

How to teach good research data management to next generation researchers?

Syed Ashfaq Hussain Shah
Prof. Dr. Ing. Frank Petzold (AI, TUM)
NFDI4Ing Conference 2023
September 27, 2023



Agenda

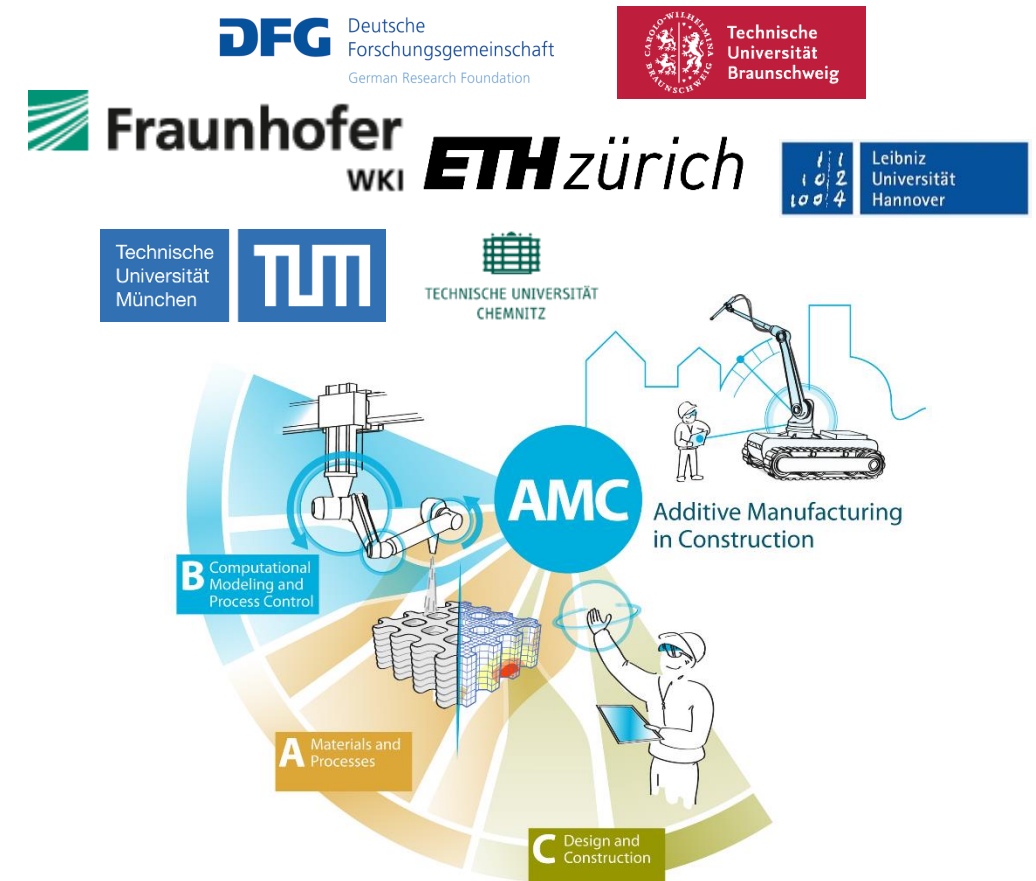
- Introduction
- Agenda of Data Management
- Teaching and guiding methodologies
- Guiding materials
- Means of distributions
- Interactive sessions
- On boarding and migration processes
- Processes to create guiding contents
- Events to improve guidance
- Actions to improve RDM practices
- Tactics to enforce compliance for RDM
- Conclusion



Image source: <https://unsplash.com/>

Introduction: TRR277 AMC (Collaborative Research Center)

- Information infrastructure project
- Hosted by **2 universities** at **2 distinct locations**, multiple participating institutes
- **3 focus** areas **22 interdisciplinary** research groups, over **120 members**
 - Materials and processes
 - Computational modelling and control
 - Design and construction
- Collaborated with **Leibniz-Rechenzentrum (LRZ)**, **Gauß-IT center**, **libraries** of the both universities and **ProLehre** at TU Munich
- Approaches developed for **in person**, **digital medium** as well as for **hybrid** situations e.g. in Lockdown time.



Credit: <https://amc-trr277.de/>

Introduction: Good Research Data Management (RDM)

- Compliant with: -
 - FAIR principles
 - Open data/ Open Science practices
 - Reproducible results (where apply)
- Resilient and mitigate changing compliance
 - (Inter) National funding body/ research council requirements
 - DFG Good Research Practice
 - Other stake holders specific requirements
 - Research organising body
 - Publishers, (Scientific) Community

DFG Deutsche
Forschungsgemeinschaft
German Research Foundation

Adopted by the Senate of the DFG at September 30, 2015

Deutsche
Forschungsgemeinschaft

DFG Guidelines on the Handling of Research Data

Release date: 21th December 2021

Handling of research data

Checklist for planning and description of handling of research data in research projects

Introduction: Samples of questions

- **What is** research data management (RDM) and why one should care?
- **What are** these terms **vocabulary, metadata, ontology, DMP?**
- **What makes research data** for project?
- **What information relating to practices** is part of research work?

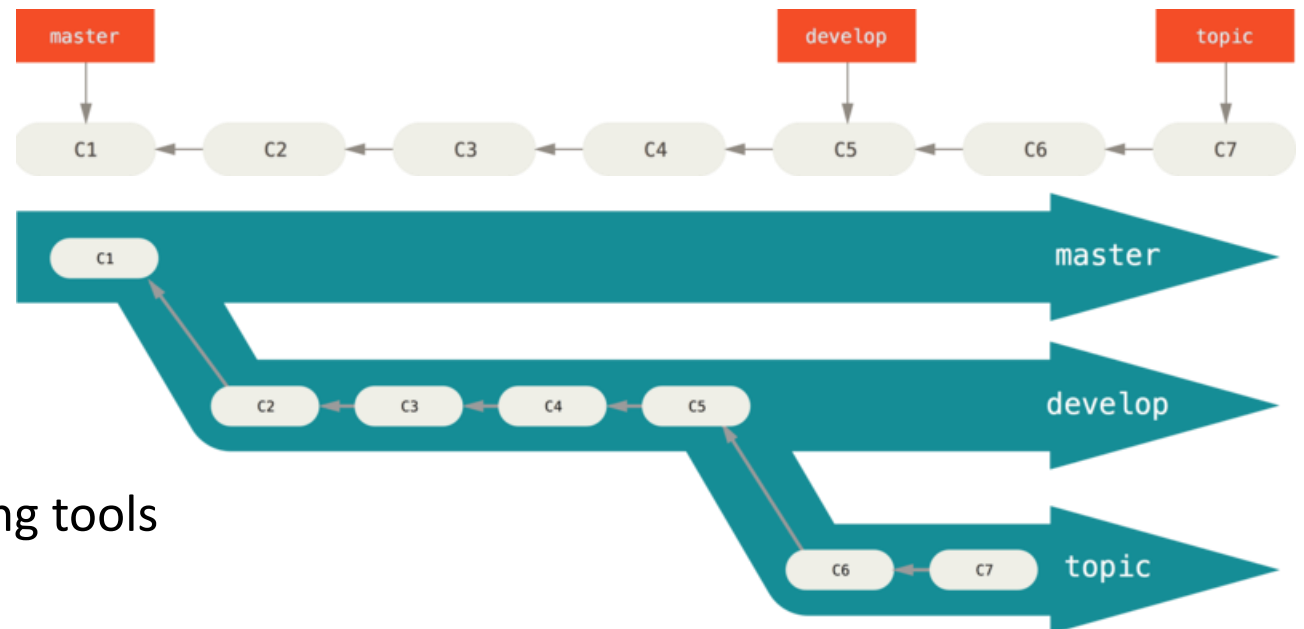
- Are there ways to **speed up the upload?**
- Is the **cable network** better than **wireless network?**
- **How to install** required software?

Introduction: Samples of questions contd...

