# Repositories for Research Data and Trusted Research Environments

## Martin Weise

with adapted slides from Andreas Rauber

Data Science Group
Technical University of Vienna

TECHNISCHE UNIVERSITÄT WIEN
Vienna University of Technology

Open Science and Research Data Management Innovative and Distributed Training Programme

SAPIENZA
UNIVERSITÀ DI ROMA

# Motivation

Data is the new oil
... or "the new water"
... or "the new light"

Martin Weise, Technical University of Vienna

# Motivation

Data is the new oil ... if properly managed!
... or "the new water"
... or "the new light"

Otherwise, it's an oilspill
... or flood
... or blinding flash of lighting



v Brussels, Theresa May, 2017. The World's Most Valuable Resource, in *The Economist*. Edition May 6th.
Martin Weise, Technical University of Vienna

# Motivation

Research depends on data in virtually all disciplines:

- **Value** of data determinded through
  - Exhaustive collection
  - Data (pre-)processing
  - Volume, e.g. meta-studies
  - Reproducability as core principle of science
- Proper **data management** enables
  - Speedup of research (avoiding repeated collection, processing)
  - Robust resarch (larger data pools)
  - Increased quality (reproducability, comparability)

# Agenda

1. Introduction
2. Background
3. Repositories for Research Data
   - DSpace
   - Gitlab
   - InvenioRDM
   - DBRepo
4. Trusted Research Environments
   - RemoteNEPS
   - SAIL Gateway
   - DEXHELPP
   - OSSDIP
5. Future Work
6. Conclusion

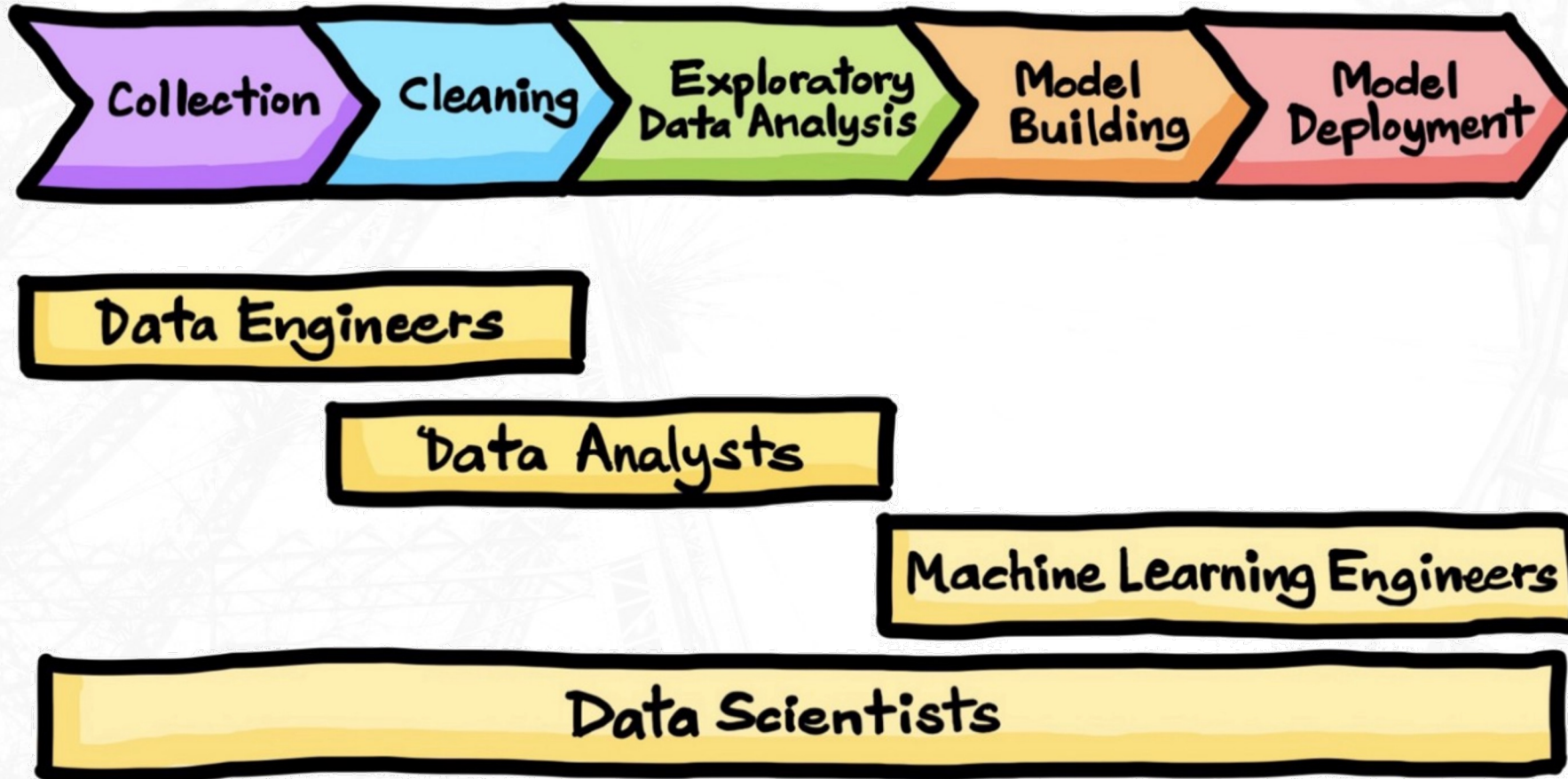Martin Weise, Technical University of Vienna

# 1. Introduction

PreDoc Researcher in the first year

- **MSc** in *Software Engineering & Internet Computing*
  2022 (Technical University of Vienna)

- **BSc** in *Software & Informaiton Engineering*
  2019 (Technical University of Vienna)
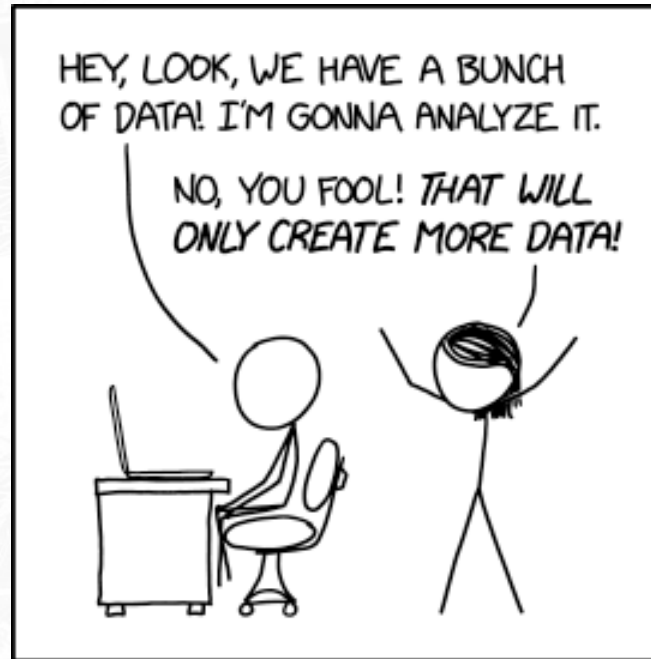
Research interests:

