**RUHR-UNIVERSITÄT** BOCHUM

# CREATING TRUST IN RESEARCH DATA REPOSITORIES

Case study of an institutional research data repository at Ruhr University Bochum
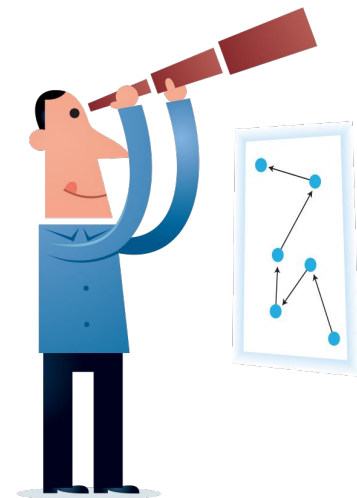
# Talk Outline

**Part I: Background**

- Research Data Management in Germany

**Part II: Who we are**

- Ruhr University Bochum
- Research Use Case from Neuroscience

**Part III: Our Approach**

- From requirements to implementation
- Collaboration with Use Case
- Current status and future aims

Digitalbevaring.dk
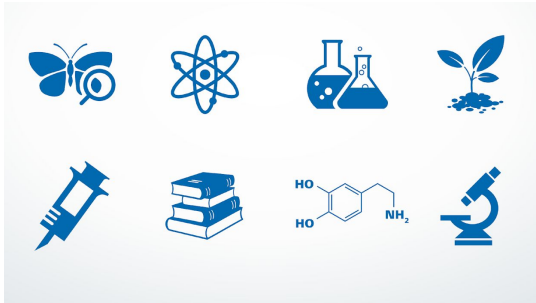
RUHR
UNIVERSITÄT
BOCHUM

RUB

# Part I

Background

# Research Data...

- are all data produced or used by research

- exist in (almost) all disciplines

- are created and processed by a variety of (subject-specific) tools, methods and devices

- usually lack a uniform description

- are very often produced in large quantities (file sizes and number of files)

- are unique and of high value in many cases (astronomy, medicine, history, …)

RUHR
UNIVERSITÄT
BOCHUM

RUB

# Research Data Repositories



Discipline-specific repositories are available,
**however the are not accepted as expected**



re3data.org lists over 3000 research data
repositories from around the world

Creating Trust in Research Data Repositories | 05.06.2024 Open Repositories, Göteborg

RUHR
UNIVERSITÄT
BOCHUM

RUB

# RDM: Challenges

**Legal**

- Sensitive data (psychology, medicine, …)

- Intellectual property (esp. when project partners from industry are involved)

**Practical**

- Lack of standards for data handling

- Unpropriate workflows in existing repositories (complex data models, peer-review, large data sets)

- Lack of time (→ high pressure in scientific careers / „up or out")

**Cultural:**

- Fear of sharing or publishing data

- Lack of trust in service providers

Digitalbevaring.dk

**However:** Research data management is key to sustainable research, reducing effort of future research and enabling innovative data analysis

RUHR
UNIVERSITÄT
BOCHUM

RUB

# Funding for RDM

Main funder: German Research Foundation (DFG)

- Revision of „Guidelines for Safeguarding Good Scientific Practice" in 2019:

  **RDM is mandatory**

- **RDM infrastructure must be provided by host institution**

However, **incentives for standardization are given**, especially in
Collaborative Research Centers (CRCs)

- CRCs are one of the larges funding lines of DFG (~ 24% of the total DFG budget (€848 million) for 268 CRCs)
- Information Infrastructure projects (INF) within CRCs fund staff capacities

  for training, consulting, policy development regarding research data

# Part II

Who we are

# Ruhr University Bochum (RUB)

*„ Deep in the west, where the sun gathers dust. It's better. Much better than you think."*
Herbert Grönemeyer

**Location:**

- Near the river Ruhr: Former industrial center → most dense populated area in Germany (5 million people in the area)

- Near (within a 30-min ride) to two other full universities collaborating in University Alliance Ruhr

- In the German state of North-Rhine Westphalia: One of the two German states with the largest higher education sector in German (about 40 universities, thereof 12 full universities)

**Facts and Figures:**

- One of the 10 largest universities in Germany (40,000 students, 6000 employees, 500 professors, 21 faculties)

- Strong research:
  - €188 million annual third-party funding (25% of total budget)

Hamburg

Berlin

**Bochum**

Munich

Creating Trust in Research Data Repositories | 05.06.2024 Open Repositories, Göteborg