- What to do when data is on a **terminal which cannot be connected** with network/ internet?
- I **do not have storage** device with **sufficient capacity** to port data on an internet connected terminal?
- **How to organise** data in **collaborative research**?
- I have data **generated over time**. I **can not find where** it should be uploaded?
- **How** the provided **system works** e.g. to record the practice, upload the data.
- I **could not understand** what you just demonstrated/ explained.

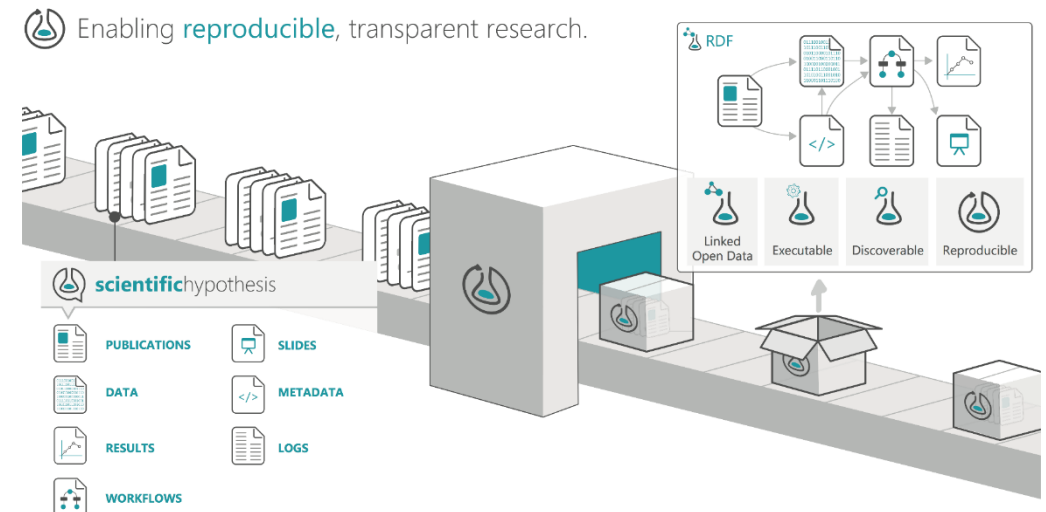
Agenda of data management

- Fundamentals of data management
 - Naming convention
 - Versioning
 - Identification of independent component and organisation
 - Documentation
 - Logging
 - README/ Docs/ Comments
 - Packaging/ Bundle
 - Unique ID
- Digital means
 - Purpose/ use based data processing tools



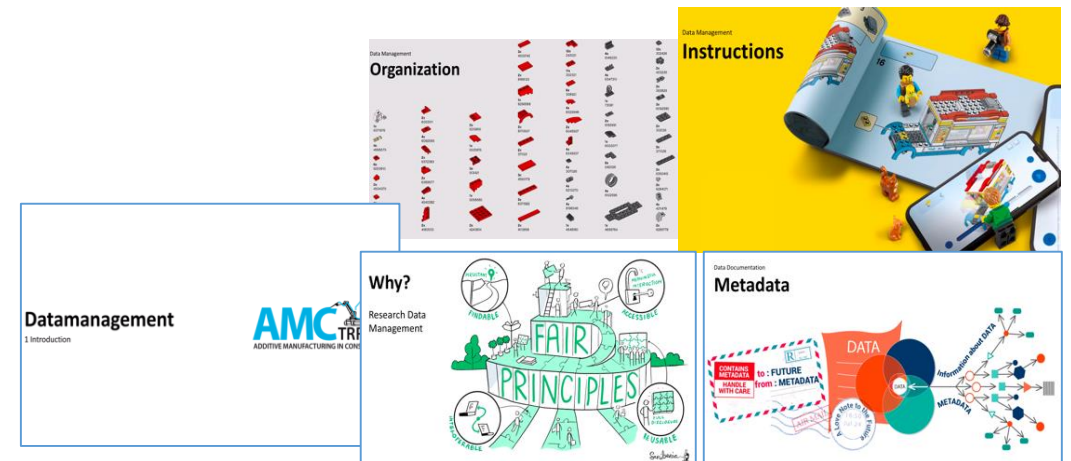
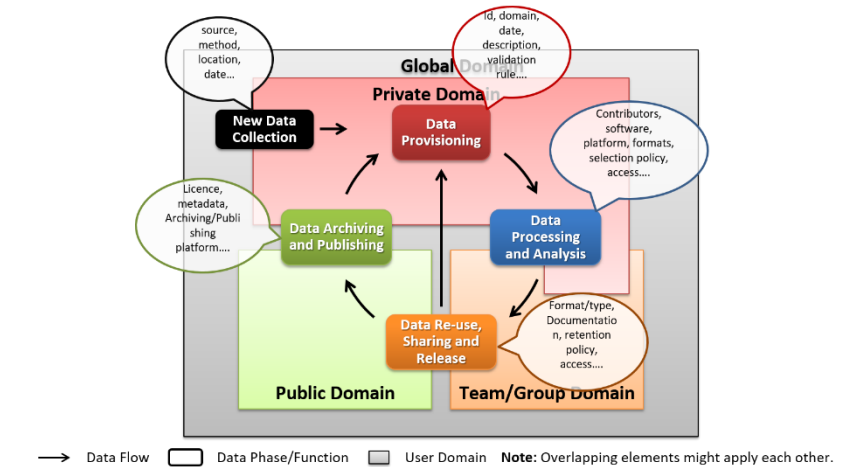
Agenda of data management contd...

- Research data management + Fundamentals of data management
 - Documentation of **practices**
 - Metadata
 - Data management plan (DMP)
 - User/ Technical guide
 - Maintenance of **provenance** and **provisioning**
 - **Compliance** with **long term** archive compliant **format**
 - Data **anonymisation** and **pseudonymisation**
 - **Universally** unique IDs.
 - ORCID ...
 - DOI ...
 - **Attribution** e.g. Licencing
 - **Compliant** system
- A well thought **change control** agenda



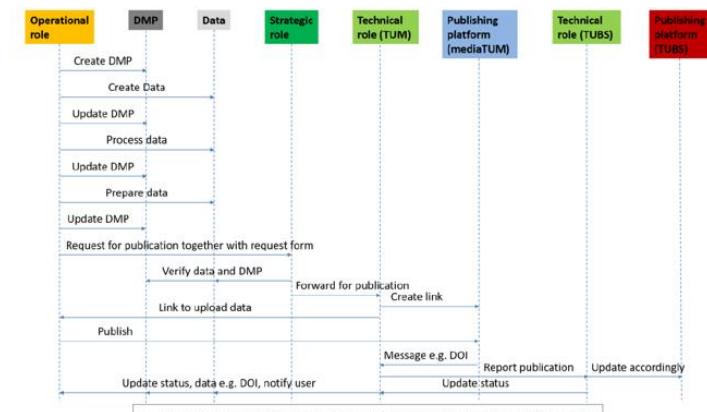
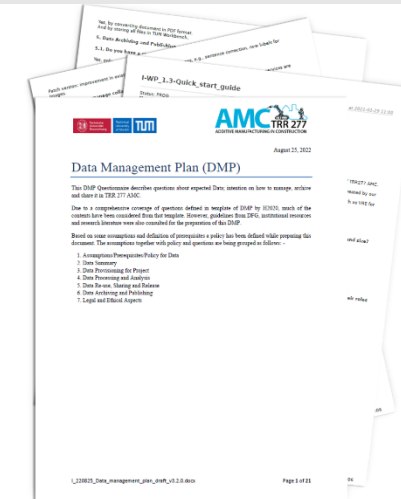
Teaching and guiding methodologies

- **Policy** and definition of general rules
- Custom Data Management Plan (**DMP**)
- **Guiding** materials
 - Printable materials
 - Multimedia contents
- (Simplified) data **models, templates**
- Conducted one to one/ group (in person/ online) **sessions**
- Dedicated **weekly hours** for support
- **Evaluation, feedback** and **follow ups**
- **Complemented** with: -
 - Screenshots, screencasts
 - diagrams, illustrations and animations
 - Glossaries and definitions