- Making sensitive data FAIR

Martin Weise, Technical University of Vienna

# 2. Background

Martin Weise, Technical University of Vienna

**Cleaning and preparing** data takes about 80% of the total engineering effort

- Real-world **data** may be
  - *incomplete*, lacking attribute values, contains only aggregated data,
  - *noisy,* containing errors or outliers,
  - inconsistent, discrepancy in names
- **Preparation** generates a subset of the data, potentially increasing utility
  - *attribute selection,* relevant data, anomaly removal, duplicate elimination
  - *reducing data*, sampling or instance selection
- **Outcome**
  - *recovery* of incomplete data
  - *purify data*, correcting errors, removing outliers
  - *resolve conflicts* using domain knowledge, expert decisions

Zhang, S., Zhang, C. & Yang, Q., 2003. Data preparation for data mining, in *Applied Artificial Intelligence*, 17(5-6), p.375-381, DOI: 10.1080/713827180
Martin Weise, Technical University of Vienna

After all this effort, people will be happy to re-use data, **right**?

ALT: "It's important to make sure your analysis destroys as much information as it produces."
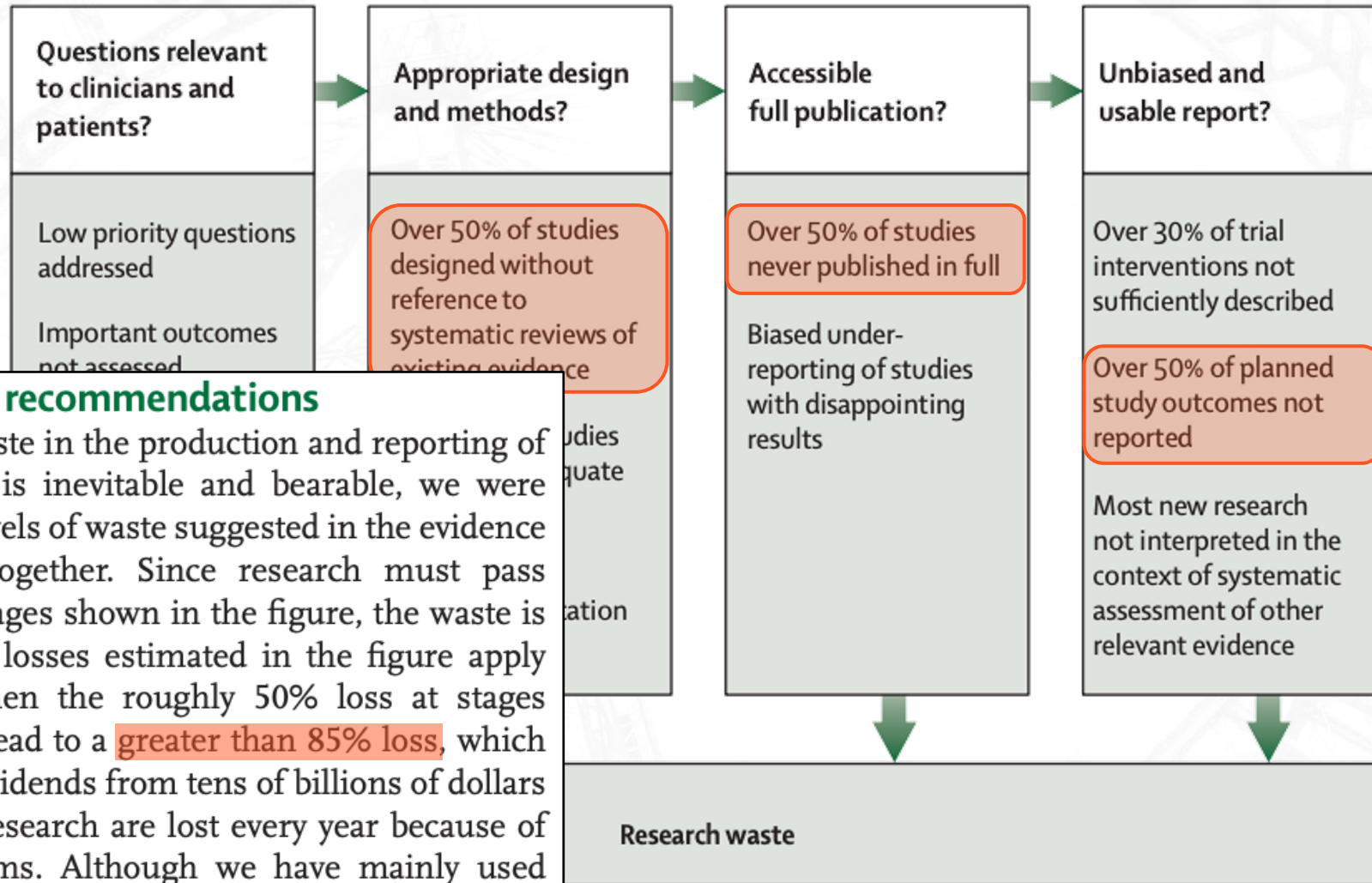
Martin Weise, Technical University of Vienna

# 2. Background



| Questions relevant to clinicians and patients? | Appropriate design and methods? | Accessible full publication? | Unbiased and usable report? |
|---|---|---|---|
| Low priority questions addressed

Important outcomes not assessed

Clinicians and patients not involved in setting research agendas | Over 50% of studies designed without reference to systematic reviews of existing evidence

Over 50% of studies fail to take adequate steps to reduce biases—eg, unconcealed treatment allocation | Over 50% of studies never published in full

Biased under-reporting of studies with disappointing results | Over 30% of trial interventions not sufficiently described

Over 50% of planned study outcomes not reported

Most new research not interpreted in the context of systematic assessment of other relevant evidence |

**Research waste**

Martin Weise, Technical University of Vienna

**Questions relevant to clinicians and patients?**

Low priority questions addressed

Important outcomes not assessed

**Appropriate design and methods?**

Over 50% of studies designed without reference to systematic reviews of existing evidence

**Accessible full publication?**

Over 50% of studies never published in full

Biased under-reporting of studies with disappointing results

**Unbiased and usable report?**

Over 30% of trial interventions not sufficiently described

Over 50% of planned study outcomes not reported

Most new research not interpreted in the context of systematic assessment of other relevant evidence

Research waste

**Conclusions and recommendations**

Although some waste in the production and reporting of research evidence is inevitable and bearable, we were surprised by the levels of waste suggested in the evidence we have pieced together. Since research must pass through all four stages shown in the figure, the waste is cumulative. If the losses estimated in the figure apply more generally, then the roughly 50% loss at stages 2, 3, and 4 would lead to a greater than 85% loss, which implies that the dividends from tens of billions of dollars of investment in research are lost every year because of correctable problems. Although we have mainly used

Martin Weise, Technical University of Vienna

Why even share research data?

