

Introduction to RDM

RDM in HPC - Challenge or Chance?

Online | 26.02.2024

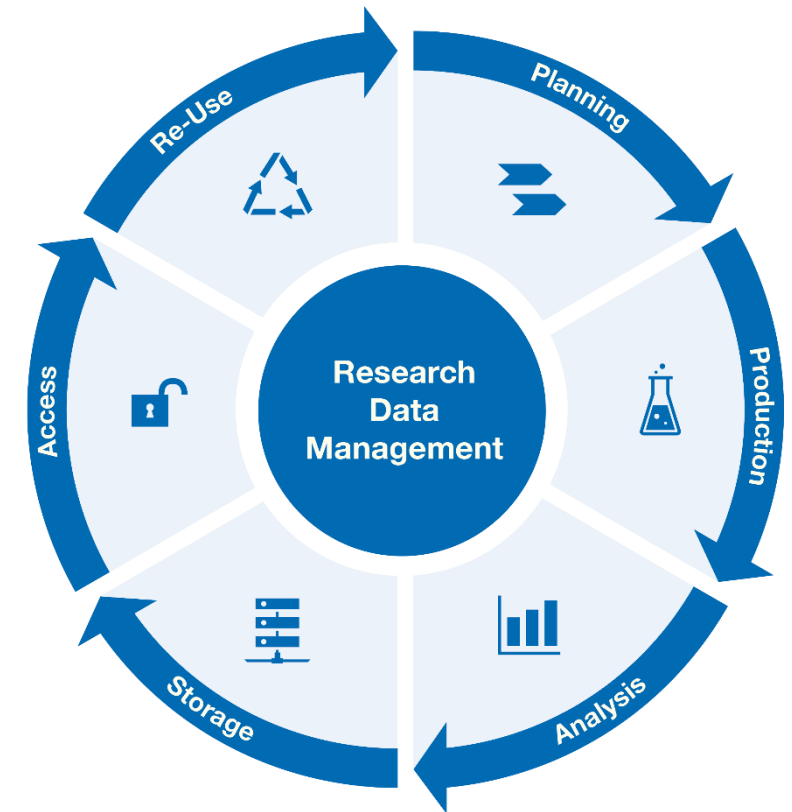
Introduction to RDM

Marcel Nellesen (RWTH Aachen University)

RDM in HPC - Challenge or Chance?

- When should I start with RDM?
 - Short Answer: NOW!!!
 - Long Answer:
Ideally before you start with the project. However
it is never too late.

- Research projects are a long term commitment
- Data is constantly created
 - First draft / proposal
 - Setup
 - Measurement
 - Analysis
 - Reports
 - Reusage



- FAIR
 - Findable
 - For humans and for machines
 - Accessible
 - Providing information and infrastructure to access the data
 - Interoperable
 - Interaction with applications, for analysis, storage and processing
 - Reusable
 - Metadata and data need to be well described to allow reproducibility and reuse in other contexts

– Findability

- (meta)data are assigned a globally unique and persistent identifier
- Data are described with rich metadata
- Metadata clearly and explicitly include the identifier of the data it describes
- (meta)data are registered or indexed in a searchable resource

