

2025

Basics of Research Data Management

Lib4RI

Felder, Fabian Ulrich-Nath, Moushumi, Bachofner, Anusch Förster, Christian





Agenda

Topic	Time	
Introduction	9:00 – 9:20	
Open Science, FAIR, RDM		
Data Collection, Processing and Analysis	9:20 — 10:00	
File folders, naming, versioning, formats		
Pause	10:00 — 10:10	
Documentation	10:10 – 10:30	
README & Metadata		
Storage, Preservation, and Sharing	10:30 — 10:50	
Repositories, data availability statements, licensing		
Pause	10:50 — 11:00	
RDM Services and Support		
• Eawag	11:00 – 11:30	
• Empa	11:30 – 12:00	



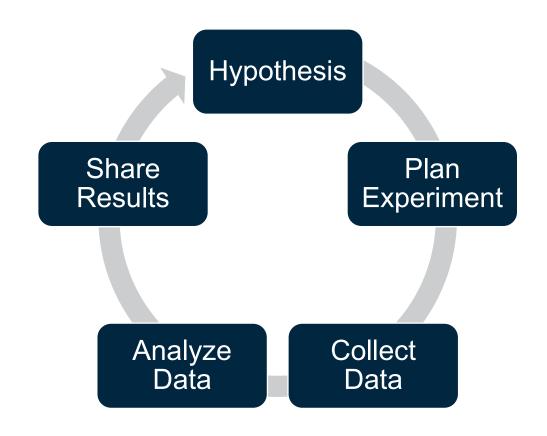
Welcome Open Science, FAIR, and RDM







Scientific Method vs Research Data Management



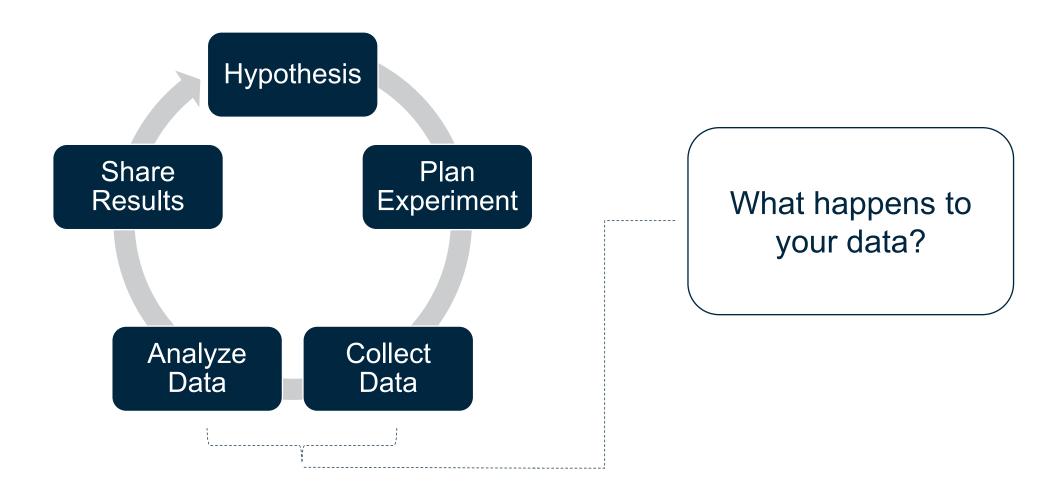








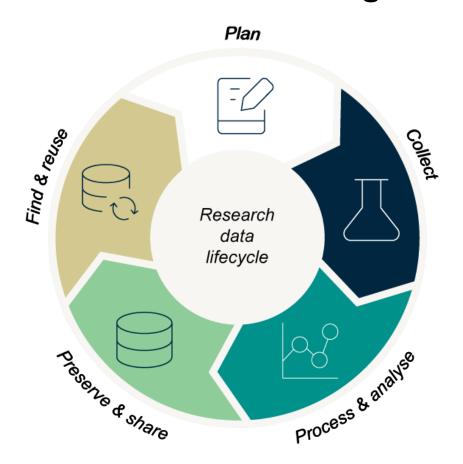
Scientific Method vs Research Data Management