- **Increase trust** in the work, allow reproduce and validate findings

- Information is **valuable** to the research community

- **Contribute** work beyond the original findings

- Allow others to **re-use** and build on top of their data

Martin Weise, Technical University of Vienna

## 2. Background

Why even share research data?

- **Increase trust** in the work, allow reproduce and validate findings

- Information is **valuable** to the research community

- **Contribute** work beyond the original findings

- Allow others to **re-use** and build on top of their data

"If I have seen further it is by standing on the shoulders of Giants."

— Isaac Newton, 1675.

We contiously use the understanding gained by major thinkers in order to make intellectual progress.

Martin Weise, Technical University of Vienna

Is this necessary, why not just **reasonably request** data from researchers?

Martin Weise, Technical University of Vienna

Is this necessary, why not just **reasonably request** data from researchers?

**DATA-SHARING BEHAVIOUR**

Of almost 1,800 manuscripts for which the authors stated they were willing to share their data, more than 90% of corresponding authors either declined or did not respond to requests for data. Only about 7% of authors actually handed over data.

Manuscripts with statement indicating data available on request — **1,792** manuscripts

Authors who did not respond or declined requests for data — **1,670**

Authors who provided usable data — **120**

Percentage (0, 20, 40, 60, 80, 100)

©nature

Martin Weise, Technical University of Vienna

Is this necessary, why not just **reasonably request** data from researchers?



Gabelica, M., Bojčić, R. & Puljakc, L., 2022. Many Researchers were not Compliant with their Published Data Sharing Statement: a Mixed-methods Study, in *Journal of Clinical Epidemiology*, DOI: 10.1016/j.jclinepi.2022.05.019

Martin Weise, Technical University of Vienna

## Common Scenario



Martin Weise, Technical University of Vienna

## Common Scenario



Apple Cloud

Overleaf

Microsoft Cloud

Data Ownership?
Availability?
Support?
FAIRness?
Sensitivity?

**Researcher**

Google Cloud

Institute Cloud

Institute E-Mails

What are **research data** repositories?



Martin Weise, Technical University of Vienna

Data Management

- Currently largely on researcher's shoulders
- Need **separation of concerns** (excerpt)
  - *Researchers*, **work with data**, domain expertise
  - *Data Stewards*, **curation**, preservation, FAIR
  - *IT-Department*, hardware/software infrastructure, security, backup
  - *Legal-Department*, licenses, NDAs
  - *Admin*, reporting, GDPR, complicance
- **Dedicated infrastructure** to ensure data is **properly managed** and value realized
- Infrastructure for **research data management** connected to internal information systems, funders, etc.

Research Data Management is a joint effort!

# 3. Repositories for Research Data

Steps towards research infrastructure at TU Vienna:

- Involved stakeholders, regular round tables (rectorate, support offices, research departments)
- Established **policy** on RDM
- Established **Center for Research Data Management**
- Devised plan for
  - National and European projects
    - **FAIR Data Austria**
    - Austrian Data Labs and Services
    - EOSC-* projects (European Open Science Cloud)
    - Involvement in **Research Data Alliance** (RDA), **EGI**, …
  - Setting up and rolling out infrastructure

Martin Weise, Technical University of Vienna

# 3. Repositories for Research Data



Research Data. [Online]. URL: https://www.tuwien.at/forschung/fti-support/forschungsdaten, accessed 2022-09-09

Martin Weise, Technical University of Vienna

# 3. Repositories for Research Data
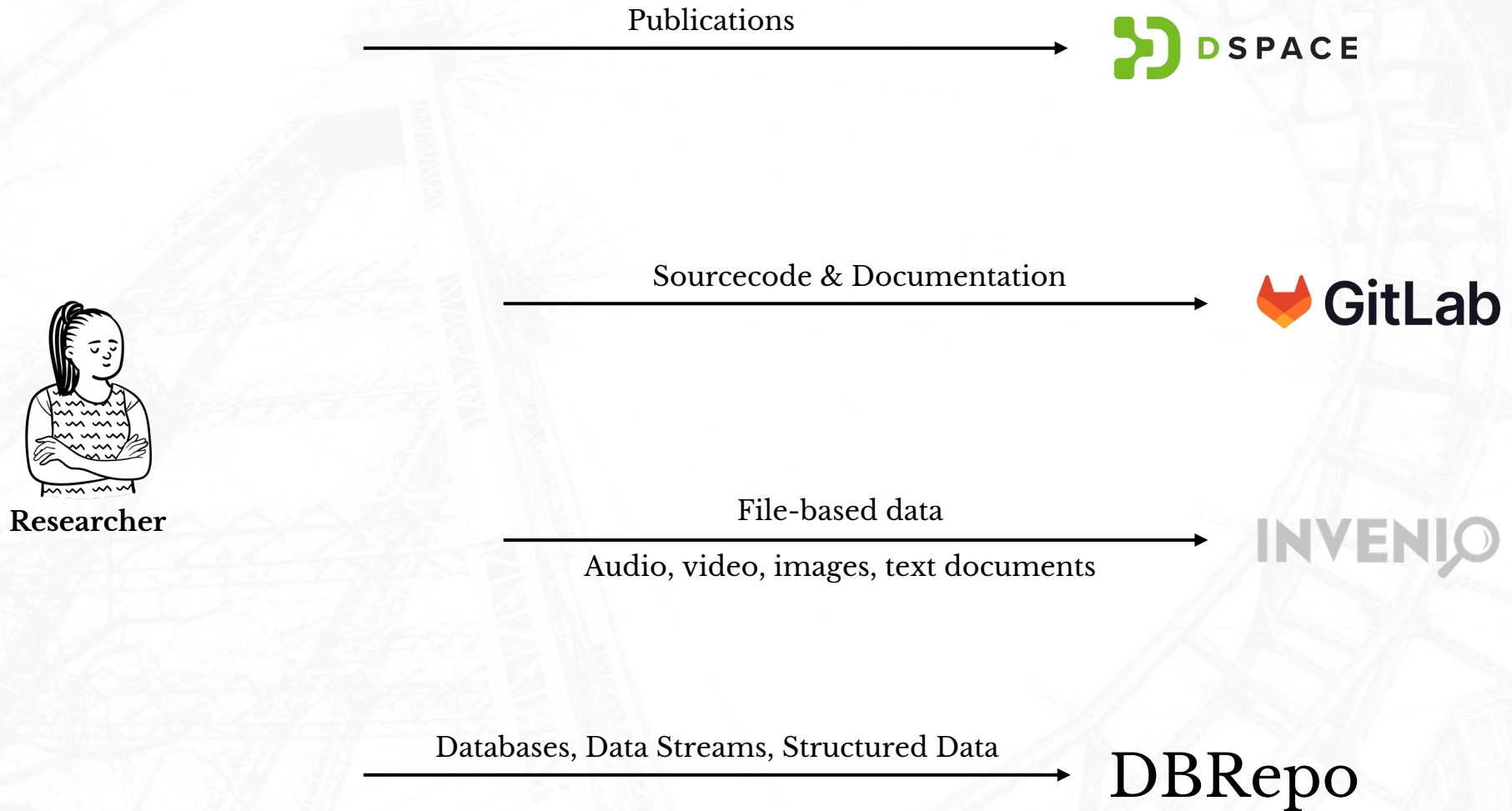
Martin Weise, Technical University of Vienna