– Accessibility

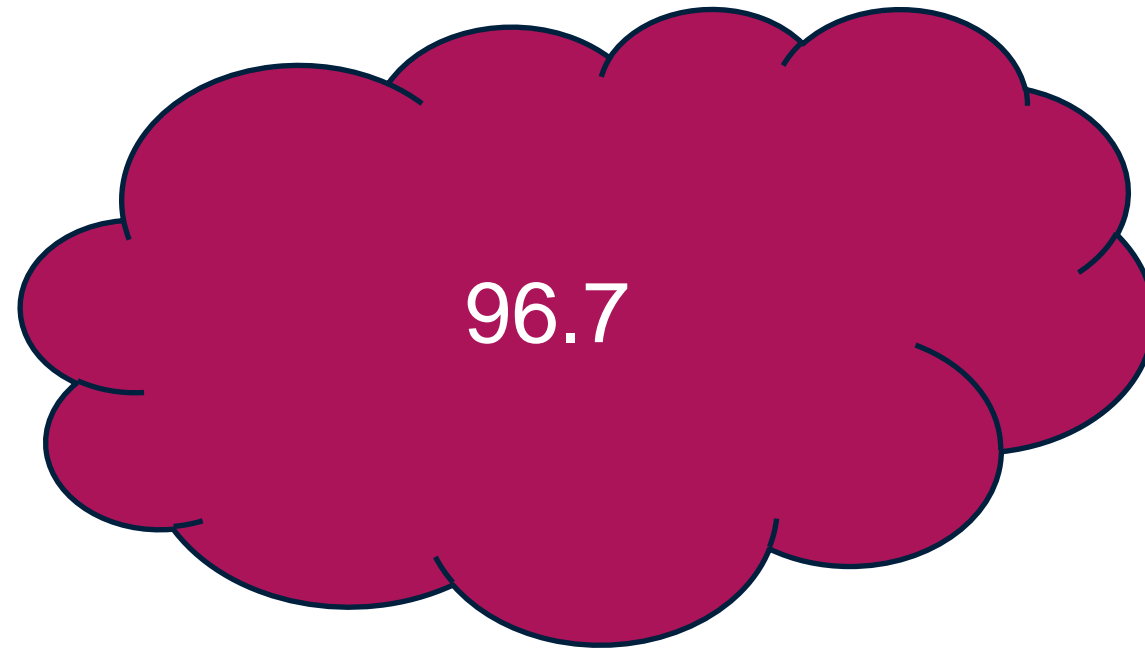
- (meta)data are retrievable by their identifier using a standardized communications protocol
 - The protocol is open, free, and universally implementable
 - The protocol allows for an authentication and authorization procedure, where necessary
- Metadata are accessible, even when the data are no longer available

– Interoperability

- (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation
- (meta)data use vocabularies that follow FAIR principles
- (meta)data include qualified references to other (meta)data

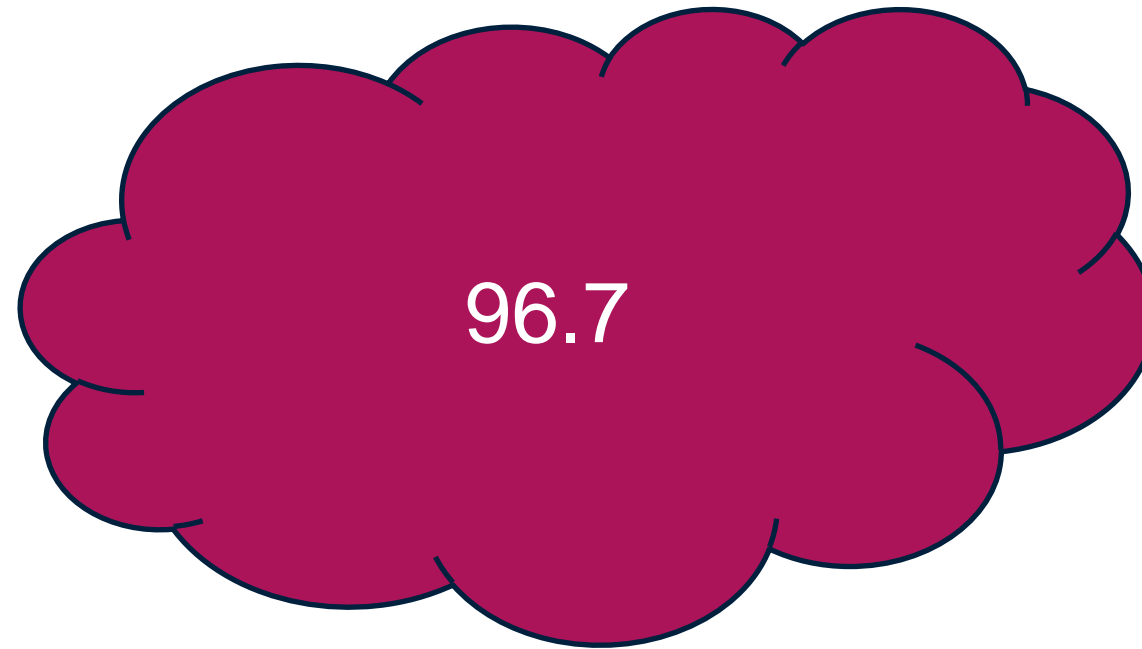
– Reusability

- (meta)data are richly described with a plurality of accurate and relevant attributes
 - (meta)data are realized with a clear and accessible data usage license
 - (meta)data are associated with detailed provenance
 - (meta)data meet domain-relevant community standards