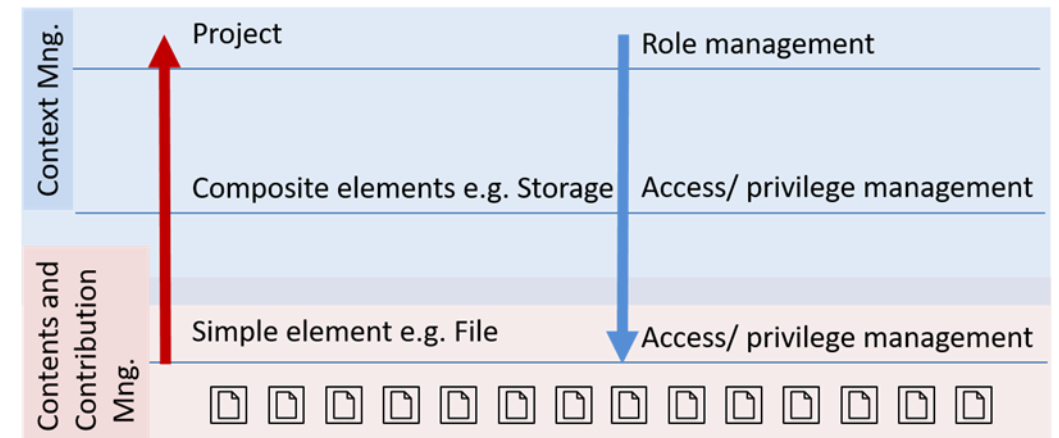
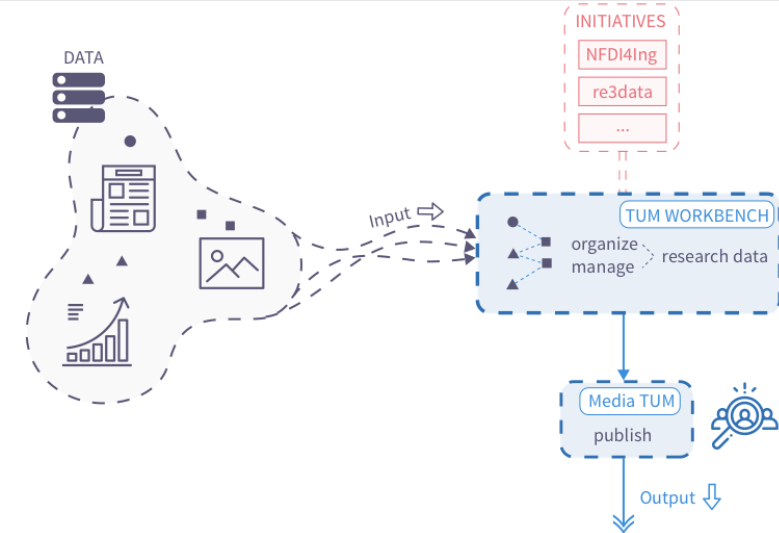
Policy and definition of general rules

- Research relevant **data** and **outcomes**
- Official **RDM platform** and **infrastructure**
- Official Data Management Plan (**DMP**)
- **Roles** and **responsibilities** of the participants
- Types of **research outcomes/ data**
- Data file **types/ formats** to communicate research outcome and long term archiving
- Common/ mandatory **metadata** standard
- **Naming conventions** including versioning
- **Use cases** of research outcomes and DMPs
- **Workflows** for publishing



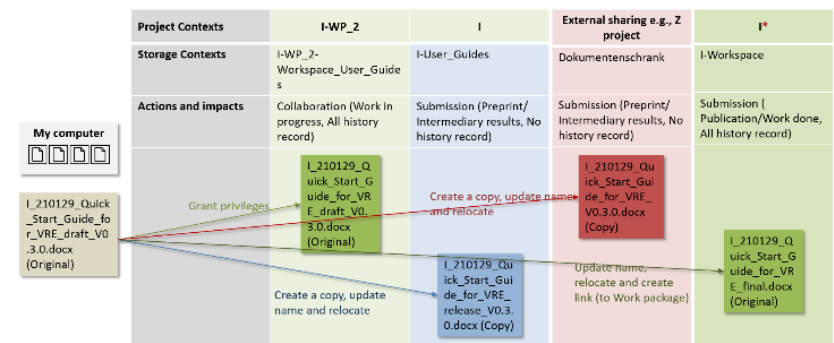
Guiding materials: Official RDM platform and infrastructure

- User **competency** requirements
- Functional **requirements**
- Hardware/ Software **requirements**
- **Conceptual/ theoretical** foundations
 - Common entities and functions
 - User interface design/ webpages
 - Navigation and exploration workflows
 - Definition of tools and controls
 - Access rights and privileges
 - Data and context management
 - Data input and persistence procedures
 - Built-in RDM features e.g. UIDs, backups, communication and messages, metadata



Guiding materials: Official RDM platform and infrastructure contd...

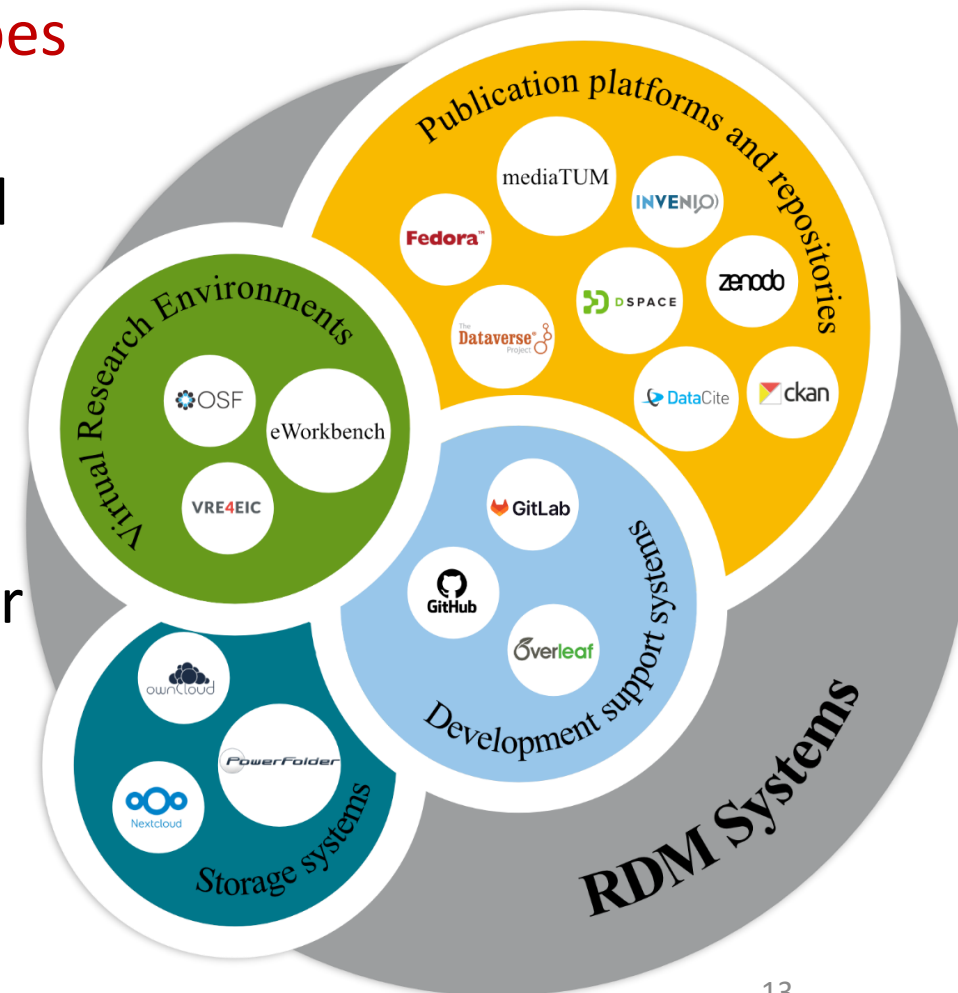
- User **scenario**
- **Practical** part
 - User interaction schemes and outcomes
 - Metadata
 - Data management plan (DMP)
 - Network and communication protocols
 - Integration and interaction with external systems
- **Conventions** and **best practices**
- **Documentation** possibilities
- Sharing and release **scenarios**
- **Configuring** functions and features
- Frequently asked questions (FAQ)
- Access to software, support and services



* It is an optional step. Keeping data under its original context is always recommended.