RUHR
UNIVERSITÄT
BOCHUM

RUB

# Research data management at RUB

**Staff capacities:**

- Research data services (library and IT)

- Started with 5 full-time equivalents of staff capacities

- Today: about 10 full-time equivalents of staff capacities



**Storage for research data:**

Two S3 object storage infrastructures in operation:

1. within a consortium of 5 universities (2 PB)

2. within the 3 universities of University Alliance Ruhr (300 TB)

**Repository development**

- 2020: €360,000 for repository implementation

- 2023: Money for 1 developer position for 6 years

Creating Trust in Research Data Repositories | 05.06.2024 Open Repositories, Göteborg

RUHR
UNIVERSITÄT
BOCHUM

RUB

# CRC 1280: Research Use Case



- Topic: "Extinction Learning" (Neuroscience)

- 81 researchers in 17 projects at 4 institutions

- Scientific disciplines: biology, psychology, medicine, and computational neuroscience

- Techniques: microscopy, single cell recording, magnetic resonance imaging, questionnaires

- Human and animal subjects → sensitive data

- Large existing data sets (32 TB, 24 million files in 2 million folders) → ingest strategy

**Common data model** applying an inheritance strategy across folder structures

RUHR
UNIVERSITÄT
BOCHUM

RUB

# Part III

Our Approach

# Requirements

- 2 data models: generic (university wide), neuroscience CRC (discipline-specific)

- Differentiated visibility of data
  → roles & permissions

- Complex (three-step) review workflows

- Data curation steps: Draft, archiving, publication, tombstoning

- Automated data import of hierarchical metadata

- Use of local S3 storage for data and metadata

- Login for project partners via ORCID

RUHR
UNIVERSITÄT
BOCHUM

RUB

# Market analysis

**Starting point:** Feature analysis of available open source repository platforms

→ modifications always necessary

**Aim:** Flexible toolbox for implementation of innovative workflows

**Challenge:** No own expertise to implement modifications

**Solution**: Tendering process for external service provider supporting

- Specification of technical requirements

- Choice of platform

- Implementation

- Roll-out, testing and training concepts

Digitalbeva

# Participation of Use Case

**Before ReSeeD implementation:**

- Cooperative development of data model and mapping to Datacite and Dublin core

- Definition of requirements for ReSeeD

**During ReSeeD implementation:**

- Participation of use case in the project team (selection of and communication with service provider, coordination of beta tests)

- Preparation of data ingest into ReSeeD by storage of (meta)data in a prescribed way

- Provision of actual (not sensitive) data for testing during development

Digitalbevaring.dk

RUHR
UNIVERSITÄT
BOCHUM

RUB

# Implementation



**Current status:**

- **Generic data model (university-wide):** Implementation finished, user acceptance tests running

- **CRC data model:** Beta tests of bulk ingest

**Lessons learned:**

- Implementation of CRC data model and workflows ($\rightarrow$ roles and rights) more challenging than expected

- Building up expertise on the operation of ReSeeD is a challenge

- Technical capabilities still not fully meet expectations of researchers (regarding flexibility and performance)



Digitalbevaring.dk

RUHR
UNIVERSITÄT
BOCHUM

RUB

# ReSeed: Next steps

- Launch publication service
- Establish supporting measures: Training, user tutorials, user survey

**Ongoing implementation and testing**

- Ingest of data records with multiple TB using bulkrax
- Usability testing

**Future aims**

- Upgrade Hyrax 3.5 → 5
- UI re-design with input from usability testing
- Specific data models for further use cases
- Data ingest API
- Flexible metadata
- German language support

Digitalbevaring.dk

RUHR
UNIVERSITÄT
BOCHUM

RUB

# Conclusion

**Challenge for acceptance of RDM infrastructures**

1.  (sensitive) data are not (yet) intended for publication or sharing among all repository users

2.  (innovative) requirements from (interdisciplinary) projects are not represented by repository workflows

3.  researchers are used to a fast evolution of RDM tools and do not trust in institutional processes

**ReSeeD fosters acceptance by**

1.  connecting to local storage infrastructure (fulfilling requirements for sensitive data storage) + roles and permissions

2.  providing individual workflows for interdisciplinary research projects

3.  generating trust via a participation of researchers in the development

Digitalbevaring.dk

Creating Trust in Research Data Repositories | 05.06.2024 Open Repositories, Göteborg

**RUHR UNIVERSITÄT BOCHUM**

**RU**B

# Thank you very much for your attention!

Nina Winter, PhD
Research Data Services @ RUB

Creating Trust in Research Data Repositories | 05.06.2024 Open Repositories, Göteborg