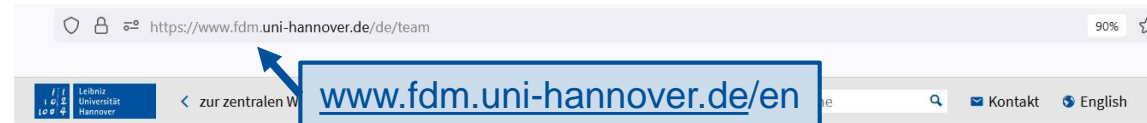# FAIR Software Management

## *How to handle research software*

**Yvana Glasenapp, LUH Research Service**

**Anna-Karina Renziehausen, TIB Publication Services**
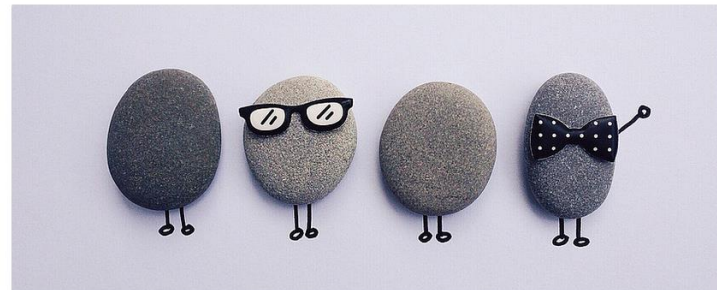
18.07.2024

# Welcome!

www.fdm.uni-hannover.de/en

Forschungsdatenmanagement

FDM an der LUH | Kurz erklärt | Förderanträge | Schulungen | Materialien | Tools
Forschungsdatenrepositorium | Team

Team

Das Service-Team Forschungsdaten



visit the website of the Research Data Support Team

**Anna-Karina Renziehausen**

TIB - Leibniz-Informations-zentrum Technik und Naturwissenschaften und Universitätsbibliothek

Publikationsdienste

**Yvana Glasenapp**

Gottfried Wilhelm Leibniz Universität Hannover

Dezernat Forschung und EU-Hochschulbüro, Technologietransfer

# Content of this course

- Introduction to Research Data Management

- How to make research software re-usable

- How to publish research software and code

- Supporting services and initiatives

# Introduction

- What are research data?

- What is research data management (RDM)?

- Why is data management important?

- The FAIR principles

- RDM is good scientific practice!

- Guidelines for handling research data at LUH

- LUIS services for data storage and backup
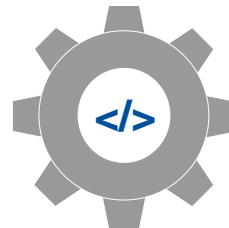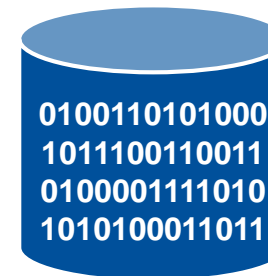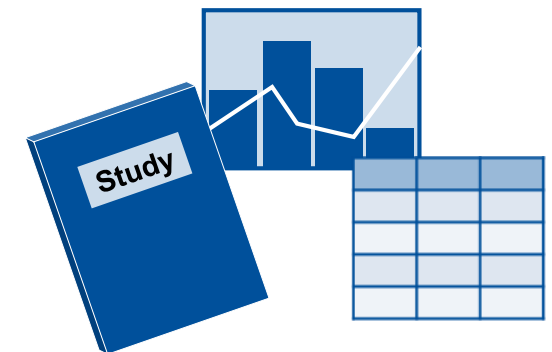
# What are research data?

Documentation

Raw data

Software code

Processed data

0100110101000
1011100110011
0100001111010
1010100011011

Results

Study

# What is research data management (RDM)?

planning

processing

documenting

administrating

selecting

structuring

archiving

publishing

```
010011010100010111
001100110100001111
010101010001101101
001101010001011100
110011010000111101
010101000110110100
```

i  RDM: Conscious and systematic handling of data from the planning stage until the end of a project and perhaps even beyond.

# Why is data management important?

✓ you keep an overview

✓ team work is easier

✓ you can safeguard high quality standards in research

✓ you save time and avoid stress

✓ you comply with official requirements

✓ Data management is the basis for Open Science

✓ You gain recognition for published data and software

**further reading**

The Thuringian Competence Network for Research Data Management compiled some "Research Data Scarytales". These are true stories about data management failures and their consequences.
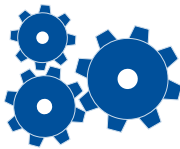
go to the „Scarytales"

# The FAIR Principles

**F**indable
- Unique and unambiguous identifier (e.g. DOI)
- Rich metadata (containing the DOI)
- Indexing in searchable directories (e.g. data repo)

**A**ccessible
- Accessible via a standardized, open, free, universal communication protocol
- Transparent access (authentication)
- Metadata accessible even if data is not

**I**nteroperable
- Standardized and widely applicable language for knowledge representation
- Vocabulary that follows FAIR principles
- References to related data

**R**e-usable
- Good documentation and precise attributes
- Clear license
- Detailed provenance
- Community standards

GO FAIR
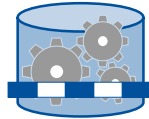Website

# RDM is good scientific practice!

RDM is especially relevant for these 6 out of 19 guidelines:

Guideline 7: Cross-phase quality assurance

Guideline 10: Legal and ethical frameworks, usage rights
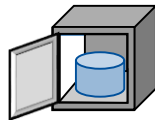
Guideline 11: Methods and standards

Guideline 12: Documentation

Guideline 13: Providing public access to research results

Guideline 17: Archiving

i You can read the guidelines themselves and accompanying disciplinary comments in the DFG portal "Research Integrity" (English version not finished, yet).

to the DFG portal

! Please note as well:

- DFG guidelines on the handling of research data
- disciplinary guidelines, requirements, policy papers etc.