What is metadata?

Unit:
Degree Celsius

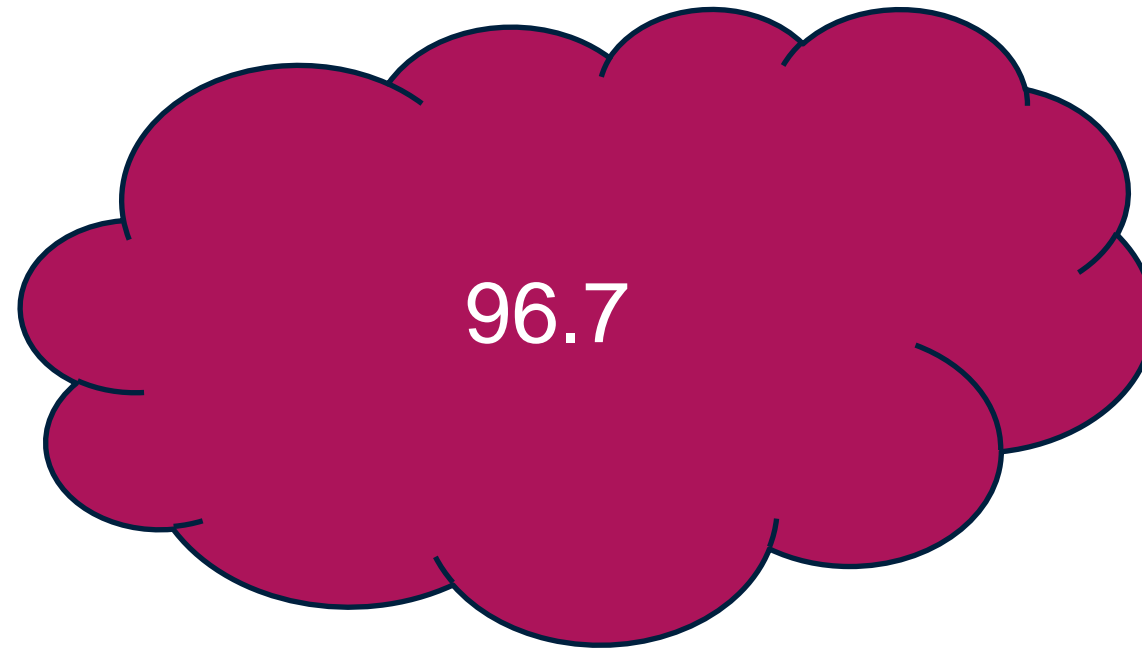


Date:
26 April 1986

Time:
01:23

What is metadata?