# 3. Repositories for Research Data

Repository Infrastructure

1. Repositories for **all kinds** of research material
   - Input, output, interim
   - Open and closed / sensitive data
2. Provide **visibility**
   - Citation, impact
   - FAIR compliant
3. Be largely **transparent** to researchers
   - Integration with TU and external infrastructure
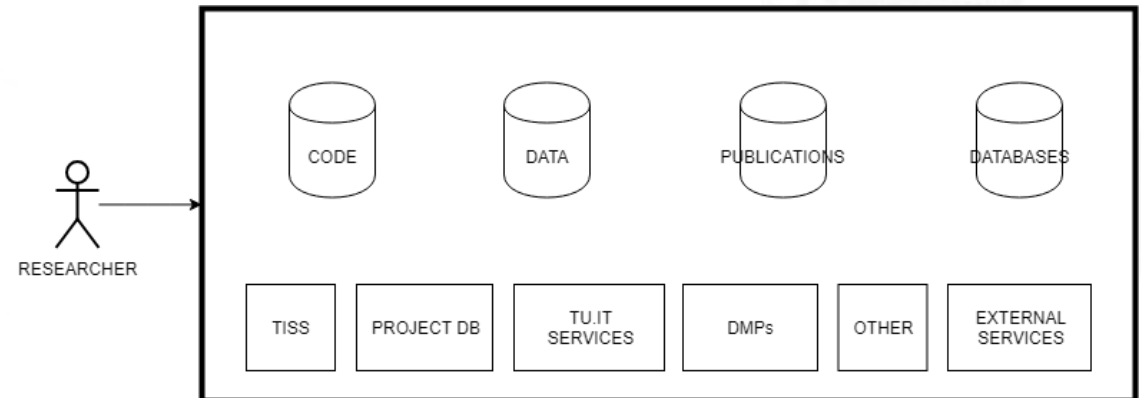4. Be **trustworthy**

Martin Weise, Technical University of Vienna

## Architecture and vision

Publications ⟶ **DSPACE**

Sourcecode & Documentation ⟶ **GitLab**

**Researcher**

File-based data
Audio, video, images, text documents ⟶ **INVENIO**

Databases, Data Streams, Structured Data ⟶ DBRepo

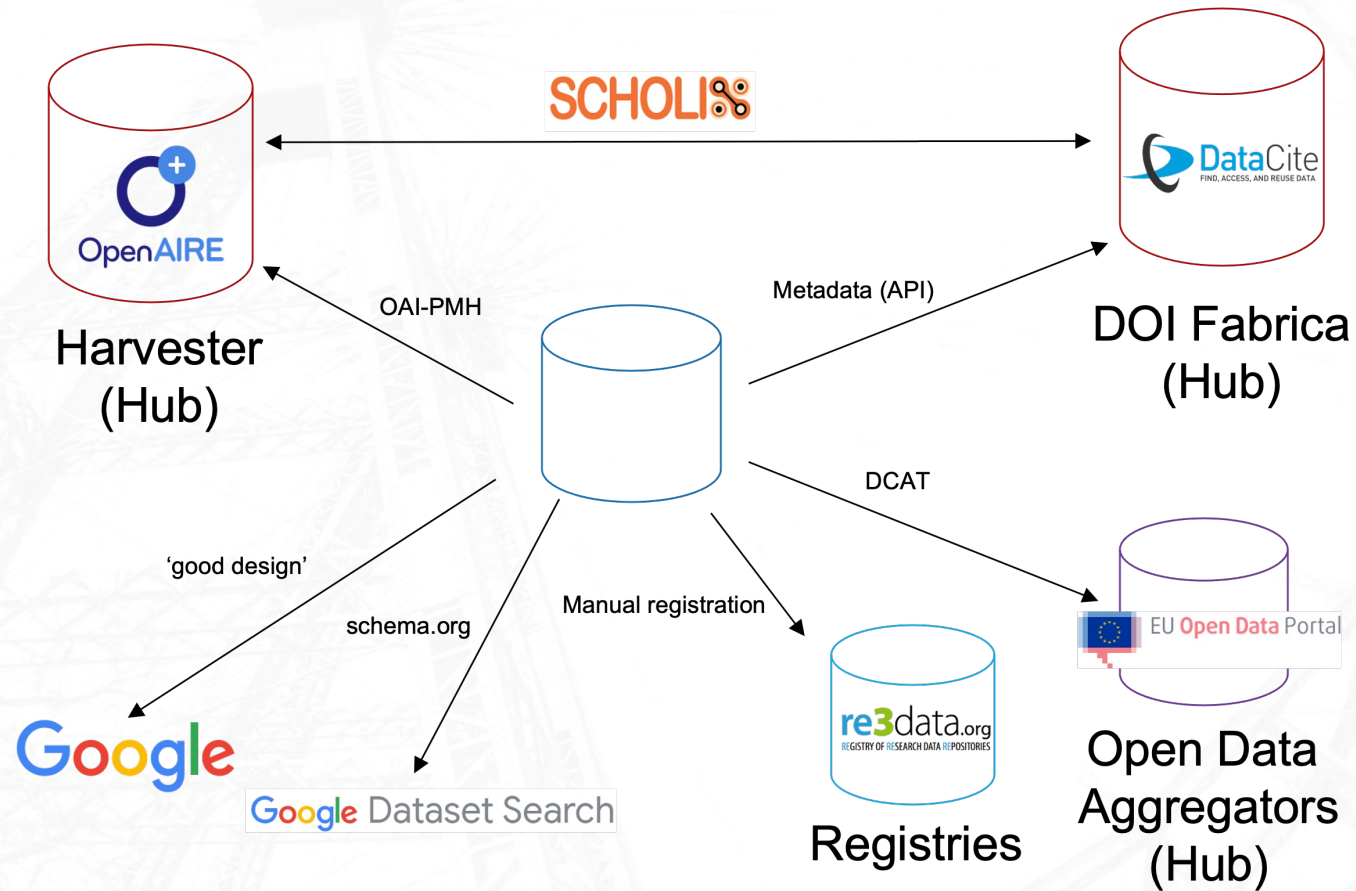Martin Weise, Technical University of Vienna