You can find both on this website:

to the DFG website „Handling of Research Data"

# Guidelines for handling research data at LUH

The four principles:

- Data management:
  - protect against losses
  - process for sustainable (re-)use
  - document
  - archive (long-term)

- publish research data according to the FAIR principles in (disciplinary) repositories

- develop project-internal RDM policies and data management plans (recommendation)

- integrate RDM into teaching (recommendation)

→ These principles do not include specific advice on research software (yet). But they do apply for software, as it is listed as a type of research data in the definition

# LUIS services for data storage and backup

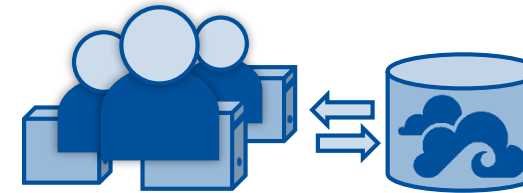personal "Cloud-Seafile"
(storage + data synchronisation)

go to service description (German)

"Projekt-Seafile"
(storage + data synchronisation)

go to service description (German)

LUH data archive

go to service description (German)

LUH data repository

go to service description (German)

Backup & Restore
(for servers of institutes)

go to service description (German)

# How to create re-usable research software

- Documentation of software projects

- The FAIR principles for research software

- Developing a software management plan (SMP)

- Services and infrastructure at LUH

- Foundation: How to set up a project

- Best practices of software projects

  - Tests and Test-driven development (TDD)

  - Refactoring

  - Clean code

  - Error messages

  - README

# Documentation of software projects

Data    is not independent from    Research software



But most researchers    ≠    Are no software developers

"Data without software are just numbers"
(Davenport & Grant, 2020)

ALL MODERN DIGITAL INFRASTRUCTURE

A PROJECT SOME RANDOM PERSON IN NEBRASKA HAS BEEN THANKLESSLY MAINTAINING SINCE 2003

https://xkcd.com/2347

# The FAIR Principles for research software

**F: Findable**

This means that the software and all important related information can be easily found by both humans and machines.

**F1.** The software, individual important components of the software and the different versions of the software all receive persistent and unique identifiers (PID), e.g. DOI.

**F2.** The software is described in detail with metadata, which means that all important information for use and subsequent use is included.

**F3.** Metadata clearly and explicitly include the PID of the software they describe.

**F4.** Metadata is also FAIR, searchable and can be found and used by metasearch engines.

✓ Use persistent identifiers and rich metadata

Findable

# The FAIR Principles for research software

**A: Accessible**

Which means that the software and its metadata is retrievable via standardized protocols

**A1.** For this, the protocol should fulfil some requirements, it should be open and free as well as universally implementable. It should also enable an authentication and authorisation procedure, if necessary.

**A2.** In addition the metadata are accessible, even when the Software is no longer available.

✓ Make software and metadata retrievable

Accessible

## The FAIR Principles for research software

**I: Interoperable**

For software, being interoperable means the possibility to exchange data between independent software. (It differs from data in the way that independent software can not be combined like data to form a new data set.)

**I1**. When software interacts to exchange data, the exchange protocol need to be clearly described. Domain-relevant standards should be used and APIs should be documented in human and machine readable form.

**I2**. Software should include references to external data, software or objects, which are required to execute the software. Example: software X is implemented using software A (a programming language)

✓ Use APIs, standards and references

Interoperable

# The FAIR Principles for research software

**R: Reusable**

This means that you can not only use the software as it is, but have further information to understand, modify and built new versions upon this software.

**R1.** The software is described with relevant attributes.

**R1.1.** A license is given and should be as unrestrictive as possible.

**R1.2.** Software is associated with detailed provenance: This is a metadata which includes the history of the software, how it came to be and who contributed to it.

**R2.** Include references to other software which is necessary to compile and run the software.

**R3.** Your software should meet domain-relevant community standards and coding practices. This can be the choice of programming language, testing methods or file formats.

✓ Document, licence, and follow community standards

Reusable

# Developing a software management plan (SMP)

Can help to:

- Create a basis for establishing best practices

- Make the research software reusable and sustainable

- Plan for the necessary resources (financial, human, infrastructure)

- Make it easy to introduce new developers into the project

Ideally an SMP should be drafted at the beginning of a research project. The extent of information in the SMP varies with level of management intensity of the software (low, medium, high).