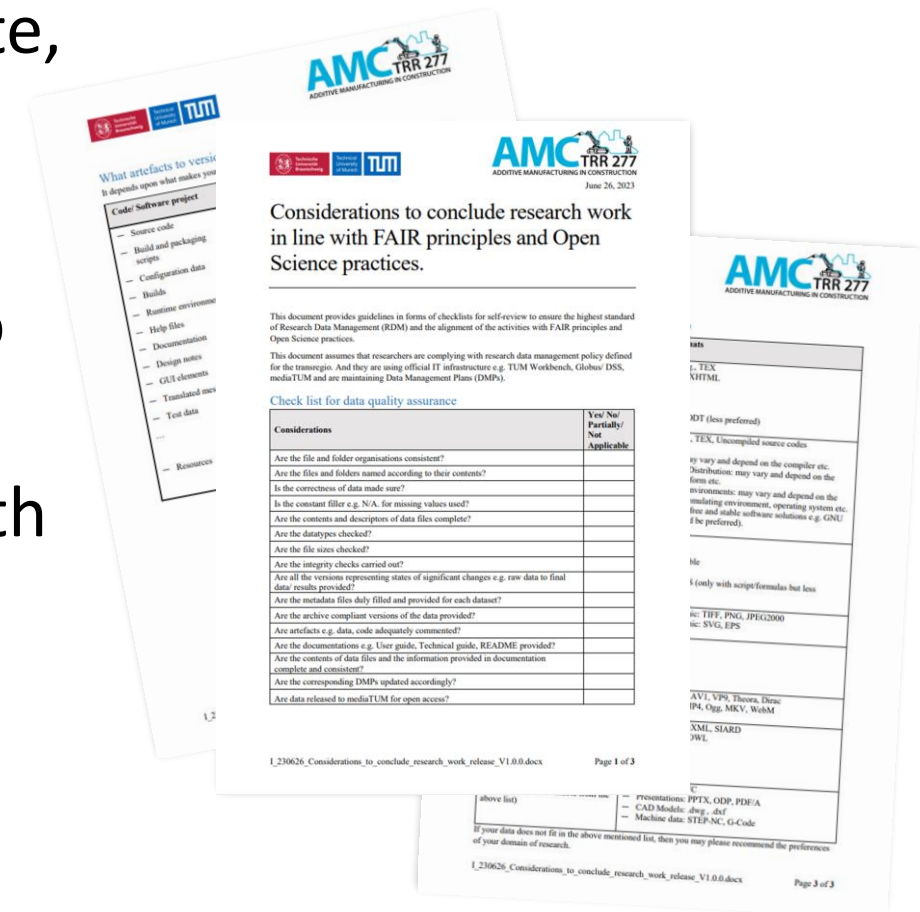
Guiding materials: Research practices

- Data **identification** based of its **roles, states, types**
- Folder/ file **organisations** and **structures**
- **Identification** of atomic/ independent workload
- **Collaborative** work e.g. task based
- **Distribution** of larger/ collaborative work packages
- Specific and general rules for **decision**: -
 - e.g. research data, information necessary for reproducibility, to register data in RDM system, significant change in data to record
- Research and data support **systems** and their **roles**



Guiding materials: Research practices contd...

- **Adoption** of common standards e.g. DataCite, DOI, ORCID...
- **Filling** DMPs
- **Applications** of DMP for atomic workload to the whole CRC
- **Concluding** research task/ project in line with **FAIR** principles and **Open Science** practices.



Guiding materials: Supplementary materials

- Standard **reusable** contents
 - e.g. metadata templates, folder/ file structures, publication templates, labels to improve naming convention, DMP
- Practical relevant research data **examples**
 - e.g. lab experiment, material mixing, survey, simulation, code, journal/ conference publication data management plan
- Collections of standard external **contents** and **resources**
 - e.g. List of licences, tools, metadata standards

<ul style="list-style-type: none"> ▪ draft ▪ final ▪ preprint ▪ release ▪ Publish 	<ul style="list-style-type: none"> ▪ raw ▪ process/ process_name ▪ release/ publish ▪ release_candidate ▪ stable_release 	<ul style="list-style-type: none"> ▪ dev ▪ master ▪ build ▪ release ▪ release_candidate ▪ stable_release
---	--	--

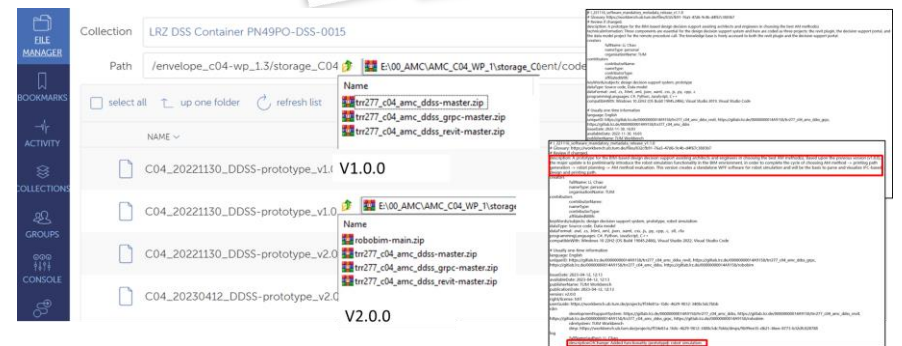
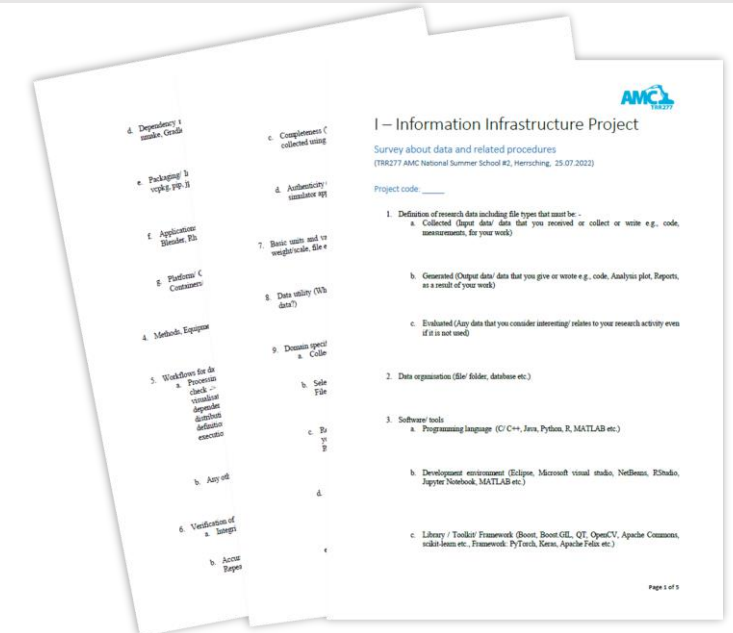
What are relevant Research Datasets? - TRR227 practical Data Examples

- **Experiments on material mixing:** details about "ingredients", results and procedure together with dependencies e.g. environmental data
- **Survey:** survey template, responses, audio/ visuals of interviews, survey results/ report, details about result compilation tools and methods, basis of template design
- **Simulation:** input, output data, details about software, methods and workflow
- **Code:** source form, guide/ readme, details and parameters about packaging system/ framework/ compiler/ dependencies/ version system
- **Paper publication (conference/ journal):** paper (source and published form e.g. PDF), bib data, images, presentation slides, video, supporting data e.g. as per simulation data
- ...