Research Data Management

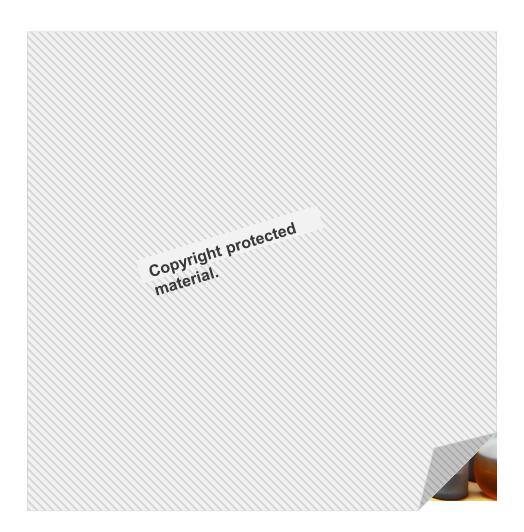


Research data management (RDM) is the process of organizing, storing, preserving, and sharing data that is generated or used in a research project.





Why bother?







What is Open Science?

- O Q1: What's the first word or phrase that comes to mind when you hear Open Science.
- Q2: I believe Open Science is the future of research.

www.menti.com



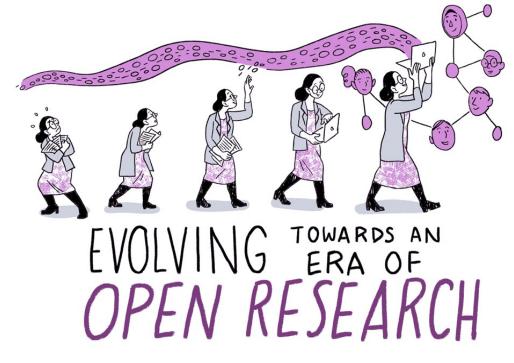


Open Science is a culture change.

Accessible

Transparent

Collaborative



The Turing Way project. CC-BY 4.0. DOI:10.5281/zenodo.3332807.











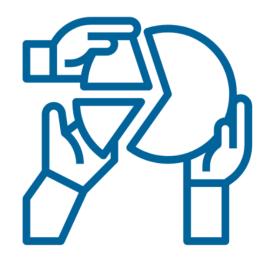


Open Science is a culture change.

Reproducibility/Replicability Crisis



Global Collaborations



Arm Okay, CC BY 3.0, via Wikimedia Commons





Open Science is the future.

Policy

Incentives

Trainings

Infrastructure

Empa

SNSF

<u>Data</u> <u>Management</u>

Campus

Zenodo

Eawag

Horizon

ERIC

PSI

SciCat

<u>WSL</u>

EnviDat





Open Science is not just a data dump.



Cezary p, CC BY-SA 4.0, via Wikimedia Commons



Fotolia/TrudiDesign. <u>Tackling trash – DW – 04/23/2012</u>







What are some ways you can make your data FAIR?



Fotolia/TrudiDesign. <u>Tackling trash – DW – 04/23/2012</u>

www.menti.com



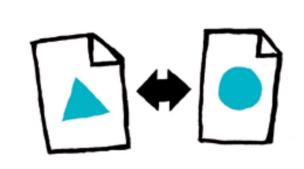




FAIR are guiding principles for data management.









Metadata – where can people find the data?

DOIs

Metadata – how can others access your data?

Metadata is accessible.

Data may be restricted.

File formats – nonproprietary, standard.

Vocabulary – standard.

Metadata – provides sufficient context to understand your data.

Licenses – determines how others can reuse your data.

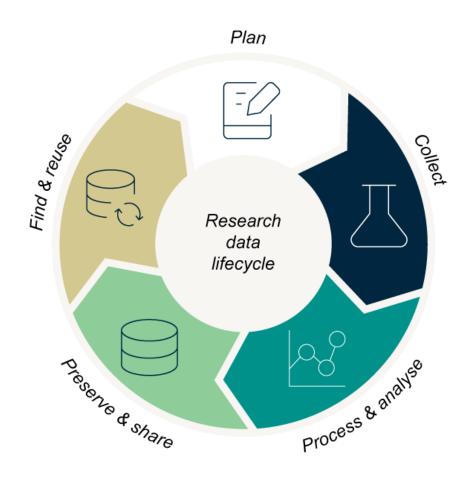
Studio 4 minutes 34, Studio Lendroit.com, CC BY-SA 4.0, via Wikimedia Commons