**Software Management Plan (SMP)**

A SMP does not need to be long. The following guide offers three templates adapted for different software management levels:

> Practical guide to software management plans

The Max Planck Society has published a comprehensive SMP template, which can be integrated in RDMO:

> MPG-Template "Software management plan for researcher"

# Developing a software management plan (SMP) part I

- Purpose

- Version Control

- Repository

- User Documentation

- Software licensing and compatibility

- Deployment documentation



**Engineering Focus**

Version control · Packaging · Testing · Software Engineering quality

Purpose

Project management focus

Documentation: User documentation · Deployment documentation · Developer documentation

Software licensing and compatibility · Maintenance · Repository · Citation · Support/ Resources (during the project) · Risk analysis

**Figure 1.** Software Management Plan requirements grouped by their focus.

# Developing a software management plan (SMP) part II

- Citation

- Developer documentation

- Testing

- Software Engineering quality

- Packaging

- Maintenance

- Support

- Risk analysis

### 6.1.4. Summary of SMP templates developed for three management levels

| Core requirement (Section 5.1) | Software management level (Section 6.1) | | |
|---|---|---|---|
| | Management level: Low (6.1.1) | Management level: Medium (6.1.2) | Management level: High (6.1.3) |
| Purpose | × | × | × |
| Version control | × | × | × |
| Repository | | × | × |
| User documentation | | × | × |
| Software licencing and compatibility | | × | × |
| Deployment documentation | | × | × |
| Citation | | × | × |
| Developer documentation | | × | × |
| Testing | | × | × |
| Software Engineering quality | | × | × |
| Packaging | | × | × |
| Maintenance | | × | × |
| Support | | | × |
| Risk analysis | | | × |

**Table 4.** Core requirements of an SMP for software grouped by management level.

# Services and infrastructure at LUH

**For you**

- <u>Version control</u> and software development: GitLab

- <u>Courses</u> for employees and students:

  - Statistics with R

  - Python: Introductory course, Object-oriented programming, Numerical calculation
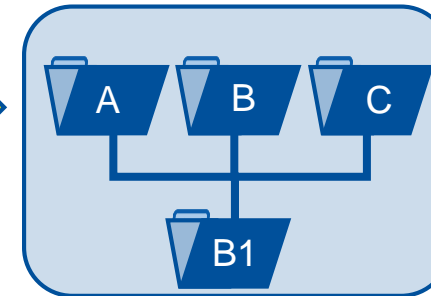
**For your institution**

- <u>Research cluster housing service:</u> The hardware of an institution can be integrated into the existing cluster system of LUIS

- <u>Scientific Computing:</u> For computationally intensive research project, LUIS provides a cluster system with massively parallel computer systems

# Foundation: How to set up a project

**Create a similar folder layout:**

- Licence
- Readme
- Experiments (chronological)
- Provenance
- Scripts

A  B  C

B1

**Logical order in subfolders:**

- Main script
- Data-generating scripts
- Raw data
- Post-processed data

**Modular naming scheme:**

Date JJMMTT ── … ── Parameter ── Version

**Example:** 230701_PuppyCuteness_Treats_Cheese_Camembert_v03.fff

**Data format:** Better not choose "Fancy File Format"

**Find a file format**



https://xkcd.com/1459

Untitled 138.docx
Untitled 241.doc
Untitled 138 copy.docx
Untitled 138 copy 2.docx
Untitled 139.docx
Untitled 40 MOM ADDRESS.jpg
Untitled 242.doc
Untitled 243.doc
Untitled 243 IMPORTANT.doc

OH MY GOD.

PROTIP: NEVER LOOK IN SOMEONE ELSE'S DOCUMENTS FOLDER.

## Best practices of software projects



How to write good code (xkcd flowchart)

There are a number of published „Best practice" guides for software development in different scientific fields. We have collected some general practices that help you to write FAIR software.
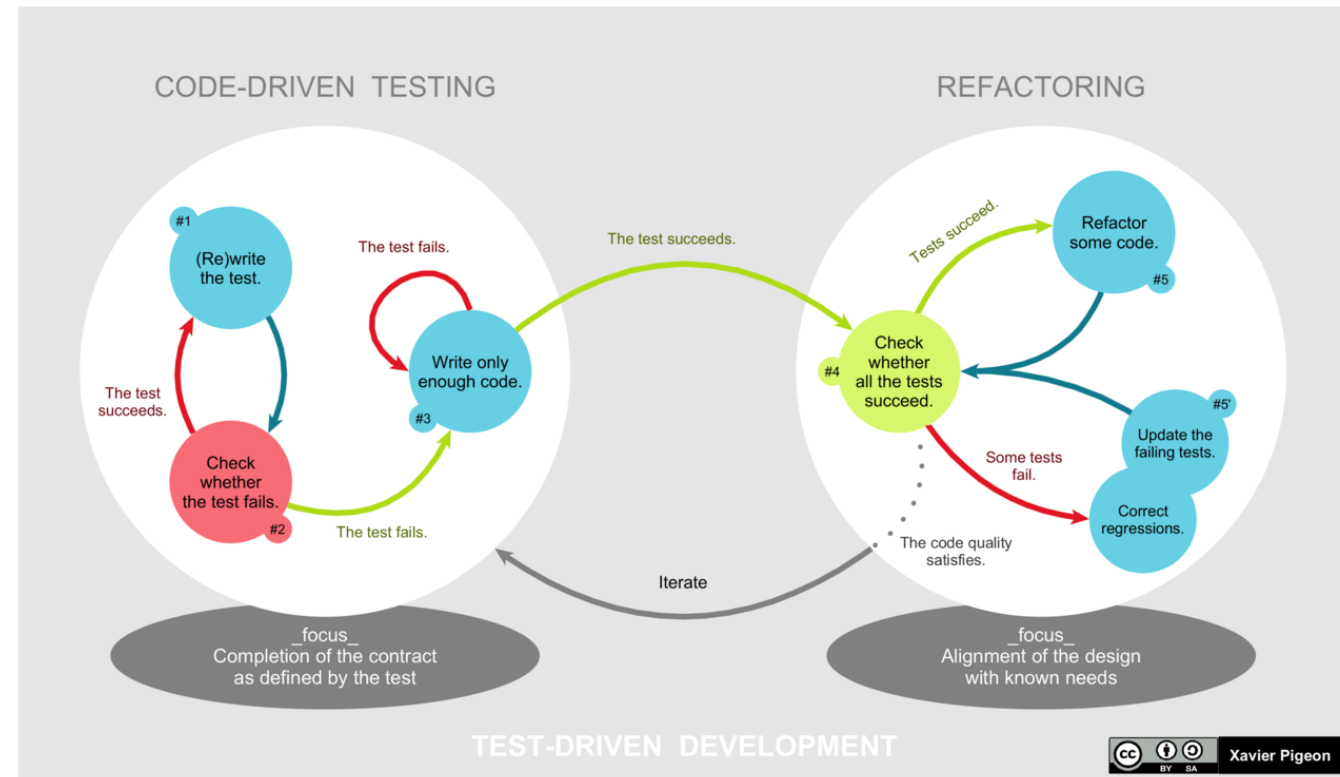
Ten simple rules for documenting scientific software