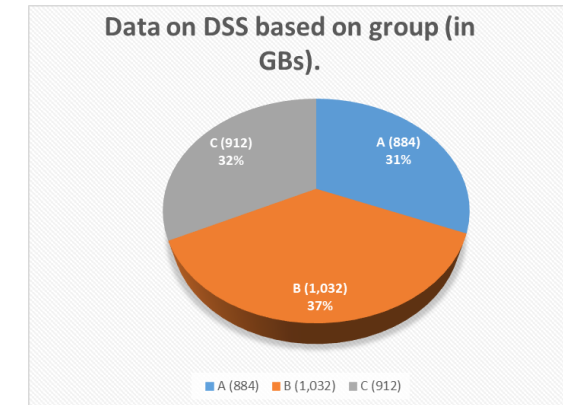
Guiding materials: Supplementary materials contd ...

- Comprehensive **survey**
- Recommendation and hints about **tools, best practices** and **workflows**
 - Coding e.g. build and packaging, dependency management, compiler, execution environments
 - Data/ metadata creation e.g. collection, selection, retention, transformation, basic units and value ranges
 - Data verification/ quality assurance measures and strategies e.g. integrity, accuracy, completeness, authenticity checks
 - Data analysis, experiments, digital representation
 - Data acquisition, integration, anonymization, pseudonymization, release, archiving

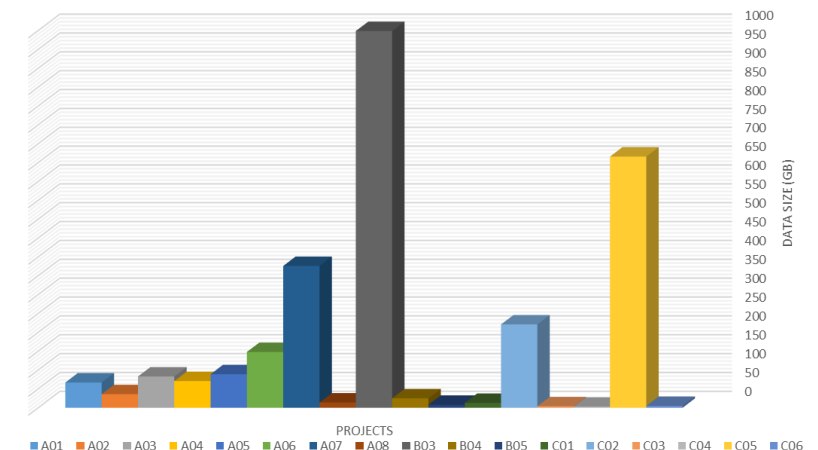


Guiding materials: Supplementary materials contd ...

- **Use cases** of collaboration and cooperation
- Systems and tools for **automation**
 - Guides for additional tools and systems
- Detailed **review**, **evaluation** and **feedback** reports and **follow ups**
- Considerations/ **Checklist** to conclude research work

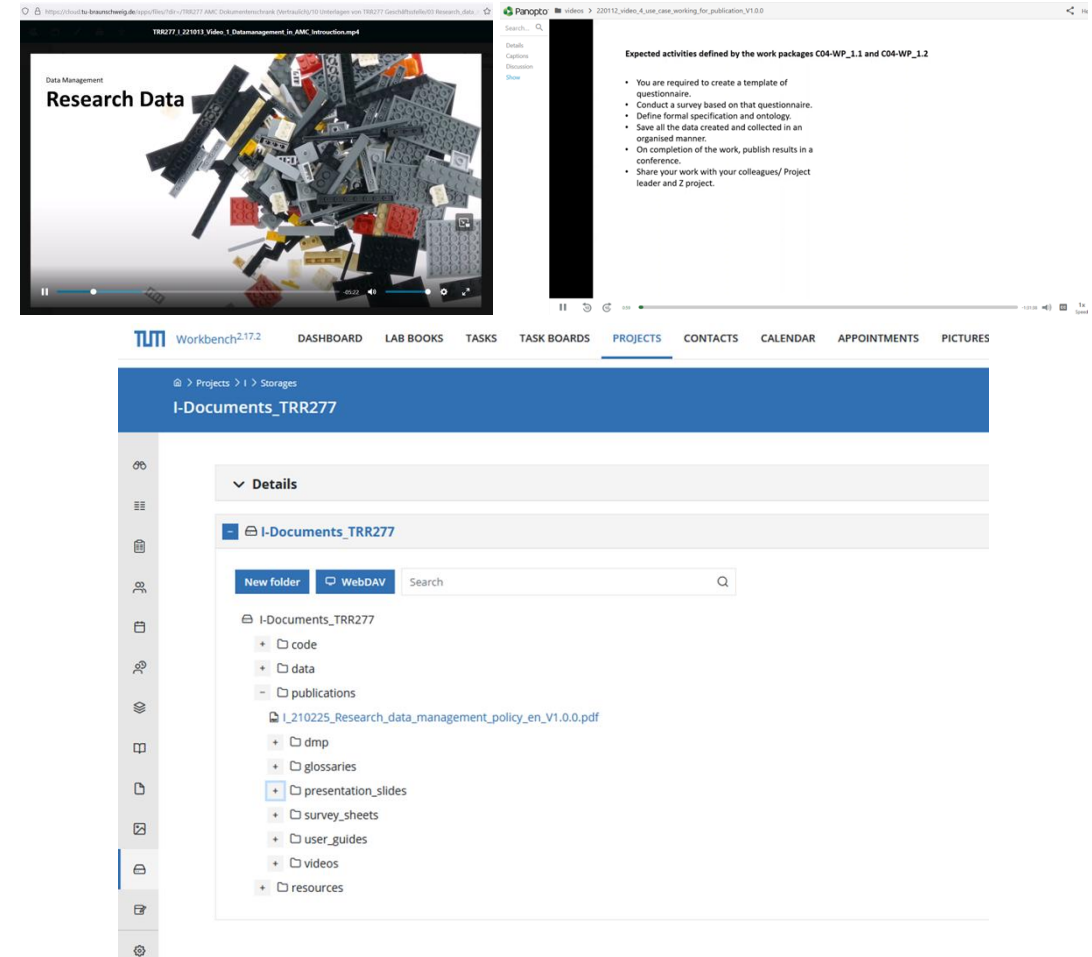


Data on DSS (total: 2,828 GBs)



Means of distribution, communication and imparting trainings

- Online **streaming**
 - Panopto
 - TUBS Cloud/ NextCloud
- File **hosting** (downloadable)
 - Official platform (TUM Workbench)
 - TUBS Cloud/ NextCloud
- **Live** sessions (One to one/ Group/ Collective)
 - Online e.g. via Zoom
 - In person
- Official platform based communication tools
- Institutional email and communication services