Repository Infrastructure

- Provide unified services for research data management
- Cooperation
  - Center for Research Data Management
  - TU.it
  - Campus Development Services
  - TU library
  - Research Unit Data Science (formerly IFS)
  - Central administration, legal department



Systems integrated appear as a 'single repository'

Martin Weise, Technical University of Vienna

External Visibility: nobody searches in your repo to find data at first

# 3.1. ReposiTUm (DSpace)

**Query Search**

**Accessability Filter**

**Faceted Browsing**

Martin Weise, Technical University of Vienna

# 3.1. ReposiTUm (DSpace)

Martin Weise, Technical University of Vienna

# 3.1. ReposiTUm (DSpace)



Controlled Vocabulary (=machine-readable)

Groupings

Direct .pdf ressource access

Weise, M, 2021. A QR-Code Optical Covert Channel in an Air-Gapped Secure Data Infrastructure. [Thesis], p.97, DOI: 10.34726/hss.2022.84700

Martin Weise, Technical University of Vienna

# 3.2. Gitlab



Single Sign-On (e.g. eduroam is also SSO)

Martin Weise, Technical University of Vienna

# 3.2. Gitlab



OSSDIP. [Online]. URL: https://gitlab.tuwien.ac.at/martin.weise/ossdip, accessed 2022-09-09

Martin Weise, Technical University of Vienna

# 3.3. Research Data Repository (InvenioRDM)

- Make digital objects FAIR
- Suitable for **research data**
- Not for publications
  - Other system exists (ReposiTUm)
- Running since December 2020

More details will be presented later today!

Martin Weise, Technical University of Vienna

# 3.4. DBRepo

- Invenio can handle collection of files well

- How about relational data in databases?
  - Releasing a data dump every x amount of time?
  - Adding continous data streams, e.g. IoT?
  - How to update / correct data in those databases?
  - Allow reproduction of any subset?

- Private cloud-based environment?

- Dump the data after the end of a project into some repository – to fulfill some grant agreement?

Martin Weise, Technical University of Vienna

# 3.4. DBRepo

- Cloud hosted repository for structured research data
- Supports data versioning & citeability via query store and dynamic data citations (Recommendations of RDA WGDC)
- **Microservice architecture**
- Each database encapsulated in a Docker container
- Central **metadata database**
- APIs for different levels of SQL-knowledge:
  - Web interface, support for CSV import,...
  - REST, message queue for data streams

Martin Weise, Technical University of Vienna

# 3.4. DBRepo



Database Name

Database Visibility

Table information

AMQP information

.csv import

Create a subset of the data

Add metadata

Martin Weise, Technical University of Vienna

# 3.4. DBRepo

Query Builder for simple subsets

Export .csv

Export DataCite Metadata



DataCite Metadata →

Weise, M., 2022. Early stage Researchers' Training Week Subset. [Dataset]. URL: https://dbrepo.ossdip.at/pid/55, accessed 2022-09-10

Martin Weise, Technical University of Vienna

# 3.4. DBRepo

Metadata makes databases findable (as in FAIR)

- Container name

- Database name, description, -license

- Table name

  - Column name, -type, -uniqueness, -nullability, -date format

  - Column measurement unit (controlled vocabulary)

  - Statistical properties

Foodvoc. [Online]. URL: http://www.ontology-of-units-of-measure.org/resource/om-2/second-Time, accessed 2022-09-10

Martin Weise, Technical University of Vienna

# 3.4. DBRepo

Persistent identification of **arbitrary subsets of data**

- Each query issued to the database is saved in the Query Store
- Attaching **metadata** to a query statement, following the DataCite schema
- Mirror the query metadata to DBRepo's central database
- Allows precise data citation
- **Implements recommendations** of the RDA WGDC

Martin Weise, Technical University of Vienna

# 3.4. DBRepo

## Simple use-case:



https://luft.umweltbundesamt.at/

AMQP Tuples

DBRepo

HTTP API

HTTP API

https://s125.dl.hpc.tuwien.ac.at/user/retropotato/lab

http://s125.dl.hpc.tuwien.ac.at:8080/grafana/d/
R4B-UWZVz/dbrepo-airquality-austria

Martin Weise, Technical University of Vienna

# 3.4. DBRepo Further Reading

Material

- https://indico.egi.eu/event/5882/contributions/16724/ (EGI'22 Poster)

- https://doi.org/10.5281/zenodo.6637333 (IDCC'22 paper)

- https://doi.org/10.17605/OSF.IO/B7NX5 (iPRES'21 paper)

Resources

- https://dbrepo.ossdip.at (sandbox)

- https://dbrepo-docs.ossdip.at (documentation, getting started guide)

- https://gitlab.phaidra.org/fair-data-austria-db-repository/fda-services (source code)

Martin Weise, Technical University of Vienna

# 4. Trusted Research Environments

"five safes" dimensions:

1. *Safe projects*, **appropriateness** of the usage of the data
2. *Safe people*, **identify** users that access senstive data, legal bindings
3. *Safe data*, **appropriate** data de-identification, research questions formulated
4. *Safe settings*, **necessity** of security and transparency
5. *Safe outputs*, **approved**, aggregated research results can be exported



Safe People

Safe Projects

Safe Settings

Safe Data

Safe Output

Desai, T., 2016. Five Safes: Designing Data Access for Research. [Online]. URL: https://uwe-repository.worktribe.com/output/914745/five-safes-designing-data-access-for-research, accessed 2022-09-10