Unit:
Degree Celsius



Date:
26 April 1986

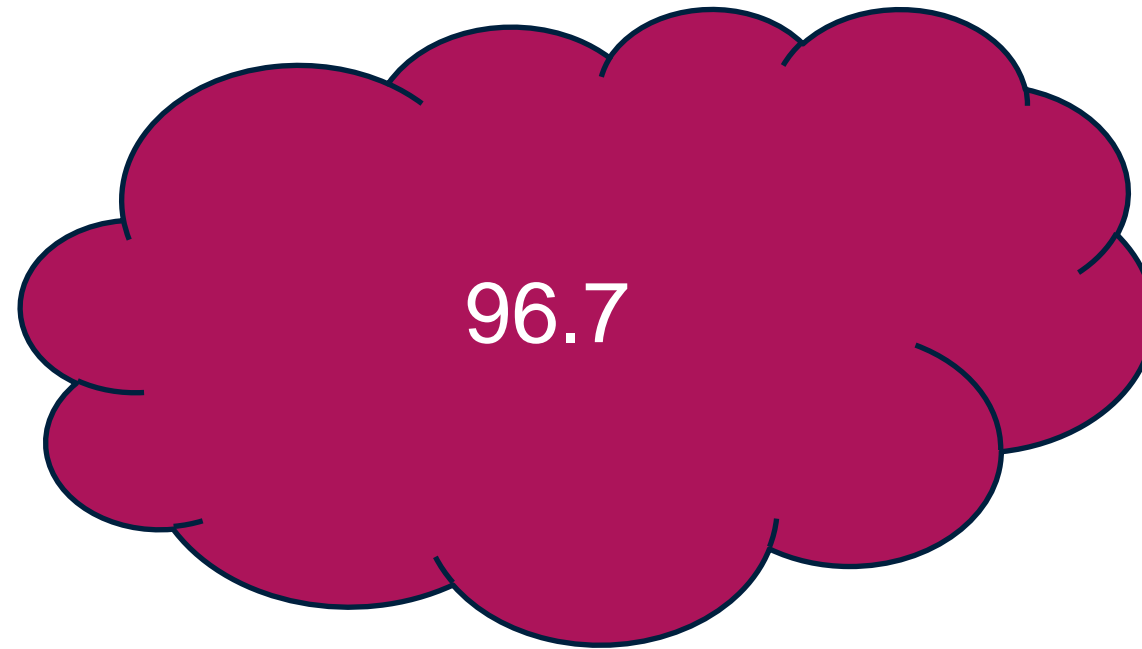
Time:
01:23

Sensor:
CPU Temperature

What is metadata?

Unit:
Degree Celsius

Location:
Chernobyl
Reactor Block 4



Date:
26 April 1986

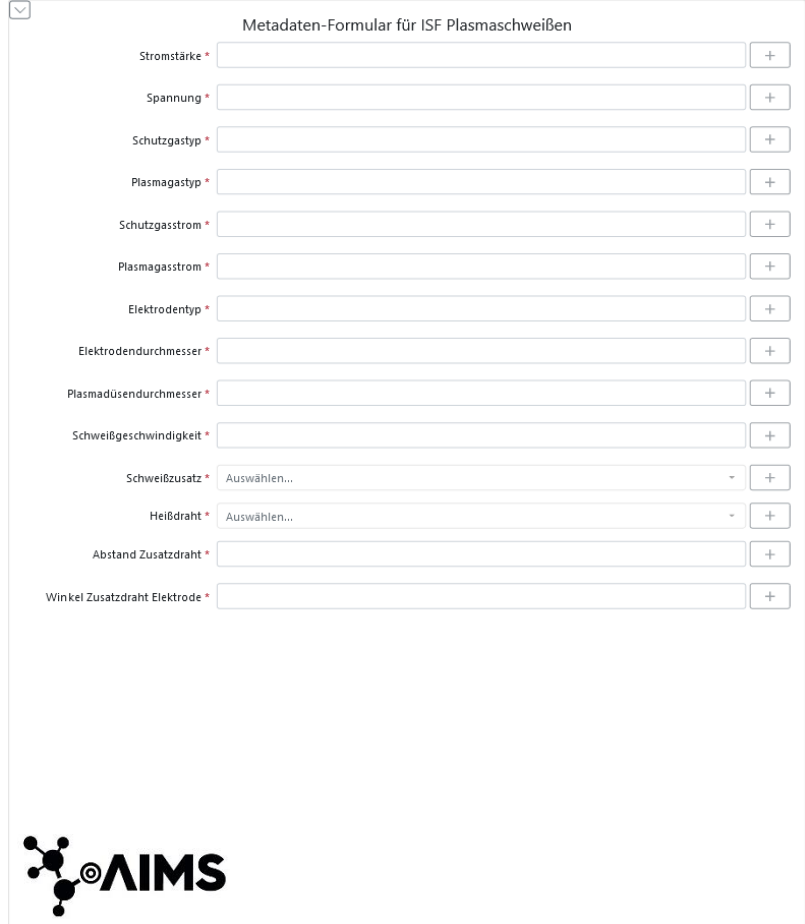
Time:
01:23

Sensor:
CPU Temperature

Metadata gives meaning to data!