Oliver Melchert's best practices for small computational projects



Suresoft — SUSTAINABLE RESEARCH SOFTWARE

# Best practice: Tests and Test-driven development (TDD)

➢ Include testing in the development process

➢ Each new piece of code is tested for functionality before it is added to the main branch

➢ Best practice: add incremental changes!

➢ TDD: First create a test for the desired functionality, then write the code that is able to fulfill the test requirements

# Best practice: Refactoring

Code refactoring is the process of restructuring code without changing its original functionality.
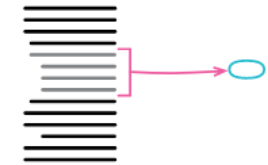
This helps to ensure that

- the code is more clearly structured and therefore easier to read and understand

- errors and weak points are easier to find and rectify

- the code is easier to extend and better to maintain

> Martin Fowler, Refactoring: improving the design of existing code (German)

> Robert C. Martin, Clean Code: Refactoring, Patterns, Testen und Techniken für sauberen Code

**Extract Function**



```
function printOwing(invoice) {
  printBanner();
  let outstanding  = calculateOutstanding();

  //print details
  console.log(`name: ${invoice.customer}`);
  console.log(`amount: ${outstanding}`);
}
```

```
function printOwing(invoice) {
  printBanner();
  let outstanding  = calculateOutstanding();
  printDetails(outstanding);

  function printDetails(outstanding) {
    console.log(`name: ${invoice.customer}`);
    console.log(`amount: ${outstanding}`);
  }
}
```

Martin Fowler, Refactoring catalog: https://refactoring.com/catalog/extractFunction.html

# Best practice: Clean code

➢ There are different options how to write names, comments, functions and classes in your code

➢ Even though a "messy" written code might work, it is not re-usable

➢ Changing small habits can make a big difference!

**In this example, it is better to define a function which contains the for-loop:**

| Do not | Do |
|---|---|

```python
input_number = int(input("Enter a number"))

# check if input_number is a prime number
is_prime = True
for i in range(2, input_number):
    if input_number % i == 0:
        is_prime = False
        break

print(input_number, "is a prime:", is_prime)
```

```python
def is_prime(number):
    prime_flag = True

    for i in range(2, number):
        if number % i == 0:
            prime_flag = False
            break

    return prime_flag


input_number = int(input("Enter a number"))
print(input_number, "is a prime:", is_prime(input_number))
```
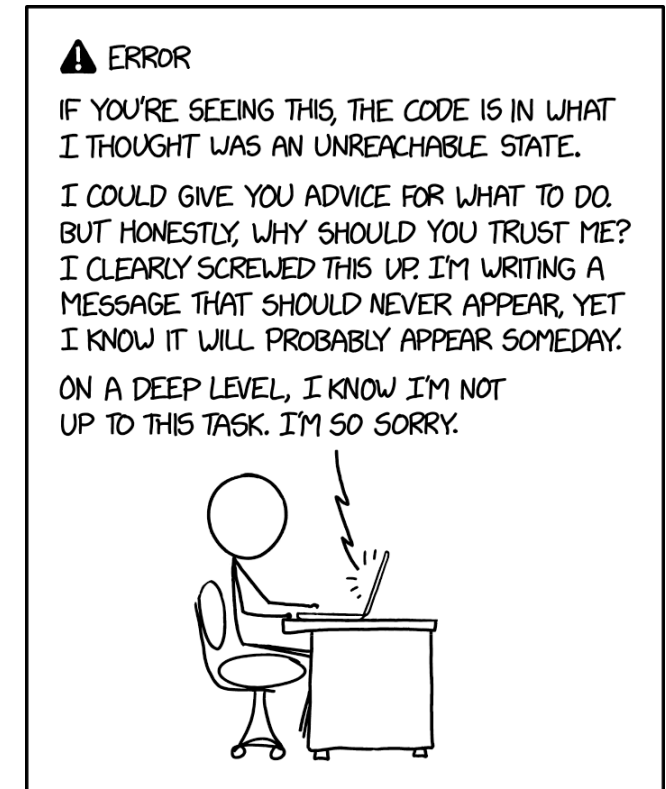
https://suresoft.dev