Martin Weise, Technical University of Vienna

# 4. Trusted Research Environments

UK Health Data
Research Alliance

- United Kingdom Health Data Research Alliance (UKHDRA)
- Confederation of leading organizations in the healthcare field
- Extend "five safes"
    6. Safe return, allow de-identified research results to be re-identified and securely mapped back to the original data set



Ensure public trust by implementing UKHDRA's recommendations on TRE+ and independent accreditation and audit

UK Health Data Research Alliance, & NHSX. (2021). Building Trusted Research Environments - Principles and Best Practices; Towards TRE ecosystems. DOI: 10.5281/zenodo.5767586

Martin Weise, Technical University of Vienna

# 4. Trusted Research Environments

1.  *Safe people*

    – Control & ensure interests of people where data

    ▪ Collected

    ▪ Analyzed

    – Sign legally-binding documents (e.g. NDAs)

    – Analysts must undergo information governance training (once approved, access all)

Martin Weise, Technical University of Vienna

**UK Health Data Research Alliance**

## 2. *Safe projects*

- **Appropriate** use of sensitive data
- Mandatory possibility to external audits
- **Ethics board** (review the project proposal and gives a clearance)
- Want to **improve the maturity** of the project management processes

Martin Weise, Technical University of Vienna

**UK Health Data Research Alliance**

## 3. *Safe setting*

- – Ensure a straightforward to use, secure environment for the sensitive data to reside in
- – Defined and transparent process
- – Trusted administrators with permission to exchange data through the air-gap
- – System-internal barriers
- – Protected individual data cannot be exported
- – Actions of authorized analysts are monitored

Martin Weise, Technical University of Vienna

UK Health Data
Research Alliance

4. *Safe computing*

- Outsourcing infrastructure (e.g. content delivery networks)

- Overcome the risk of exposing sensitive data to public cloud providers

- Additional safeguards that disallow any outsourced hardware or software

- Access the sensitive data at any time must be implemented

- This approach is well-studied and supported from major public cloud providers.



UK Health Data Research Alliance, & NHSX. (2021). Building Trusted Research Environments - Principles and Best Practices; Towards TRE ecosystems. DOI: 10.5281/zenodo.5767586

Martin Weise, Technical University of Vienna

5. *Safe data*

– Minimize risk of (accidental) re-identification of protected individuals when importing data into the TRE

– De-identification software tools and encrypted (virtual) disks

– Prioritizes the interests of the protected individual over the analyst through technical & organizational measures

Martin Weise, Technical University of Vienna

**UK Health Data Research Alliance**

6. *Safe outputs*

– Barrier ("air-gap") of the safe setting and the open internet

– Data Ingress/egress process

■ Manual interaction to control risk of disclosure

– Manages the communication to the outside world

Martin Weise, Technical University of Vienna

## 7. *Safe return*

- De-identified research output to be returned into the TRE where the sensitive data comes from and the identity of the protected individual is known

- Gain information about the protected individual itself

- Profit from allowing analysts to work with the data

- Researchers can only access de-identified sensitive data while all their actions are overseen by a committee



Upon approval of both the patient and the ethics board, the research output can be re-identified and mapped back to the original data set to enrich value to it.

UK Health Data Research Alliance, & NHSX. (2021). Building Trusted Research Environments - Principles and Best Practices; Towards TRE ecosystems. DOI: 10.5281/zenodo.5767586

Martin Weise, Technical University of Vienna

National Education Panel Study (NEPS)

– Remote access to data to **utilize** them better than local access

– Better than conventional methods

  ▪ Remote execution

  ▪ Job-submission systems

  Queues, input-output (heavily) **delayed**

– Hosts a full-fledged secure data infrastructure

– Data access is **moderately** anonymized, the Analyst must sign an additional supplementary agreement

– Data export is possible via link after signing data us agreements (**heavily** anonymized and aggregated)

Skopek, J., 2016. RemoteNEPS - An Innovative Research Environment, in *Methodological Issues of Longitudinal Surveys: The Example of the National Educational Panel Study*, p.611-626. DOI: 10.1007/978-3-658-11994-2_34

Martin Weise, Technical University of Vienna

**NEPS**

## Social Architecture

## Second-factor
(prevent credential sharing)

## SPSS, R, Office, Windows 7

Martin Weise, Technical University of Vienna

## dexhelpp

Facilitates research in Austria for almost 10 years

- Provide **analysts** with a secure and controlled environment **without the need** to exfiltrate data out of the system

- Data owners on the other side **deposit** their data from **heterogeneous** sources in an encrypted vault and specify **fine-grained access** rights, e.g.
  - Entire data assets
  - Just specific subsets

Windows Environment          Linux Environment

Martin Weise, Technical University of Vienna

- Monitoring Node continuously **monitors** the access to the Data Endpoint (allows for **auditing** and inspection of the usage of the data at any time)

- **Docker environment** (e.g. PostgreSQL, RStudio, Web Applications)

- Analysts **working** on the Remote Desktop Node, connected through the VPN Client using two-factors

- Special hardware and hypervisor for GPU execution

Accelerated Environment

Special Hardware

Martin Weise, Technical University of Vienna

# 4.3. OSSDIP

- Sensitive Data (=due privacy issues, commercial interests, …)

- Provide access for analysis, but ensure data is not leaked or misused

- Standard approach
  - Pseudonymization
  - Anonymization
  - $k$-anonymity
  - $l$-diversity
  - $t$-closeness



Martin Weise, Technical University of Vienna

- Sensitive Data (=due privacy issues, commercial interests, ...)

- Provide access for analysis, but ensure data is not leaked or misused

- Standard approach

  – Pseudonymization

  – Anonymization

  – *k*-anonymity

  – *l*-diversity

  – *t*-closeness

**sensitive attribute**

| | ZIP Code | Age | Disease |
|---|---|---|---|
| 1 | 47677 | 29 | Heart Disease |
| 2 | 47602 | 22 | Heart Disease |
| 3 | 47678 | 27 | Heart Disease |
| 4 | 47905 | 43 | Flu |
| 5 | 47909 | 52 | Heart Disease |
| 6 | 47906 | 47 | Cancer |
| 7 | 47605 | 30 | Heart Disease |
| 8 | 47673 | 36 | Cancer |
| 9 | 47607 | 32 | Cancer |

**3-anonymity**

**Alice** knows **Bob** is 27yo, lives in 47678 and is in the first three entries

⊆ **Bob** has a **heart diseaese**

| | ZIP Code | Age | Disease |
|---|---|---|---|
| 1 | 476** | 2* | Heart Disease |
| 2 | 476** | 2* | Heart Disease |
| 3 | 476** | 2* | Heart Disease |
| 4 | 4790* | ≥ 40 | Flu |
| 5 | 4790* | ≥ 40 | Heart Disease |
| 6 | 4790* | ≥ 40 | Cancer |
| 7 | 476** | 3* | Heart Disease |
| 8 | 476** | 3* | Cancer |
| 9 | 476** | 3* | Cancer |