- Every field of research is different
- Application profiles must be customizable
- Researchers and institutes must be able to create their own profiles
 - If the application profile does not fit the workflow, it will not be used correctly
 - If the application profile is not relevant, it will not be adapted
- Using the application profiles must be easy
 - Entering metadata for ~50 different fields by hand is tedious
 - Support of automation
- Reduce the amount of work for the researchers

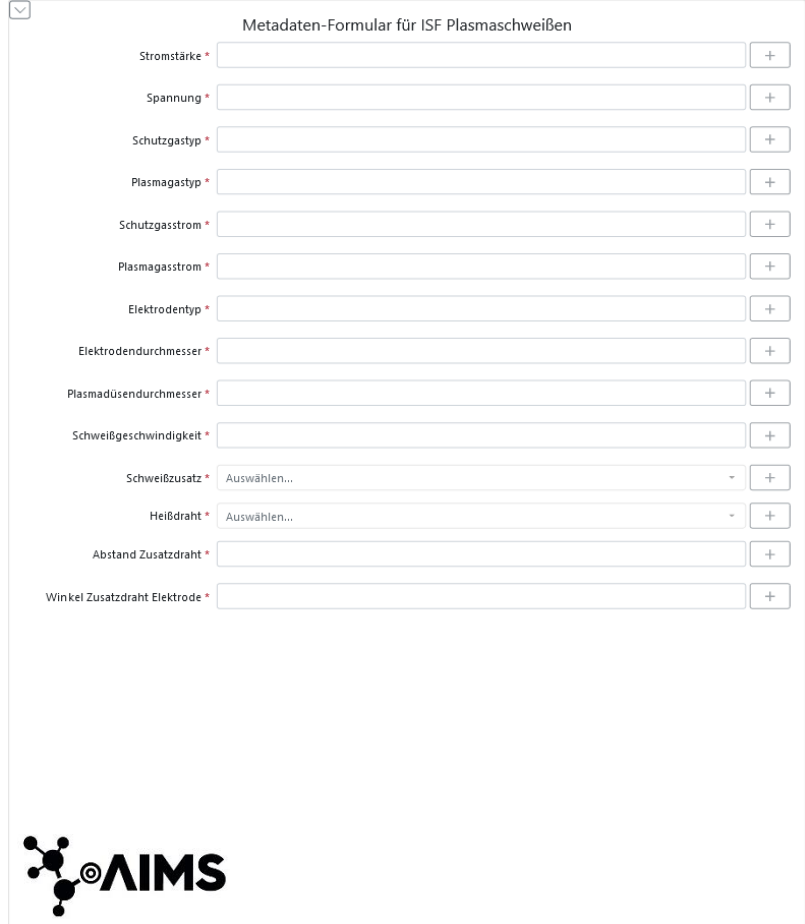


Metadaten-Formular für ISF Plasmaschweißen

Stromstärke *	<input type="text"/>	+
Spannung *	<input type="text"/>	+
Schutzgastyp *	<input type="text"/>	+
Plasmagastyp *	<input type="text"/>	+
Schutzgasstrom *	<input type="text"/>	+
Plasmagasstrom *	<input type="text"/>	+
Elektrodentyp *	<input type="text"/>	+
Elektrorendurchmesser *	<input type="text"/>	+
Plasmadüsendurchmesser *	<input type="text"/>	+
Schweißgeschwindigkeit *	<input type="text"/>	+
Schweißzusatz *	Auswählen...	- +
Heißdraht *	Auswählen...	- +
Abstand Zusatzdraht *	<input type="text"/>	+
Winkel Zusatzdraht Elektrode *	<input type="text"/>	+