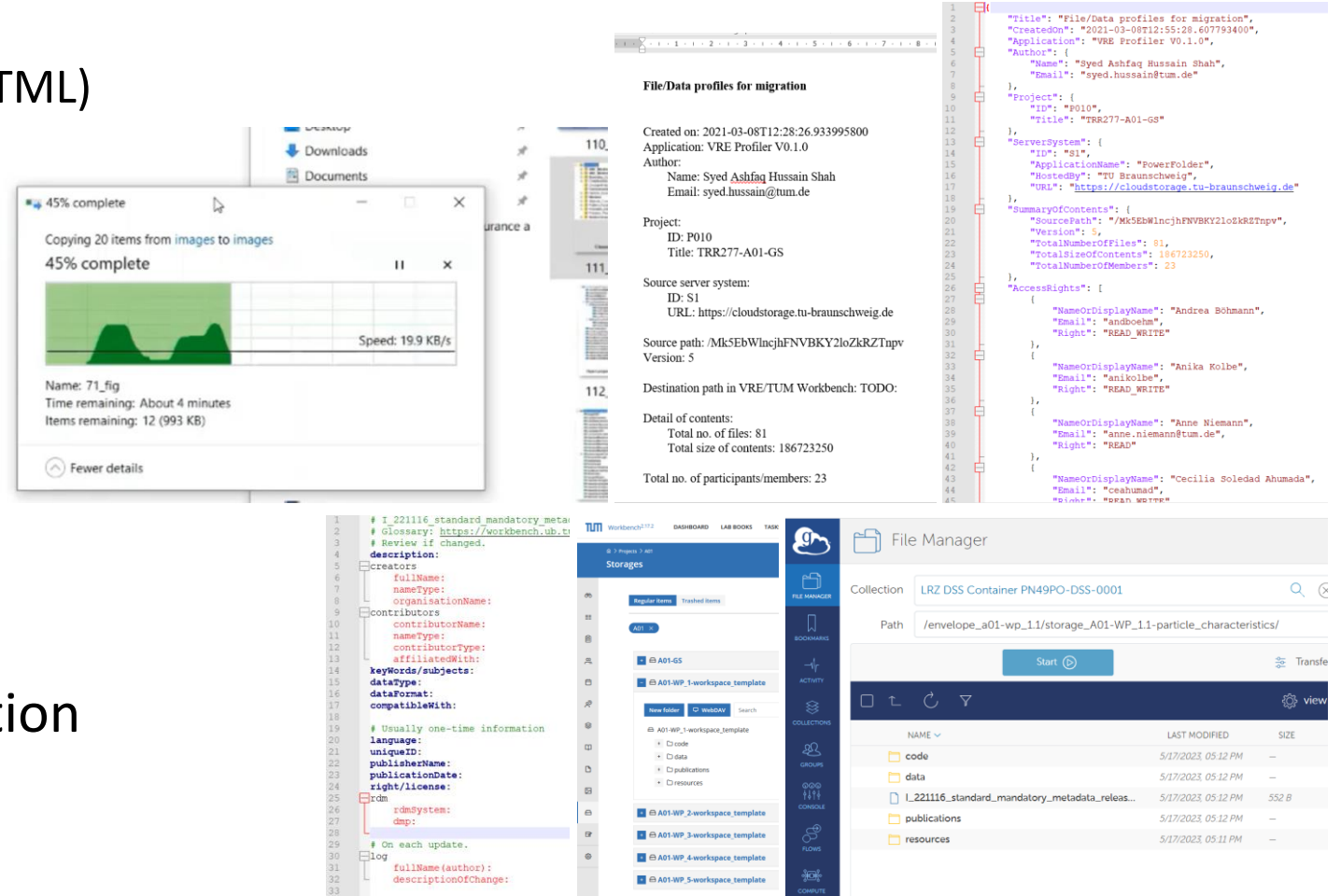


Interactive sessions and workshops' themes

- **Introduction** to the VRE system and concepts
- User and data **migration**
- VRE for **data managers/ RDM**
- RDM and its **basics**
- RDM **tools** and **systems**
- **How to use** (tools specific sessions)
- Software/ tools **test** sessions
- **Workflows** e.g. for publication, experiments, collaboration
- Practical research data **examples**
- **Do it together** session
- **Weekly** consultation session
- Issue/ Topic/ Task specific **consultation** sessions
- **Q & A** sessions
- **Reviews** and **feedbacks** sessions
- Consultations based on the reviews and feedbacks
- **Quarterly progress meeting**
- **Summer school**

On boarding and migration processes

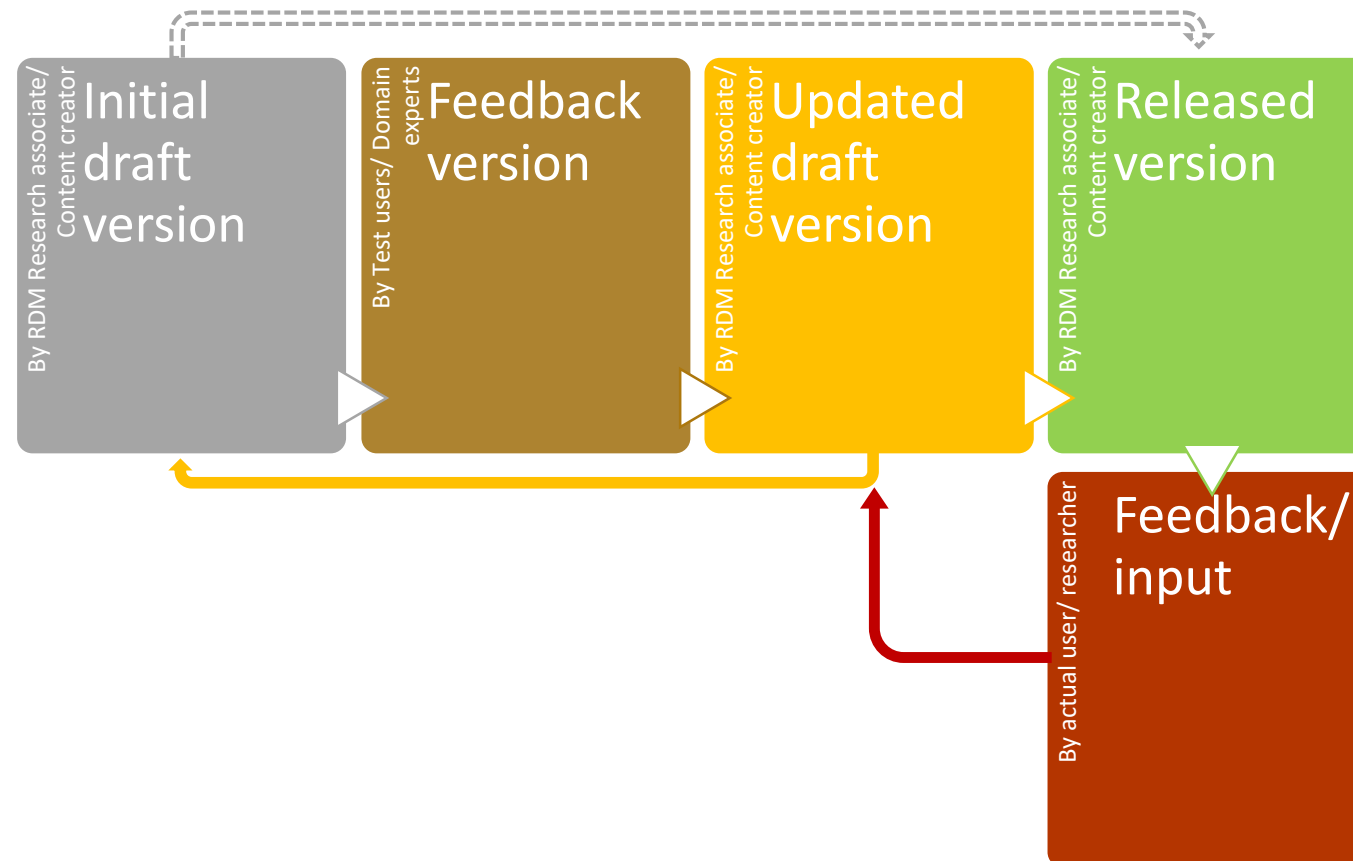
- **Profiling** of existing data
 - Human readable format (.docx/HTML)
 - Machine readable format (JSON)
- **Verification** of data profiles
- **Test and allocation** of resources
 - Project structures
 - Storage structures
 - Access rights
 - Metadata templates
- **Interactive** do it together migration sessions