Research Data Management











- Folder structures
- File naming
- File formats
- Versioning



Documentation (activity)

README & Metadata



- Storage
- Repositories



- Data availability statements
- Licenses

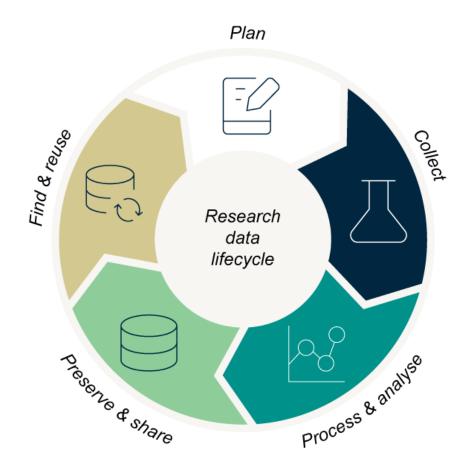


How to do RDM





Planning







Planning – the data management plan (DMP)

What types of data will be collected or generated?

What type documentation will you provide with the data?

How will data be stored during and after the project?

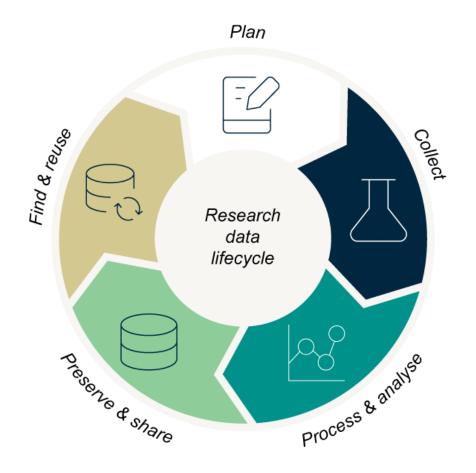
How will your data be shared?

How can it be accessed?

Are there any ethical, legal, or security issues to address?



Data Collection, Processing, and Analysis



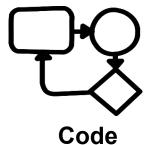




Data Collection

Data

observational, experimental, simulation...



Applications, scripts...



Metadata

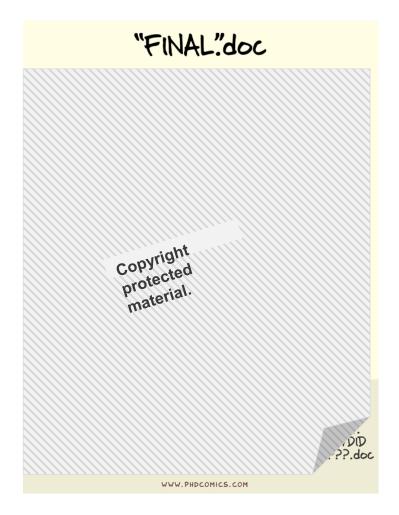
Structured information associated with data (and code)

The Who, What, Where, Why & How of data





What to do with all this data, code, metadata?

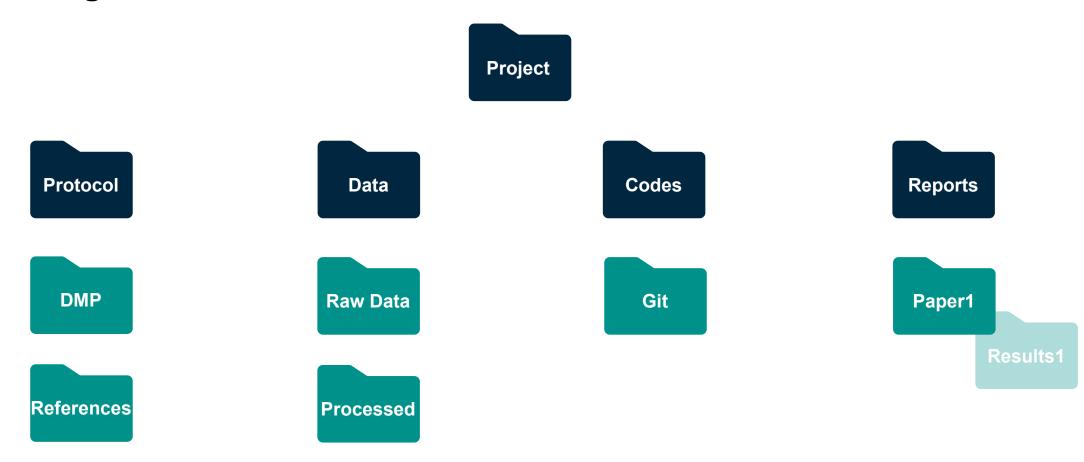


File Organization File Naming Versioning





Organize data: File/Folder Structures







File Naming

Date ProjectName DocumentType Version FileExtension

Be consistent

Include Date: YYYYMMDD

Document Type: data, code, results, paper...