Samarati, P., 2001. Protecting Respondents Identities in Microdata Release, in *IEEE Transactions on Knowledge and Data Engineering, 13(6)*. DOI: *10.1109/69.971193*

Martin Weise, Technical University of Vienna

- Sensitive Data (=due privacy issues, commercial interests, …)

- Provide access for analysis, but ensure data is not leaked or misused

- Standard approach

    - Pseudonymization

    - Anonymization

    - $k$-anonymity

    - $l$-diversity

    - $t$-closeness

**Alice** knows **Bob** is 27yo, lives in 47678 and is in the first three entries

⊆ **Bob** has a **low salary** and some **stomach disease**

sensitive attributes

| | ZIP Code | Age | Salary | Disease |
|---|---|---|---|---|
| 1 | 47677 | 29 | 3K | gastric ulcer |
| 2 | 47602 | 22 | 4K | gastritis |
| 3 | 47678 | 27 | 5K | stomach cancer |
| 4 | 47905 | 43 | 6K | gastritis |
| 5 | 47909 | 52 | 11K | flu |
| 6 | 47906 | 47 | 8K | bronchitis |
| 7 | 47605 | 30 | 7K | bronchitis |
| 8 | 47673 | 36 | 9K | pneumonia |
| 9 | 47607 | 32 | 10K | stomach cancer |

3-diverse

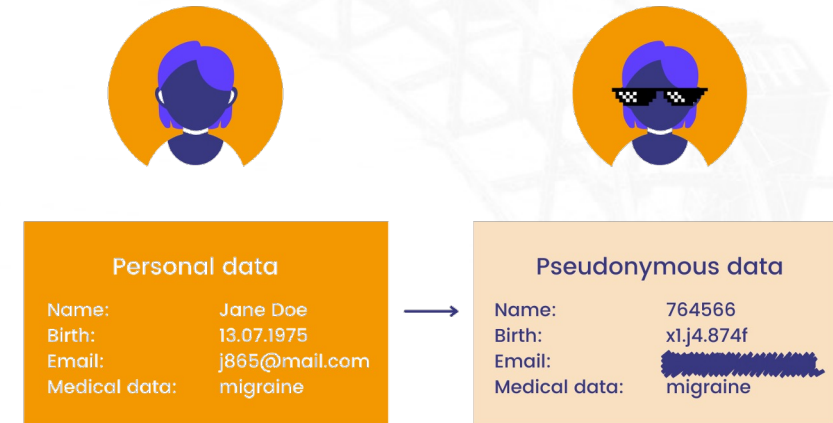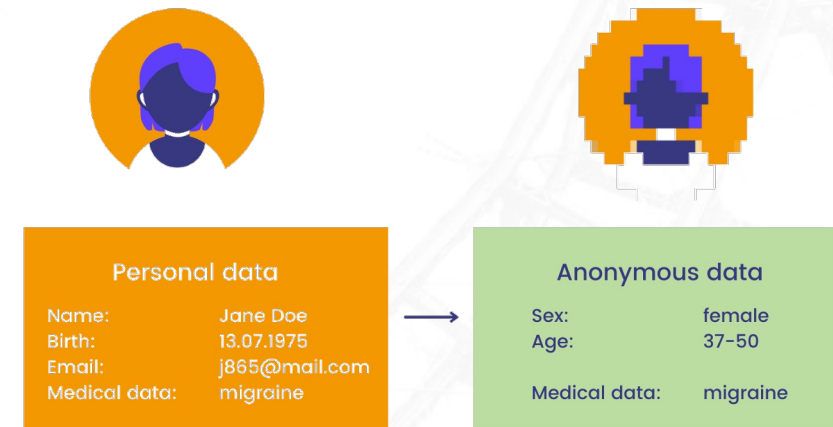| | ZIP Code | Age | Salary | Disease |
|---|---|---|---|---|
| 1 | 476** | 2* | 3K | gastric ulcer |
| 2 | 476** | 2* | 4K | gastritis |
| 3 | 476** | 2* | 5K | stomach cancer |
| 4 | 4790* | ≥ 40 | 6K | gastritis |
| 5 | 4790* | ≥ 40 | 11K | flu |
| 6 | 4790* | ≥ 40 | 8K | bronchitis |
| 7 | 476** | 3* | 7K | bronchitis |
| 8 | 476** | 3* | 9K | pneumonia |
| 9 | 476** | 3* | 10K | stomach cancer |

Martin Weise, Technical University of Vienna

- Sensitive Data (=due privacy issues, commercial interests, …)

- Provide access for analysis, but ensure data is not leaked or misused

- Standard approach
  - Pseudonymization
  - Anonymization
  - $k$-anonymity
  - $l$-diversity
  - $t$-closeness

**sensitive attributes**

| | ZIP Code | Age | Salary | Disease |
|---|---|---|---|---|
| 1 | 47677 | 29 | 3K | gastric ulcer |
| 2 | 47602 | 22 | 4K | gastritis |
| 3 | 47678 | 27 | 5K | stomach cancer |
| 4 | 47905 | 43 | 6K | gastritis |
| 5 | 47909 | 52 | 11K | flu |
| 6 | 47906 | 47 | 8K | bronchitis |
| 7 | 47605 | 30 | 7K | bronchitis |
| 8 | 47673 | 36 | 9K | pneumonia |
| 9 | 47607 | 32 | 10K | stomach cancer |

0.167-closeness (Salary)
0.278-closeness (Disease)

**Alice** knows **Bob** is 27yo, lives in 47678 and is in the first three entries

⊆ **Bob** has a ???

| | ZIP Code | Age | Salary | Disease |
|---|---|---|---|---|
| 1 | 4767* | ≤ 40 | 3K | gastric ulcer |
| 3 | 4767* | ≤ 40 | 5K | stomach cancer |
| 8 | 4767* | ≤ 40 | 9K | pneumonia |
| 4 | 4790* | ≥ 40 | 6K | gastritis |
| 5 | 4790* | ≥ 40 | 11K | flu |
| 6 | 4790* | ≥ 40 | 8K | bronchitis |
| 2 | 4760* | ≤ 40 | 4K | gastritis |
| 7 | 4760* | ≤ 40 | 7K | bronchitis |
| 9 | 4760* | ≤ 40 | 10K | stomach cancer |

Martin Weise, Technical University of Vienna

# 4.3. OSSDIP

- **Data Owner** maintains full control over data and use:
  - Who to **allow access**,
  - Over which **period of time**,
  - For which **subset of data**,
  - To answer which research question / analysis goals,
  - While monitoring what they are doing
- Based on experience of operating DEXHELPP for nearly 10 years

Martin Weise, Technical University of Vienna

# 4.3. OSSDIP

## Technical Architecture:

Martin Weise, Technical University of Vienna

(highlight)

1. Researcher sends **request** to Data Owner (*Person, question, required data*)

2. Granted: **subset of data**, at specific **aggregation level**, potentially with **fingerprint** is extracted onto a VM for a dedicated **researcher** for a dedicated **time period** to address the **question** posed

3. Expose metadata of data subsets (**FAIRness**)

4. ...