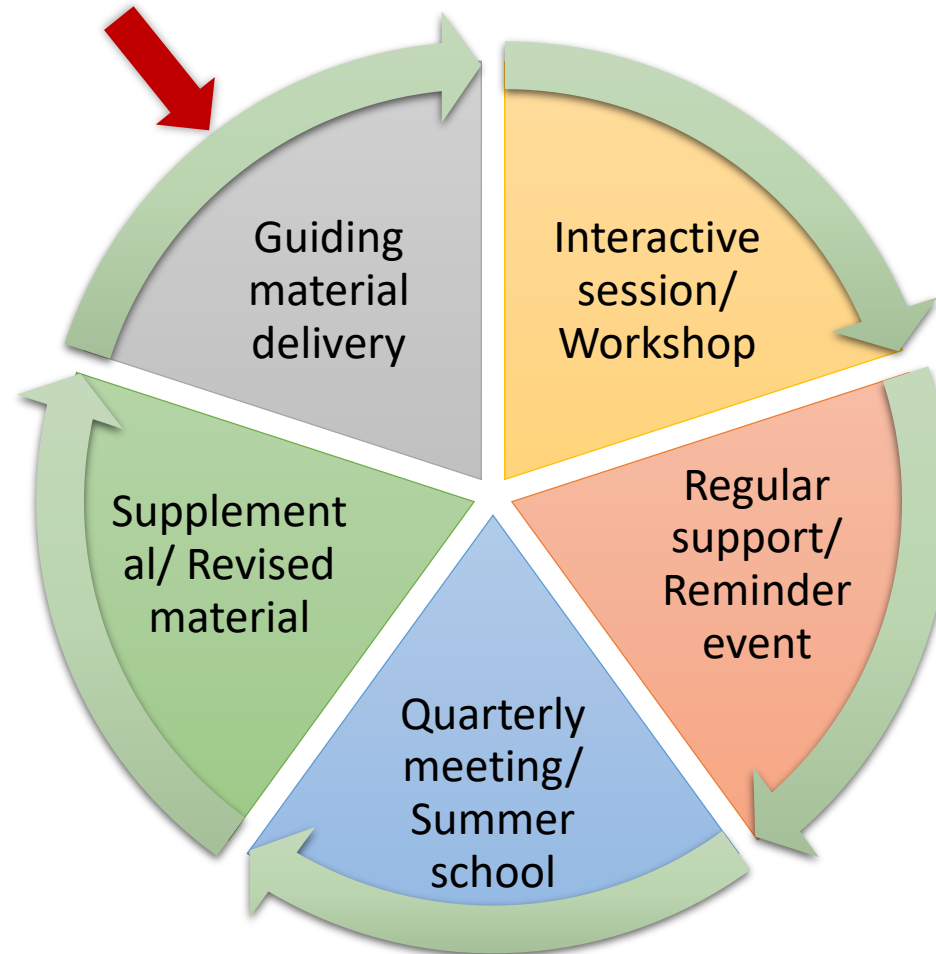
```
1 {"title": "File/Data profiles for migration",  
2 "creator": "2021-03-08T12:28:26.933995800",  
3 "application": "VRE Profiler V0.1.0",  
4 "author": {  
5 "name": "Syed Ashfaq Hussain Shah",  
6 "email": "syed.hussain@tum.de"  
7 },  
8 },  
9 "project": {  
10 "id": "P010",  
11 "title": "TRR277-A01-GS"  
12 },  
13 "serverSystem": {  
14 "id": "S1",  
15 "applicationName": "PowersFolder",  
16 "hostedBy": "TU Braunschweig",  
17 "url": "https://cloudstorage.tu-braunschweig.de"  
18 },  
19 "summaryOfContents": {  
20 "sourcePath": "/Mk5EbWlncjhFNVBKY2loZkR2Tnpv",  
21 "version": 5,  
22 "totalNumberOfFiles": 81,  
23 "totalSizeOfContents": 186723250,  
24 "totalNumberOfMembers": 23  
25 },  
26 "accessRights": [  
27 {  
28 "nameOrDisplayName": "Andrea B\u00f6tman",  
29 "email": "andboehm",  
30 "right": "READ_WRITE"  
31 },  
32 {  
33 "nameOrDisplayName": "Anika Kolbe",  
34 "email": "anikolbe",  
35 "right": "READ_WRITE"  
36 },  
37 {  
38 "nameOrDisplayName": "Anne Niemann",  
39 "email": "anne.niemann@tum.de",  
40 "right": "READ"  
41 },  
42 {  
43 "nameOrDisplayName": "Cecilia Soledad Ahumada",  
44 "email": "ceahumada",  
45 "right": "READ_WRITE"  
46 }  
47 ]  
48 }
```

```
1 # I_221116_standard_mandatory_metadata  
2 # Glossary: https://workbench.tu-braunschweig.de  
3 # Review if changed.  
4 description:  
5 - creators:  
6 - fullName:  
7 - nameType:  
8 - organisationName:  
9 - contributors:  
10 - contributorName:  
11 - nameType:  
12 - contributorType:  
13 - affiliatedWith:  
14 - keywords/subjects:  
15 - dataType:  
16 - compatibleWith:  
17  
18 # Usually one-time information  
19 language:  
20 uniqueID:  
21 publisherName:  
22 publicationDate:  
23 right/license:  
24 rdm:  
25 - rdmsystem:  
26 - dmp:  
27  
28 # On each update.  
29 Log:  
30 - fullName (author):  
31 - descriptionOfChange:  
32  
33
```