# Best Practice: Error messages

➢ Write error messages

➢ Provide solutions or where to find the information relevant to fixing the error.



https://xkcd.com/2200

Graphic: Unknown, via: https://community.spiceworks.com/topic/2105369-10-hilarious-error-messages-facepalm-worthy-computer-prompts-that-make-no-sense

# Best practice: README Part I

- ➢ Give an overview of the project

- ➢ List the names and affiliation of people involved

- ➢ Environment and dependencies needed to run the code

### All-optical Supercontinuum Switching (code and data)

Code repository for the article

"All-optical Supercontinuum Switching"

Oliver Melchert (1,2,3), Carsten Brée (4), Ayhan Tajalli (2), Alexander Pape (2,5), Rostislav Arkhipov (6), Stephanie Willms (1,2), Ihar Babushkin (1,2), Dmitry Skryabin (7), Günter Steinmeyer (8,9), Uwe Morgner (1,2,3), and Ayhan Demircan (1,2,3)

1. Cluster of Excellence PhoenixD, Welfengarten 1, 30167, Hannover, Germany
2. Institute of Quantum Optics, Leibniz University Hannover, Welfengarten 1 30167, Hannover, Germany
3. Hannover Centre for Optical Technologies, Nienburgerstr. 17, 30167, Hannover, Germany
4. Weierstraß Institute for Applied Analysis and Stochastics, Mohrenstraße 39, 10117 Berlin, Germany
5. VALO Innovations GmbH, Hollerithallee 17, 30419, Hannover, Germany
6. St. Petersburg State University, Universitetskaya nab. 7/9, St. Petersburg 199034, Russia
7. Department of Physics, University of Bath, Bath, BA2 7AY, UK
8. Max-Born-Institute (MBI), Max-Born-Str. 2a, 12489 Berlin
9. Institut für Physik, Humboldt-Universität zu Berlin, Newtonstraße 15, 12489 Berlin, Germany

This repository contains code and data analysis scripts for reproduction of simulated data and a draft versions of Figures 3 of the article.

The provided code implements the nonlinear propagation of optical pulses in a NLPM750 photonic crystal fiber in terms of the generalized nonlinear Schrödinger equation, including the effects of dispersion, self-phase modulation, self-steepening, Raman effect, and quantum noise.

### Environment and dependencies

The provided code is written using Python (2.7.15+) and requires the functionality of

- numpy (>=1.8.0rc1)
- scipy (>=0.13.0b1)
- matplotlib (>=1.2.1)

O. Melchert, Research data management in practice

## Best practice: README Part II

➢ Details about included material

➢ License, Citation & Acknowledgements

Subfolder `/scripts` contains:

- `pyGNLSE.py` : driver script implementing the genaralized nonlinear Schrödinger equation.

### License

This project is licensed under the MIT License - see the LICENSE file for details.

### Acknowledgments

This work received funding from the Deutsche Forschungsgemeinschaft (DFG) under Germany's Excellence Strategy within the Cluster of Excellence PhoenixD (Photonics, Optics, and Engineering – Innovation Across Disciplines) (EXC 2122, projectID 390833453).

**Included materials**

The repository follows a modular structure:

```
numExp01_FIG03_propagationDynamics
├── data
├── FIG03_subfigures
│   ├── FIGS
│   │   ├── GNLSE_BlowWood_tMax8000.000000_Nt32768_zMax600000.000000_Nz15
│   │   ├── GNLSE_BlowWood_tMax8000.000000_Nt32768_zMax600000.000000_Nz15
│   │   └── GNLSE_BlowWood_tMax8000.000000_Nt32768_zMax600000.000000_Nz15
│   ├── getFigures.sh
│   ├── main_FIG03_subfigure.py
│   └── spectrogram.py
└── main_NLPM750_propagationDynamics_qNoise.py

provenance_data
├── dataProvenance.sh
├── get_info.py
└── provenance_numExp01.log

scripts
└── pyGNLSE.py

src
├── data_handler.py
├── raman_response.py
├── solver.py
└── utils.py

LICENSE
Readme.md
```

O. Melchert, Research data management in practice

28

# How to publish research software and code

- Why publish your data and software?

