

# Best Practices in Research Data Management

*Hauke Sonnenberg & Michael Rustler*

*2018-06-08*



# Contents

<b>Introduction</b>	<b>5</b>
<b>I Best Practices</b>	<b>7</b>
<b>1 Data-Related Project Controlling</b>	<b>9</b>
<b>2 Naming Conventions</b>	<b>11</b>
2.1 Acronyms . . . . .	11
<b>3 Folder Structures</b>	<b>13</b>
3.1 Project Folder Structure . . . . .	13
3.2 Rawdata Folder Structure . . . . .	13
<b>4 Raw Data</b>	<b>15</b>
<b>5 Metadata</b>	<b>17</b>
5.1 Special Files . . . . .	17
<b>6 Data Processing</b>	<b>19</b>
6.1 Automatic Data Import into R . . . . .	19



# Introduction

This document describes best practices in research data management.

It is the outcome of the FAKIN project.

It is mainly based on the experiences of the authors.



Part I

Best Practices





# Chapter 1

## Data-Related Project Controlling

At the start of a research project

- Choose a project acronym and store it in `PROJECTS.txt`.
- Check if the organisations that you expect to get data from are listed in `ORGANISATIONS.txt` and extend this file if necessary.
- Create a subfolder for your project and subfolders for the organisations in the rawdata folder structure.

At the start of a project or if an employee or trainee enters the project

- Give an introduction to our research data management as described in this document.

Regularly during the project

- Check if the folder structure within your project's rawdata subfolder still complies with the rawdata folder structure and clean the structure, if not.



## Chapter 2

# Naming Conventions

### 2.1 Acronyms

Acronyms are unique, clear names for objects. They should

- be short but meaningful and easy to remember,
- be all lowercase,
- consist of only alphanumeric letters (**a-z**, **0-9**) or the hyphen (**-**).

#### 2.1.1 Acronyms for Projects

At the start of a project we define a project acronym. This acronym is intended to be used in file and folder names.

Whenever we want to indicate the relation to a certain project in a file or folder name, we use the project acronym in exactly the typing that was defined. This is important as we want to distinguish between raw data, processed data and project results in our Folder Structures.

The project acronyms are defined in a simple text file **PROJECTS.txt** in the **//server/projects\$** folder, see Project Folder Structure.

#### 2.1.2 Acronyms for Organisations

It is very important to know the owners of data. Therefore we define unique acronyms for the owners of data that we use. The acronyms are defined in a special file **ORGANISATIONS.txt**



## Chapter 3

# Folder Structures

### 3.1 Project Folder Structure

We said that we want to concentrate on the folder structures within the project folders. Nevertheless, we would like to give a recommendation on how the project folders could be organised within its top level folder. In this structure, there are no subfolders for the different departments any more.

```
//server/projects$  
- PROJECTS.txt  
- project-1/  
- project-2/  
- project-3/
```

In the `projects$` folder

- each subfolder name should appear in the file `PROJECTS.txt`
- there should not be any folder on the top level that does not represent a project.
- there should be no other files on the top level as the files that are described in this documentation.
- there are no subfolders representing departments any more. The mapping of projects to departments is done in the file `PROJECTS.txt`

### 3.2 Rawdata Folder Structure

We will create a network folder `//server/rawdata$` in which all files have set the read-only property. We suggest to store raw data by project first and by the organisation that owns (i.e. generated, provided) the data second. This could look like this:

```
//server/rawdata$  
- ORGANISATIONS.txt  
- PROJECTS.lnk [Symbolic Link to PROJECTS.txt in //server/projects$]  
- flusshygiene  
  - bwB  
  - kwB  
  - uba  
  - ...  
- ogre  
  - kwB  
  - bwB
```

- uba
- ...
- ...

Restrictions/Conventions:

- Each top-level folder should represent a project, i.e. should be defined in the top level file **PROJECTS.txt**.
- Each possible owner should be defined in the top level file **ORGANISATIONS.txt**.
- The naming convention for the organisations is the same as for projects.

## Chapter 4

# Raw Data

As raw data we define data that we receive from a device or from a project partner.

Most of our research results are based on data. We acknowledge the importance of raw data by

- storing them in a special place where it is specially secured
- describing them with metadata

Rawdata are stored in the rawdata folder structure





# Chapter 5

## Metadata

### 5.1 Special Files

We propose to define some special files that contain metadata related to files and folders. To indicate that these files have a special meaning, the file names are all uppercase.

#### 5.1.1 File PROJECTS.txt

This file contains the project acronyms as we want to use them e.g. as top-level folder names in our project folder structure. The projects are grouped by department.

Possible content of PROJECTS.txt:

```
# Department SUW (Surface Water)
dswt: DSWT
flusshygiene: Flusshygiene
kuras: KURAS
mia-cso: MIACSO
monitor: MONITOR
ogre: OGRE
reliable-sewer: RELIABLE_SEWER
sema: SEMA
sema-berlin: SEMA Berlin
sema-berlin-2: SEMA Berlin 2
spree-2011: SPREE2011
spree-2011-2: SPREE2011 "reloaded"

# Department GRW (Groundwater)
optiwells: OPTIWELLS
optiwells-2: OPTIWELLS 2
wellma: WELLMA

# Department WWT (Wastewater Treatment)
...
```

In the file PROJECTS.txt the project acronyms appear in alphabetical order. They map the acronym to a project name or a project title and the year of the start of the project.

**Question:** Do we already have a place where “official” metadata about projects are stored? If yes, the acronym could be included there. But then, everybody should know about it!

### 5.1.2 File ORGANISATIONS.txt

Possible content of ORGANISATIONS.txt

```
bwb: Berliner Wasserbetriebe  
kwb: Kompetenzzentrum Wasser Berlin  
uba: Umweltbundesamt
```

## Chapter 6

# Data Processing

### 6.1 Automatic Data Import into R

In the following we describe how data can be imported into the R-Programming Environment

#### 6.1.1 Import Data From One Excel File

- Save the original file in the rawdata zone.

#### 6.1.2 Import Data From Many Excel Files

##### 6.1.2.1 Files Are In the Same Format

Import Excel files of the same format by

- defining a function that is able to read the data from that file
- calling this function in a loop for each file to import.

##### 6.1.2.2 Files Are In Different Formats

We developed a general approach of importing data from many Excel files in which the formats (e.g. more than one table area within one sheet, differing numbers of header rows) differ from file to file.