5. Provisioning of VNC and Compute VMs with dedicated software and data
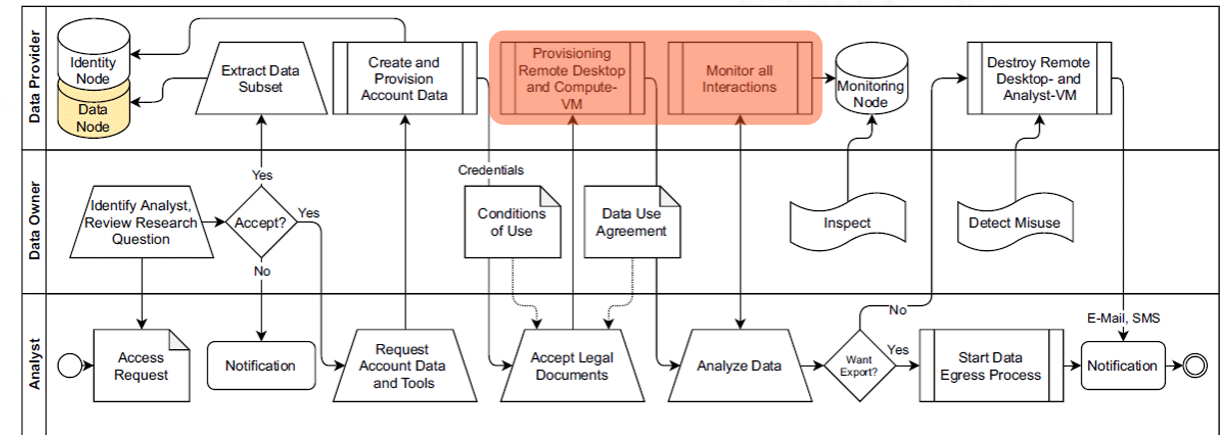
Martin Weise, Technical University of Vienna

(highlight)

1. Researcher sends **request** to Data Owner (*Person, question, required data*)

2. Granted: **subset of data**, at specific **aggregation level**, potentially with **fingerprint** is extracted onto a VM for a dedicated **researcher** for a dedicated **time period** to address the **question** posed

3. Expose metadata of data subsets (FAIRness)

4. ...

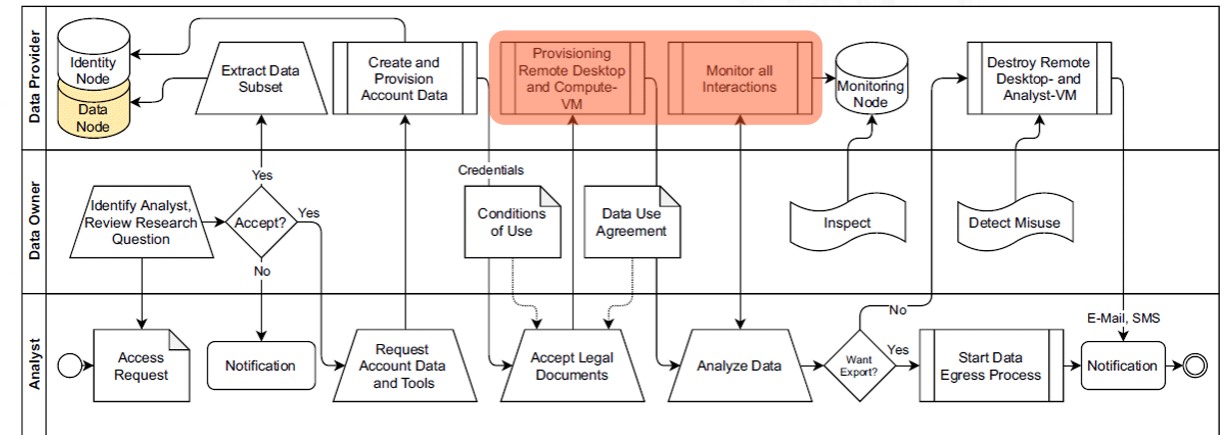5. Provisioning of VNC and Compute VMs with dedicated software and data

Martin Weise, Technical University of Vienna

# 4.3. OSSDIP Further Reading

Material

- https://doi.org/10.5334/dsj-2022-004 (journal paper)

- https://doi.org/10.34726/hss.2022.84700 (master's thesis)

- https://doi.org/10.17605/OSF.IO/VKN4R (iPRES'21 paper)


Resources

- https://ossdip.at/ (documentation, getting started guide)

- https://gitlab.tuwien.ac.at/martin.weise/ossdip (source code)

Martin Weise, Technical University of Vienna

# 5. Future Work

DBRepo

– Prepare for **test** phase Q1 2023

– Prepare for **rollout** phase Q1 2024

– Document all endpoints, methods, files, readme, changelog, etc.

– Implement OAI-PMH interface for metadata harvesting

OSSDIP

– Find (friendly) test-users that want to deploy it within their premises using **synthetic** data

– Implement DBRepo features to make sensitive data **findable** and **reusable** (<u>FAIR</u>)

– **Transform** into Virtual Research Environment (e.g. integrating *Jupyterhub, Collaborative chat*)

– HPC

Martin Weise, Technical University of Vienna

# 6. Conclusion

Different Repositories

- Must support a **wide range** of research objects
- Provide **good visibility** to them
- Will be a place where research happens

Trusted Research Environments

- Traditional privacy methods **not sufficient** for exploratory research
- **Virtual** meeting points to work with sensitive data and known tools
- Organizational / Technical / Legal

**Composition** of systems needed

- Gradual development with increasing complexity
- Integration and automation are key to facilitate adoption

Repositories and TREs are as **trustworthy** as their institutions

Martin Weise, Technical University of Vienna

# Contact

**Martin Weise**

Technical University of Vienna

martin.weise@tuwien.ac.at

0000-0003-4216-302X