- Preparation of your software publication

  - Version control your documentation

  - Use automated documentation tools

  - Software citation

  - User guide and examples

- Publishing software

- Licences

- Linking journal articles with related data and software
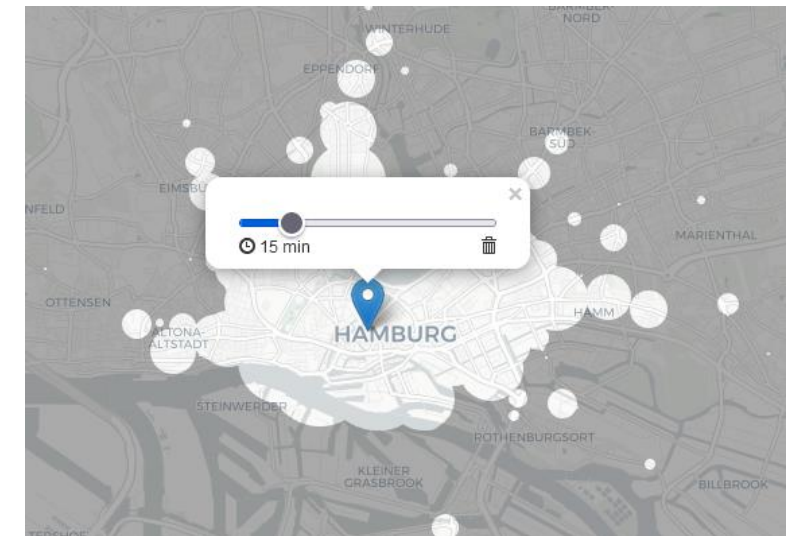
- Example: Publish software via Zenodo

# Why publish your data and software?

Open data are valuable and can be used to create new research output or tools of public interest:

- Participation: findtoilet.dk

- Empowerment of citizens: mapnificient

- Innovation: EHEC outbreak genome analysis

Reasons to publish your research software:

- To enable software citation

- To make the software FAIR

- To let software count towards evaluation

Mapnificent.net

# Software publication: Version control your documentation

➢ Make it clear which documentation applies to which software version

➢ A changelog can make this task much easier

➢ There are useful services that can help you do this.

## Changelog

All notable changes to this project will be documented in this file.

The format is based on Keep a Changelog, and this project adheres to Semantic Versioning.

### Unreleased

#### Added

- v1.1 Italian translation.
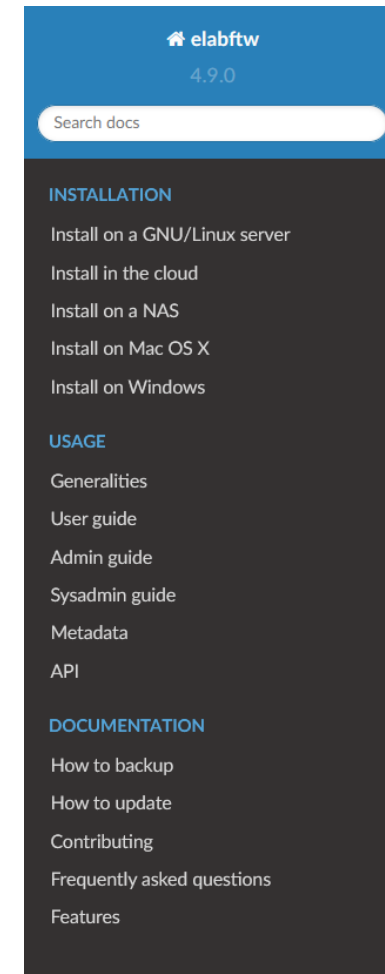- v1.1 Polish translation.

### 1.1.1 - 2023-03-05

#### Added

- Arabic translation (#444).
- v1.1 French translation.
- v1.1 Dutch translation (#371).
- v1.1 Russian translation (#410).
- v1.1 Japanese translation (#363).
- v1.1 Norwegian Bokmål translation (#383).
- v1.1 "Inconsistent Changes" Turkish translation (#347).
- Default to most recent versions available for each languages

Screenshot: https://github.com/olivierlacan/keep-a-changelog/blob/main/CHANGELOG.md

# Best Practice: Use documentation tools

➢ Documentation of a project is the alpha and omega

➢ Complete and prompt documentation is essential

➢ Use tools to create a nice documentation site which helps others to use your software

Sphinx Python Documentation Generator:

```
docs
├── build
├── make.bat
├── Makefile
└── source
    ├── conf.py
    ├── index.rst
    ├── _static
    └── _templates
```

# Software publication: Software Citation

➤ Provide information on how to cite your code

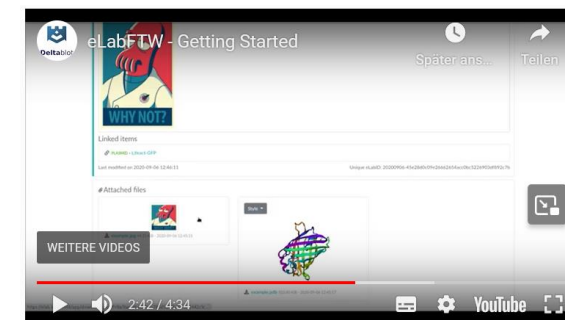➤ Be sure to always properly cite other people's code that you reuse as well

```
# This CITATION.cff file was generated with cffinit.
# Visit https://bit.ly/cffinit to generate yours today!

cff-version: 1.2.0
title: Cheese_TAX_Alert
message: >-
  If you use this software, please cite it using the
  metadata from this file.
type: software
authors:
  - given-names: Sheryl
    family-names: Mc Sniff
    email: SMS@gmail.com
    affiliation: University of Dogford
repository-code: 'https://git.eu/McSniff/cheesetaxalert'
repository: 'https://doi.org/10.123456789/mcsniff'
abstract: >-
  the software can be installed in a refrigerator so that as
  soon as the refrigerator door is opened, the Puppy
  receives a warning signal that the Cheese Tax must be
  collected.
keywords:
  - puppy
  - cheese
  - tax
license: MIT
version: '1.0'
date-released: '2023-09-01'
```
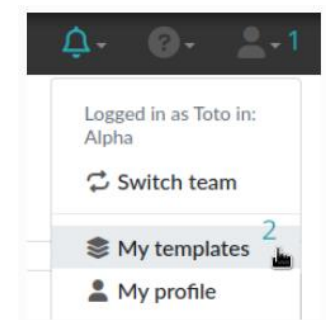
# Software publication: User guide and examples