No special characters or signs @°§°#¬@

Project Name: short (< 30 characters) and descriptive

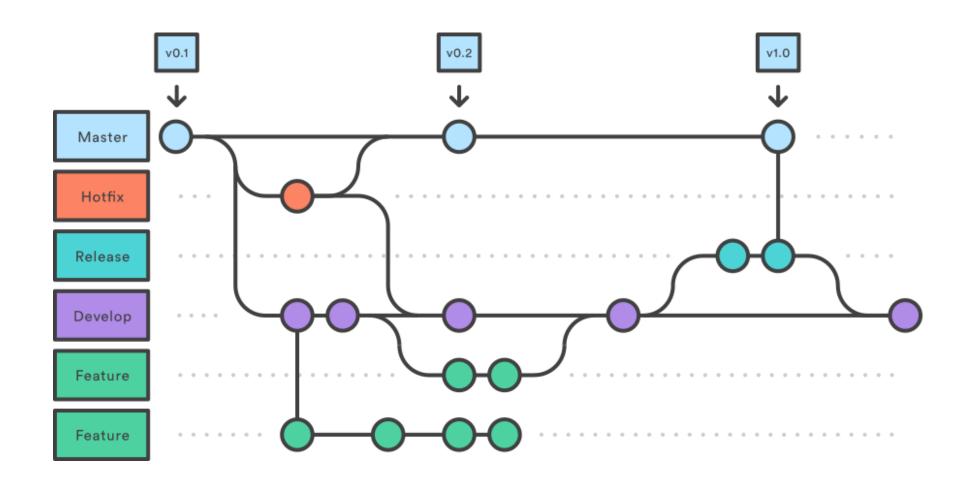
Include version: v0.8, v1.2...

Tip: Already started your project?
Use Bulk Rename Utility (Windows), Renamer 6 (Mac), Rename/Thunar Bulk Rename (GNU/Linux).





Software Version Control







Tools for Software Version Control



















(i)











• CLI (Command Line interface)

GUIs (Graphical User Interfaces) https://git-scm.com/downloads/guis





Data Versioning

Raw rev. 0

Proc. lev. 1 rev. 0

...

Proc. lev. m rev. 0

Raw rev. 1

Proc. lev. 1 rev. 1

...

Proc. lev. m rev. 1

:

:

:

Raw rev. n

Proc. lev. 1 rev. n

. . .

Proc. lev. m rev. n





Tools for Data Versioning



Data Version Control (https://dvc.org)



Git Large File Storage (https://git-lfs.com)



Lake FS(https://docs.lakefs.io)





Tool for data and code organization

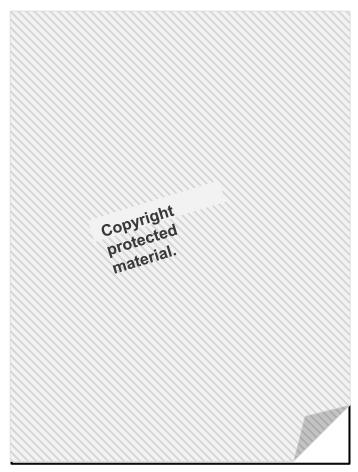


Renku (https://renku.readthedocs.io/en/stable/index.html)





What to do with all this data, code, metadata?



https://xkcd.com/2116/

Your files are organized, appropriately named, clearly versioned...but what's the point if nobody can open the #*\$)@* file??

File Formats







File formats: open and non-proprietary

Data type	Recommended file formats	
Text	 PDF/A Plain Text coded as ACII. UTF-8 or UTF-16 XML 	
Spreadsheet	CSV (NEAD)	
Images	TIFF (uncompressed or lossless compressed)PNG	
Code	Languages with free environments (e.g. Py or R UTF-8 format of ASCII text)	
Audio	FLACWav	

Aim for open and lossless formats
If you are using a proprietary format (ex. MS word extensions xlsx, docx), consider adding an additional format.