 AIMS

- Create a definition of the required metadata
 - Usually called **Application Profile** or **Metadata schema**
- General description of the metadata fields
- Additional information on
 - Expected data types
 - Length of fields
 - Ranges of input fields
- Usage of vocabularies
 - Closed list of possible values (e.g., DFG categories)
- Multitude of input fields
 - For example, each paper must have at least one author but can have multiple authors



Metadaten-Formular für ISF Plasmaschweißen


Stromstärke *	<input type="text"/>	+
Spannung *	<input type="text"/>	+
Schutzgastyp *	<input type="text"/>	+
Plasmagastyp *	<input type="text"/>	+
Schutzgasstrom *	<input type="text"/>	+
Plasmagasstrom *	<input type="text"/>	+
Elektrodentyp *	<input type="text"/>	+
Elektrorendurchmesser *	<input type="text"/>	+
Plasmadüsendurchmesser *	<input type="text"/>	+
Schweißgeschwindigkeit *	<input type="text"/>	+
Schweißzusatz *	Auswählen...	- +
Heißdraht *	Auswählen...	- +
Abstand Zusatzdraht *	<input type="text"/>	+
Winkel Zusatzdraht Elektrode *	<input type="text"/>	+



- Metadata can be validated with the application profiles
 - Complete information
 - Quality standard
- Well defined metadata
 - Interoperability
 - Allows reuse and comparison
- Extensive set of metadata for research data and processes
 - Analyze
 - Share
 - Reuse

Metadaten-Formular für ISF Plasmaschweißen

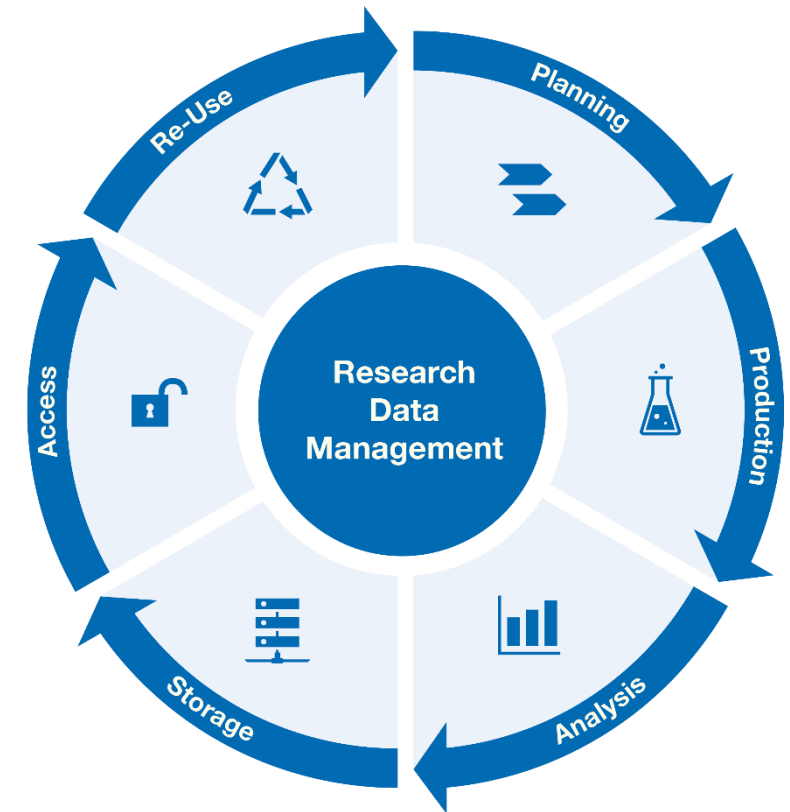
Stromstärke *	<input type="text"/>	+
Spannung *	<input type="text"/>	+
Schutzgastyp *	<input type="text"/>	+
Plasmagastyp *	<input type="text"/>	+
Schutzgasstrom *	<input type="text"/>	+
Plasmagasstrom *	<input type="text"/>	+
Elektrodentyp *	<input type="text"/>	+
Elektrorendurchmesser *	<input type="text"/>	+
Plasmadüsendurchmesser *	<input type="text"/>	+
Schweißgeschwindigkeit *	<input type="text"/>	+
Schweißzusatz *	Auswählen...	- +
Heißdraht *	Auswählen...	- +
Abstand Zusatzdraht *	<input type="text"/>	+
Winkel Zusatzdraht Elektrode *	<input type="text"/>	+



- Short definition:
 - “A data management plan (DMP) structures the handling of research data of a scientific project. It describes how the data is handled during the term and after the end of the project.”