➢ People might not use your software if it takes too much time to get started

➢ Add examples, a tutorial, videos – anything that helps to show the functionalities of your software

➢ If you have numerous examples, use a special section or directory

➢ Try your user guide on someone who does not know your software yet



Video tutorials

Screenshots

## Publishing software

- Publish software in a publication repository (like Zenodo or your institutional data repository)

  - This gives you:

    - An unique and persistent identifier ( e.g. DOI)

    - An identifier that points to the specific version of your software

- Provide citation metadata for your software (via CFF)

- License your code

- Possibility of additional publication in a software journal

Repository publication → software journal → papers

**Important: The source code repository (e.g. GitHub, GitLab) is not a publication repository!**

# Licenses

**Common software licenses:**
- MIT
- MPL 2.0
- GPL 3.0 – only

Choose an open source license

**Licenses for text and data:**

**Creative Commons licenses**
- CC0 – no name attribution, no restrictions
- CC-BY – author attribution when using the data
- CC-BY-SA – author attribution and share under the same conditions

Choose a CC license

**Example: Overview of the MIT-License**

## MIT License

A short and simple permissive license with conditions only requiring preservation of copyright and license notices. Licensed works, modifications, and larger works may be distributed under different terms and without source code.

**Permissions**
- Commercial use
- Distribution
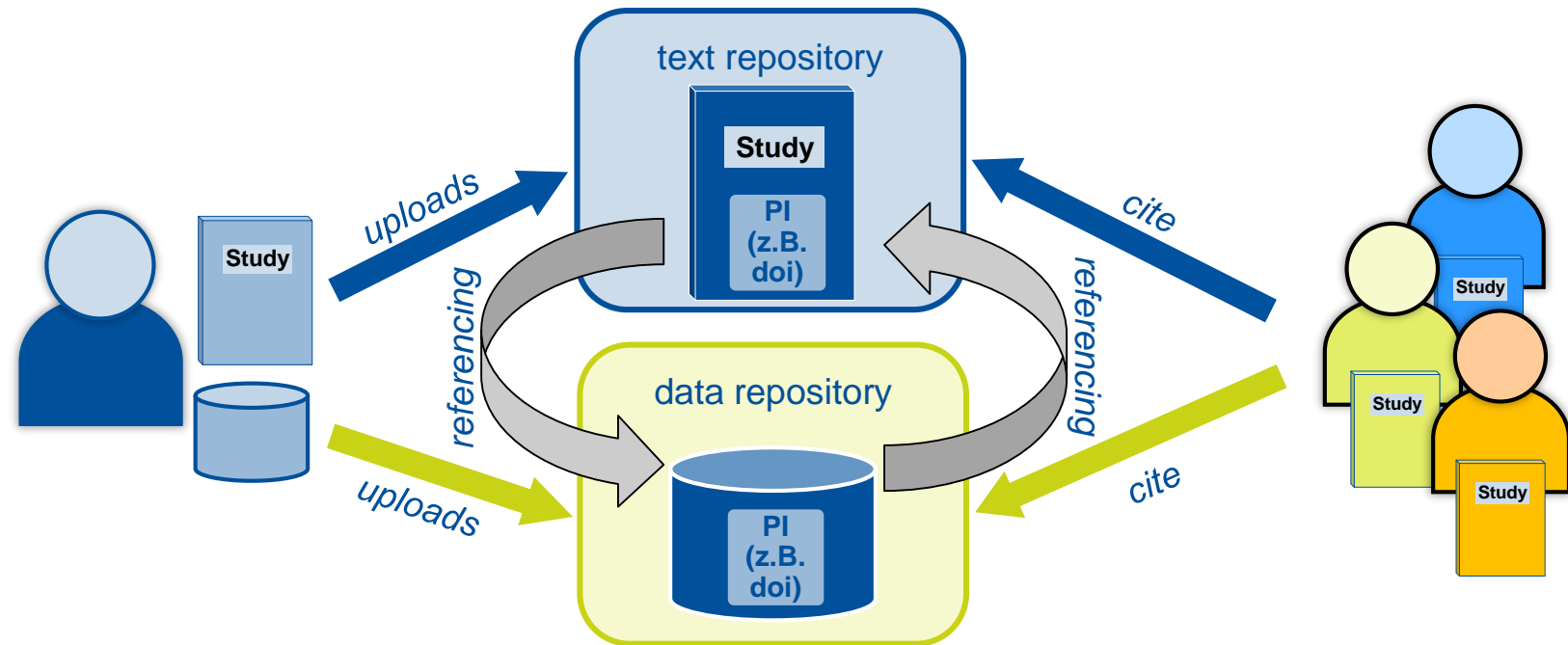- Modification
- Private use

**Conditions**
- License and copyright notice

**Limitations**
- Liability
- Warranty

https://choosealicense.com/licenses/mit/

# Linking journal articles with related data and software



**further reading**

Diana Kwon (2019): The Push to Replace Journal Supplements with Repositories. Online article, 19 August 2019, on the website of „The Scientist".