Activity: Prepare your data. (15 min, individual or pairs)

Folder Structure

- Does the folder and file organization make sense?
- Separated raw vs processed data?

File Naming

- o Is it consistent?
- Does it include:
 - YYYYMMDD
 - Project Name
 - Document Type
 - Version
 - o File Extension

Versioning

- Do you have a versioning system?
- Is there a tool you can adopt to help you?

File Formats

- Are your files in nonproprietary, open formats?
 - $\rightarrow xlsx \rightarrow csv$
 - \circ docx \rightarrow pdf
 - o images → tiff



Data vs Metadata

https://dataedo.com/cartoon/tag/data-vs-metadata







Documentation: README and Metadata

Feature	README	Metadata
Audience	Humans	Machines (and sometimes also humans)
Format	Free-form text (txt, md)	Structured (XML, JSON)
Scope	How and Why	What and Who
Standards	Flexible	Discipline-specific
Functionality	Help others understand and use data	Facilitate searchability and machine processing





Activity: README (15 min, pairs or trios)

General Information

Sharing/Access Information

Data & File Overview

Methodological Information

Data-specific information for: [filename]

Activity: Data-specific information for: [filename]

Select a dataset/folder/file and fill in the following:

- Number of variables:
- Number of cases/rows:
- Variable List: list variable name(s), description(s), unit(s) and value labels as appropriate for each>
- Missing data codes: st code/symbol and definition
- Specialized formats or other abbreviations used

For access to entire template:

carpentries-incubator.github.io/scientific-metadata/files/AUTHOR DATASET ReadmeTemplate.txt







Metadata

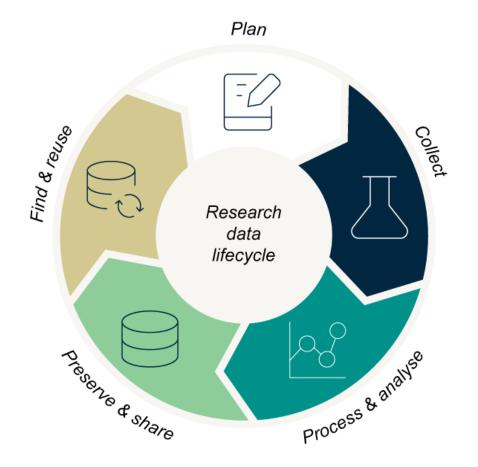
- Definition: Structured data that contains information about other data, but is not the content of the data.
- o Metadata is very subject specific. The following directories are helpful:
 - o Digital Curation Centre (https://www.dcc.ac.uk/guidance/standards)
 - o RDA Metadata Standards (https://rdamsc.bath.ac.uk)
 - o Fairsharing (https://fairsharing.org)
- Recommendation: Stick to a list of defined terms (controlled vocabulary) and don't use synonyms to describe the same object (e.g. picture or image)

Activity: Find out which metadata standard is relevant to your field (7 min)





Storage, Preservation, and Sharing







Storage

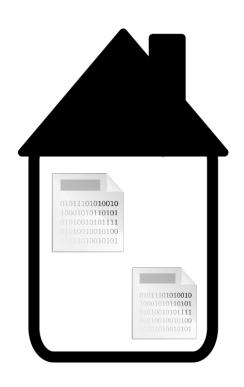
O Q: Where is your data stored and how is it backed-up?