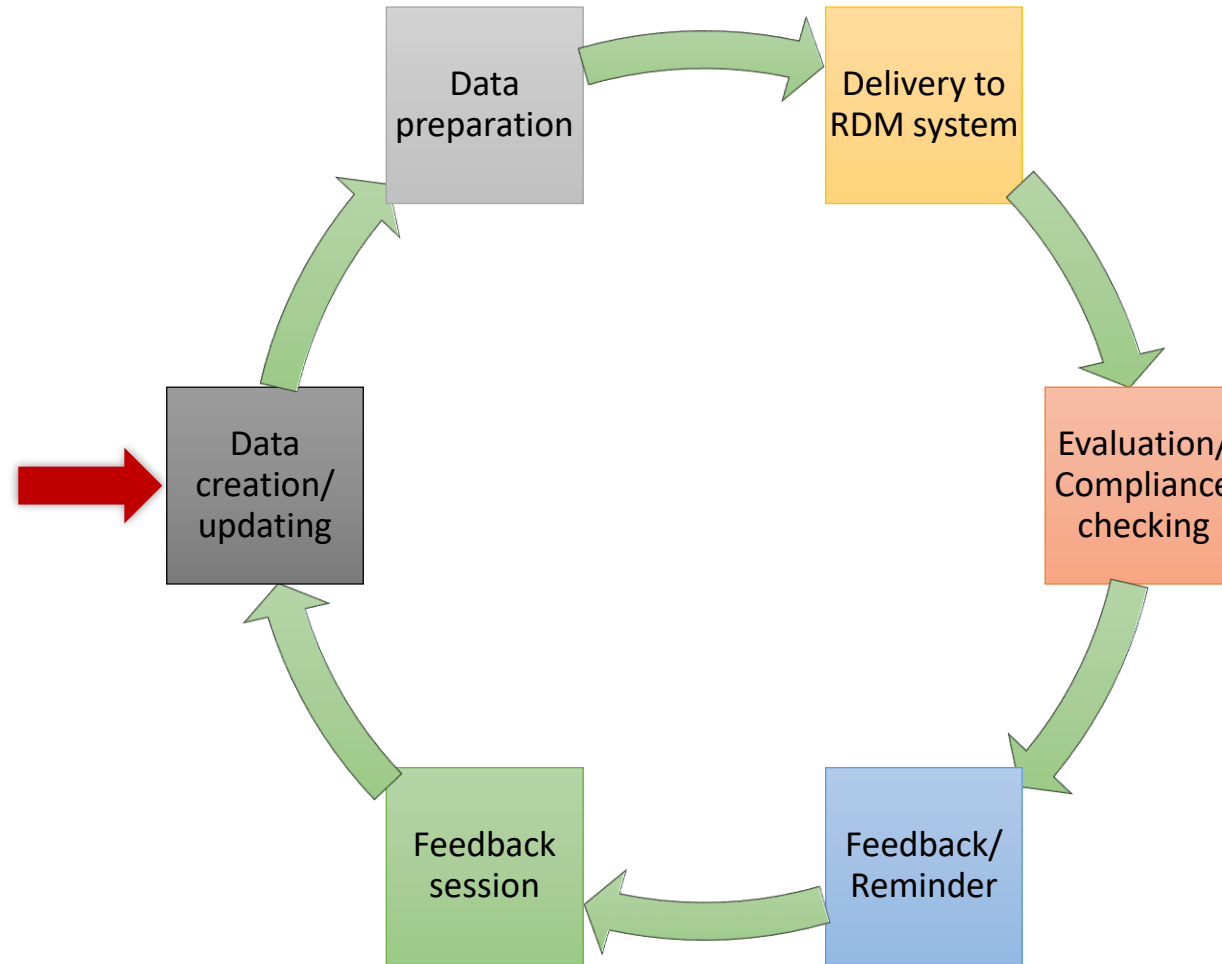
Processes and phases to create contents



Events to improve guidance and understanding

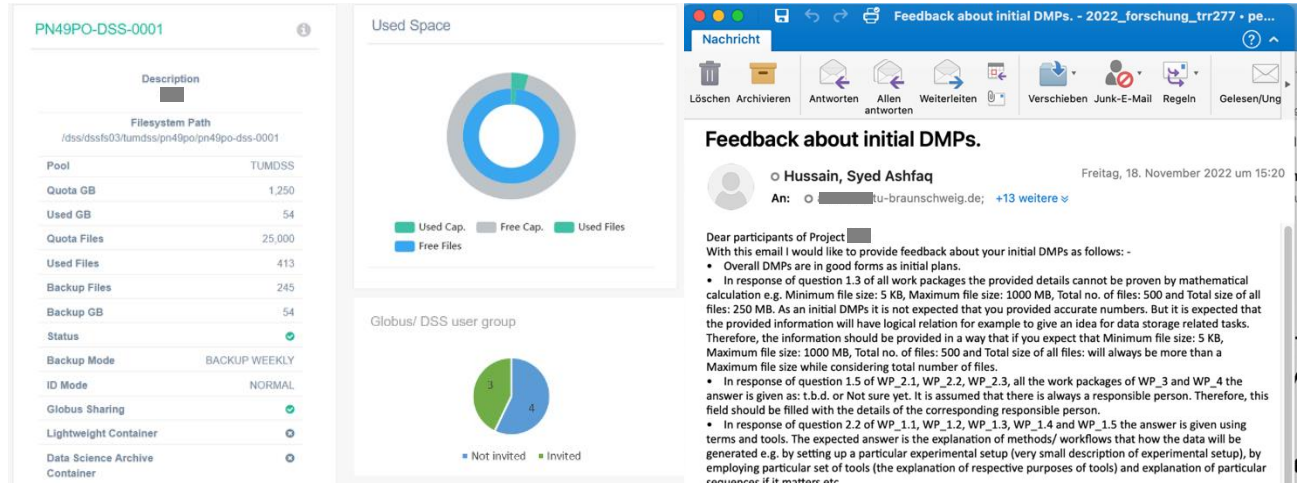


Processes and actions to improve RDM practices



Key compliance strategies

- **Provision** of initial structures and templates for a smooth and easy **kick start**
- **Examples** of best practices
- Detailed **interactive** sessions
- Interactive **do it together** sessions
- **Q & A** rounds
- Detailed **reviews** and **feedbacks**
- **Interactive review** and **feedback** sessions
- **Supplementary** materials
- **Comprehensive survey** about RDM practices
- **Collaborative work/ assistance** in case of difficulties
- **Continuously improving** guiding materials and IT solutions.



PN49PO-DSS-0001

Description

Filesystem Path: /dss/dsfs03/tumdss/pn49po/pn49po-dss-0001

Pool	TUMDSS
Quota GB	1,250
Used GB	54
Quota Files	25,000
Used Files	413
Backup Files	245
Backup GB	54
Status	✔
Backup Mode	BACKUP WEEKLY
ID Mode	NORMAL
Globus Sharing	✔
Lightweight Container	○
Data Science Archive Container	○

Used Space

Donut chart showing Used Cap., Free Cap., and Used Files.

Globus/ DSS user group

Donut chart showing Not invited (3) and Invited (4).

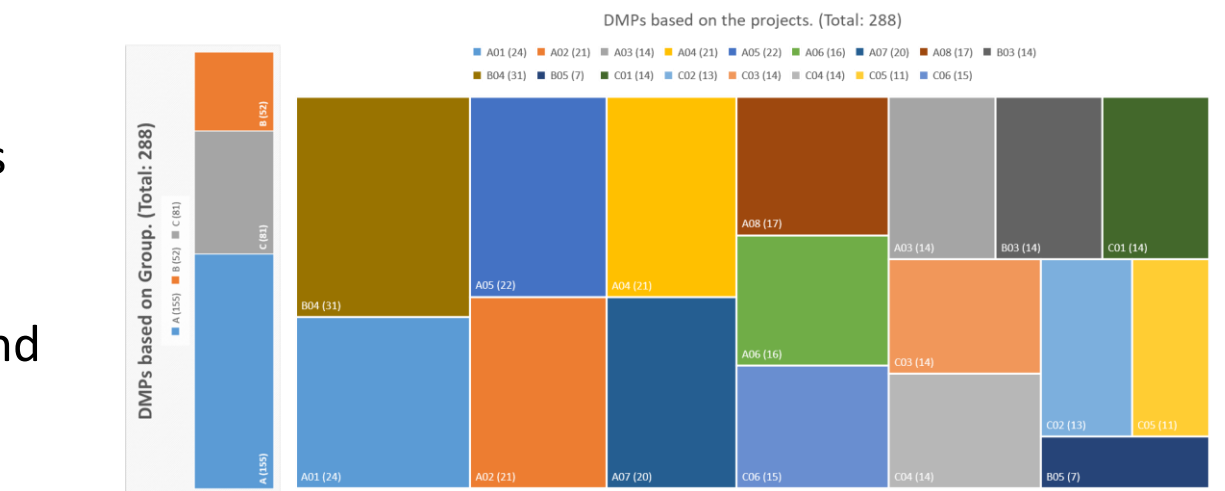
Feedback about initial DMPs.

Hussain, Syed Ashfaq

Dear participants of Project [redacted],

With this email I would like to provide feedback about your initial DMPs as follows:

- Overall DMPs are in good forms as initial plans.
- In response of question 1.3 of all work packages the provided details cannot be proven by mathematical calculation e.g. Minimum file size: 5 KB, Maximum file size: 1000 MB, Total no. of files: 500 and Total size of all files: 250 MB. As an initial DMPs it is not expected that you provided accurate numbers. But it is expected that the provided information will have logical relation for example to give an idea for data storage related tasks. Therefore, the information should be provided in a way that if you expect that Minimum file size: 5 KB, Maximum file size: 1000 MB, Total no. of files: 500 and Total size of all files: will always be more than a Maximum file size while considering total number of files.
- In response of question 1.5 of WP_2.1, WP_2.2, WP_2.3, all the work packages of WP_3 and WP_4 the answer is given as: t.b.d. or Not sure yet. It is assumed that there is always a responsible person. Therefore, this field should be filled with the details of the corresponding responsible person.
- In response of question 2.2 of WP_1.1, WP_1.2, WP_1.3, WP_1.4 and WP_1.5 the answer is given using terms and tools. The expected answer is the explanation of methods/ workflows that how the data will be generated e.g. by setting up a particular experimental setup (very small description of experimental setup), by employing particular set of tools (the explanation of respective purposes of tools) and explanation of particular processes if it matter etc.



Tactics to enforce compliance for RDM

- **Obligation** through Research data **policy**
- **Advantages** and **confidence** in official platforms
- **Influence** of research organising board/ committee
- **Influence** through research supervisor/ supervisory board
- **Reminders** during collective meeting events, through communication channels before and after the schedule is missed.
- Presentation of **results**, **summaries** and **updates** during meetings/ collective events
- Presentation of **feedback** and **evaluation** reports
- Narrating **adverse impacts** and **implications** of non compliance
 - e.g. data loss, consequences of further approvals

Future work

- Marginalising and estimating the weight of impact.
- Further features' definition and development for automation.
- Evaluation of different approaches.

- It is a process of continuous improvement with retrospectives.

Thank you for your attention

Questions?

Acknowledgements

The work is being conducted as part of the collaborative research centre 'Additive Manufacturing in Construction - The Challenge of Large Scale. Funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) - Project Number 414265976 - TRR 277".