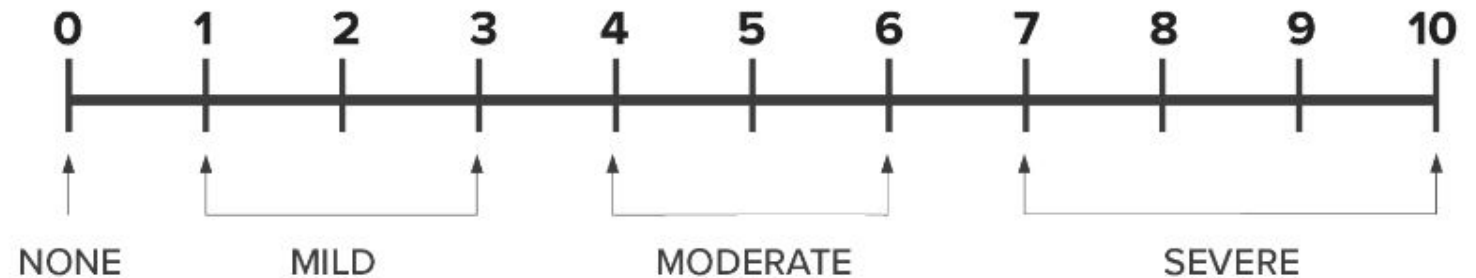# NIH Common Data Elements

A **Common Data Element** (CDE) is a standardized, precisely defined question, paired with a set of allowable responses, used systematically across different sites, studies, or clinical trials to ensure consistent data collection.



0 - 10 Numeric Pain Rating Scale

0  1  2  3  4  5  6  7  8  9  10

NONE    MILD    MODERATE    SEVERE

Categorical Scale

NO PAIN | HURTS A LITTLE | HURTS A LITTLE MORE | HURTS EVEN MORE | HURTS A WHOLE LOT | HURTS WORST

KU LEUVEN

# Quiz data standards

RDM Competence Centre

KU LEUVEN

# Lego replication game

KU LEUVEN

# Lego replication game: discussion

- Structured templates help to write out instructions
    - Standardises the process

- Brick lists help to write out instructions
    - Reduces ambiguity
    - Standardises naming
    - Brick lists could have unique numbers / codes for each brick
    - = controlled vocabulary / community standard

- Visuals help: drawing or pictures of vehicel

KU LEUVEN

# Standardisation …

| | |
|---|---|
| 6x2 brick | flag |
| 4x2 brick | 4x2 brick with slope |
| 3x2 brick | 2x2 brick with slope |
| 2x2 brick | 3x2 brick with slope |
| 4x1 brick | tall 2x1 brick with slope |

**DETAILED INSTRUCTIONS**

| Step | Parts required | Instructions |
|---|---|---|
| 1 | | |
| 2 | | |

**DETAILED INSTRUCTIONS**

| Step | Part shape | Part colour | Instructions |
|---|---|---|---|
| 1 | | | |
| 2 | | | |
| 3 | | | |

KU LEUVEN

# Take away messages

- When you start a project, design your folder structure and file naming system

- When you end your project / publish your paper, check your folder structure / file names are still in order (or fix), then zip and archive your data

- Use open / standard file formats when you can to make your data FAIR

- Use data standards where you can, to make your data interoperable and FAIR

RDM Competence Centre

KU LEUVEN