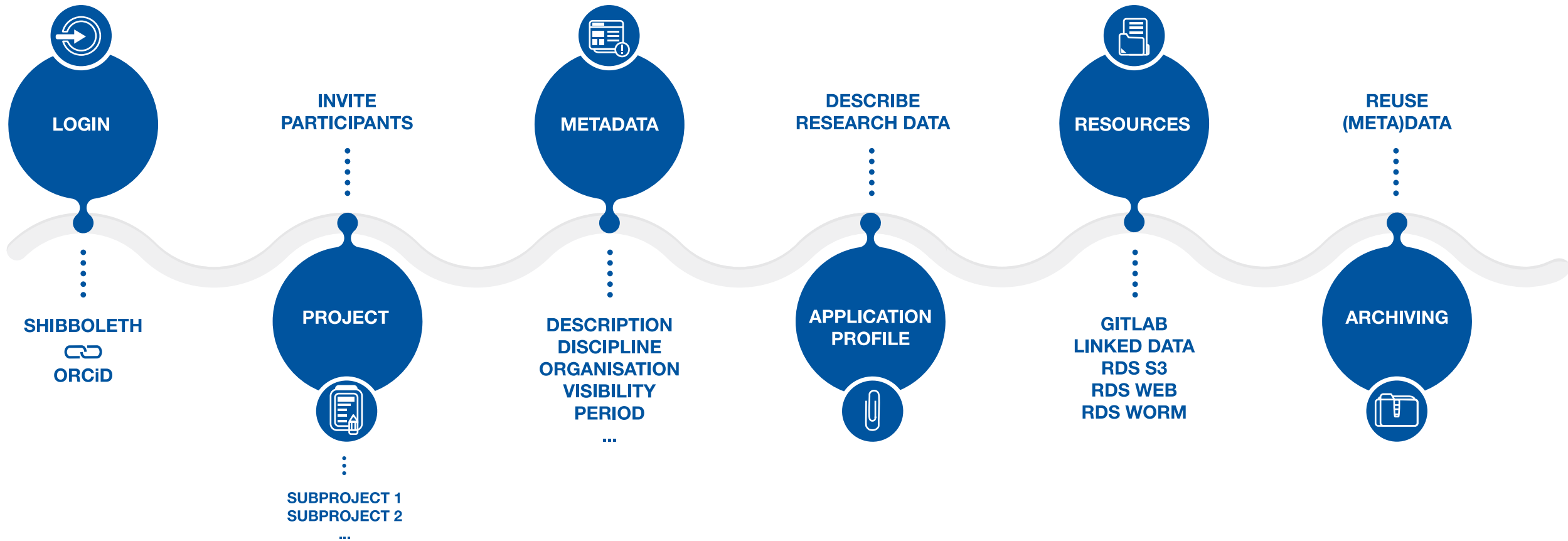
Source: <https://forschungsdaten.info/themen/informieren-und-planen/datenmanagementplan/> 2023-09-17.

- A document which...
 - ...describes the handling of data from the collection to archiving, as well measures to ensure availability and reusability
 - ...is prepared before the project begins and is updated regularly („living document“)
 - ...describes the whole data life cycle of research data



- Problems for Researchers:
 - Limited knowledge about RDM and FAIR principles
- Additional work before publishing
 - Metadata not complete
 - Repeat experiments
 - Ensure Reuse and Interoperability
- Solution:
 - Integration of RDM in daily workflow
 - Use a platform which helps to follow FAIR principles





- Regular talks on RDM topics
- General tools
 - Coscine
 - RDMO / DMP
 - Electronic Lab books
- Domain specific tools
 - Nomad
 - ...

Thank you for
your attention!