www.menti.com





Storage: 3-2-1 backup







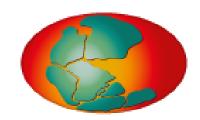


Data repositories: publication vs preservation















For alternatives: https://www.re3data.org/

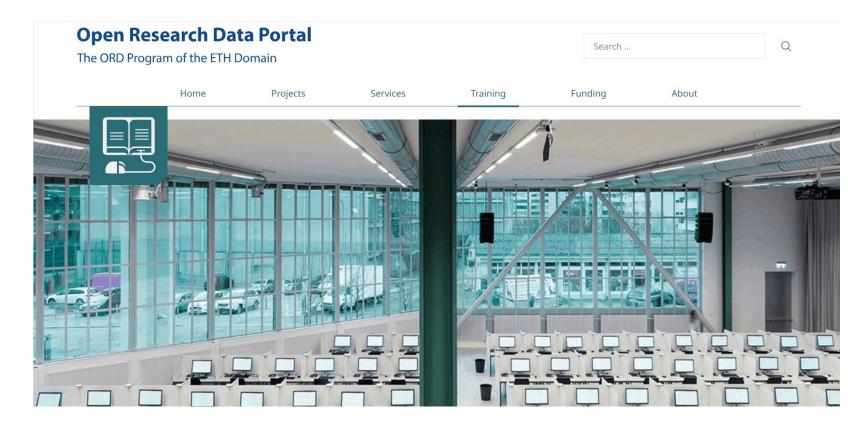








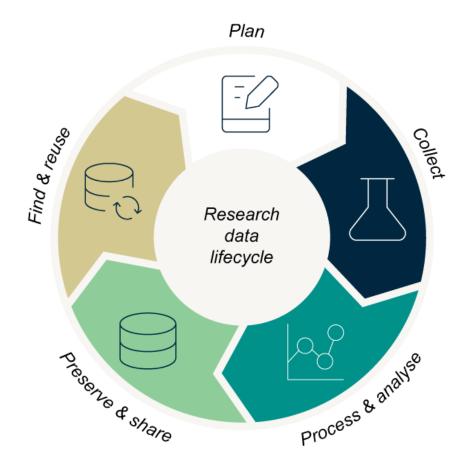
Learn more: Data Management Campus



https://open-research-data-portal.ch/training/



Finding and Reusing







Data availability statements

Where is the data stored (repository name and DOI)?

How can the data be accessed (open, restricted, available upon request)?

How can the data be reused (licensing)?

Very Famous (Open Access) Journal

Example

The data supporting this study's findings is openly available in [Repository Name] at [DOI]. The dataset includes x,y,z and is available under [License].







Licensing

Data

Code











Licensing for Code

Copyleft

- Examples: GPL, LGPL
- Use cases: for projects maintaining open-source is the priority (operating systems, applications, platforms)

Permissive

- Examples: MIT, Apache, BSD
- Use cases: for projects which encourage wide-spread adoption and commercial use (libraries, frameworks, tools)

For more information: <u>Licenses – Open Source Initiative</u>

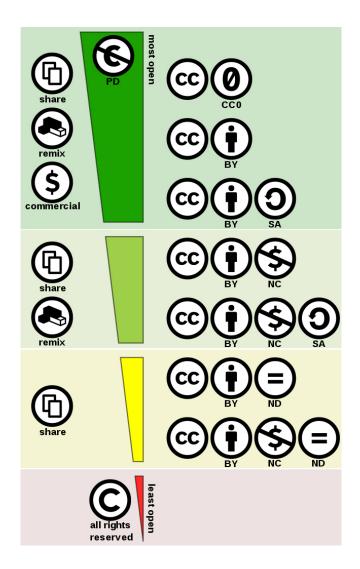




Licensing for Data

cc creative commons









Planning – the data management plan (DMP)

What types of data will be collected or generated?

What type documentation will you provide with the data?

How will data be stored during and after the project?

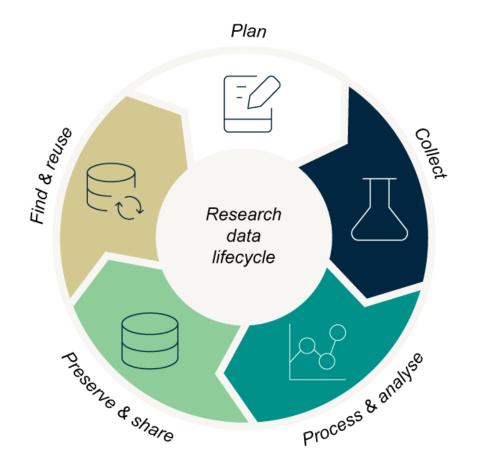
How will your data be shared? How can it be accessed? • File folders, file naming, versioning, file formats

README and Metadata

- 3-2-1 rule
- Repositories
- Data availability statements
 - Licensing



Research Data Management





Lib4RI – Excellent Services for Excellent Research.

fabian.felder@lib4ri.ch moushumi.ulrich-nath@lib4ri.ch T: + 41 58 765 57 00