go to article

i

LUH also provides its own institutional text repository for open access publications.

go to text repository

37

# Summary

- Software management happens before you even start coding:

    - Clarify responsibilities

    - Agree on standards and methods and write them down (SMP)

    - Recognize that data and software management are work tasks that will consume some of your time

- Develop a culture that values data- and software management

    - Include the best practices in your everyday work

    - Think of your future self and anyone who might want to re-use your code

    - Engage in projects of others: Create and solve issues, reuse instead of setting up new
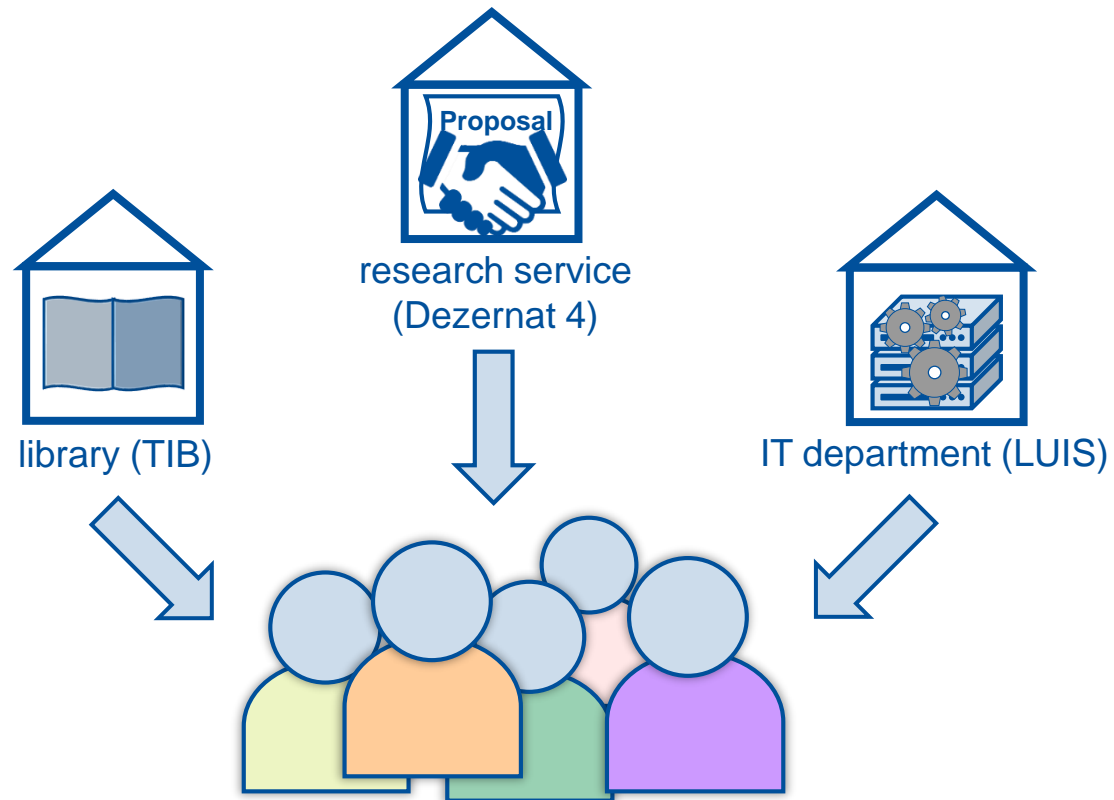


https://xkcd.com/1421/

# Supporting services and initiatives

- The Research Data Support Team of LUH

- External information and support

- RSE working groups and initiatives

# The Research Data Support Team of LUH

**i**

We train and counsel LUH members on the following topics, among others:

- documentation and publication of data
- DMP and RDM policies
- RDM statements in funding proposals
- LUH services and infrastructures for data management
- legal issues (in cooperation with the legal department and the data protection office)
- practical implementation of data management in research processes

Please note as well the additional information, guides and training courses on our website.

**Proposal**

research service
(Dezernat 4)

library (TIB)

IT department (LUIS)

Sign up for our monthly newsletter on RDM!

go to website of the
Research Data Support Team

# External information and support

| Courses, talks and schools: | Self-training courses: |
|---|---|

**Suresoft Project:** Workshop schedule with courses on best practices in research software engineering

**Suresoft Knowledge Hub:** Workshops on Git foundations, Clean Code, Software Testing, CI and more

**Hessian Research Data Infrastructures:** Data and Code Schools, Data Talks

**Software Carpentry:** You can work through the lessons on the Unix Shell, Git, Programming with Python or R by yourself

**Helmholtz Information & Data Science Academy:**
Courses, schools, workshops & talks

# RSE working groups and initiatives



**Research Software Engineers (RSEs)**
- German association (de-RSE e.V.)
- monthly call, Chat, Blog
- International conferences for and by Research Software Engineers



**Software Sustainability Institute (UK)**
- Motto: "Better software, better research"
- Research Software Healthcheck: Evaluation tool to check your software
- online Research Software Camps



**Research Data Alliance (RDA):**
International network of RDM specialists with a special focus on technical aspects.



**NFDI:** The federal government and the state governments jointly fund the development of a national research data infrastructure (NFDI). X consortia are developing disciplinary standards, establish services and offer trainings.

# Thank you for taking this course!



by Auke Herrema