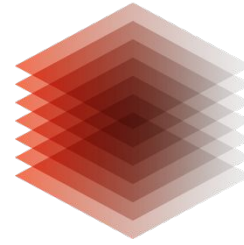LEIBNIZ INFORMATION CENTRE
FOR SCIENCE AND TECHNOLOGY
UNIVERSITY LIBRARY

# Findable

Angelina Kraft, Katrin Leinweber
TIB, 9. July 2018                    Recording: doi.org/10.5446/37823
FAIR Data & Software (Carpentries-based workshop) **#TIBFDS**

Leibniz
Association

# to be **F**indable 🔍

F1. (meta)data are assigned a globally unique and eternally persistent identifier

F2. data are described with rich metadata

F3. (meta)data are registered or indexed in a searchable resource

F4. metadata specify the data identifier

# Your institution's / repository's role

- assign a globally unique PID upon publication (or draft upload)

- provide metadata schema in human- & machine-readable format

  - PID, author names, subject areas, etc.

  - support structured input of metadata (submission forms or XML schema)

  - index (meta)data to enable effective searching

  - allow metadata upload & assign corresponding PID

# Your role as a scientist

$F$indable $Q$

- check datasets that you use for a PID & cite it

- ensure that your datasets get published with a PID

  - choose repositories that automate this

  - report this requirement to repos that don't

- add rich metadata (describe dataset's context, quality, condition & characteristics)

  - should be understandable by researchers from different discipline (ask a friend to proofread)

# Example of paper - data citation using PIDs

## Paper:

Koen Kole, Rik G.H. Lindeboom, Marijke P.A. Baltissen, Pascal W.T.C. Jansen, Michiel Vermeulen, Paul Tiesinga, Tansu Celikel (2017):
**Proteomic landscape of the primary somatosensory cortex upon sensory deprivation**, *GigaScience*, Volume 6, Issue 10, 1 October 2017, Pages 1–10. DOI
https://doi.org/10.1093/gigascience/gix082

## Note in the paper:

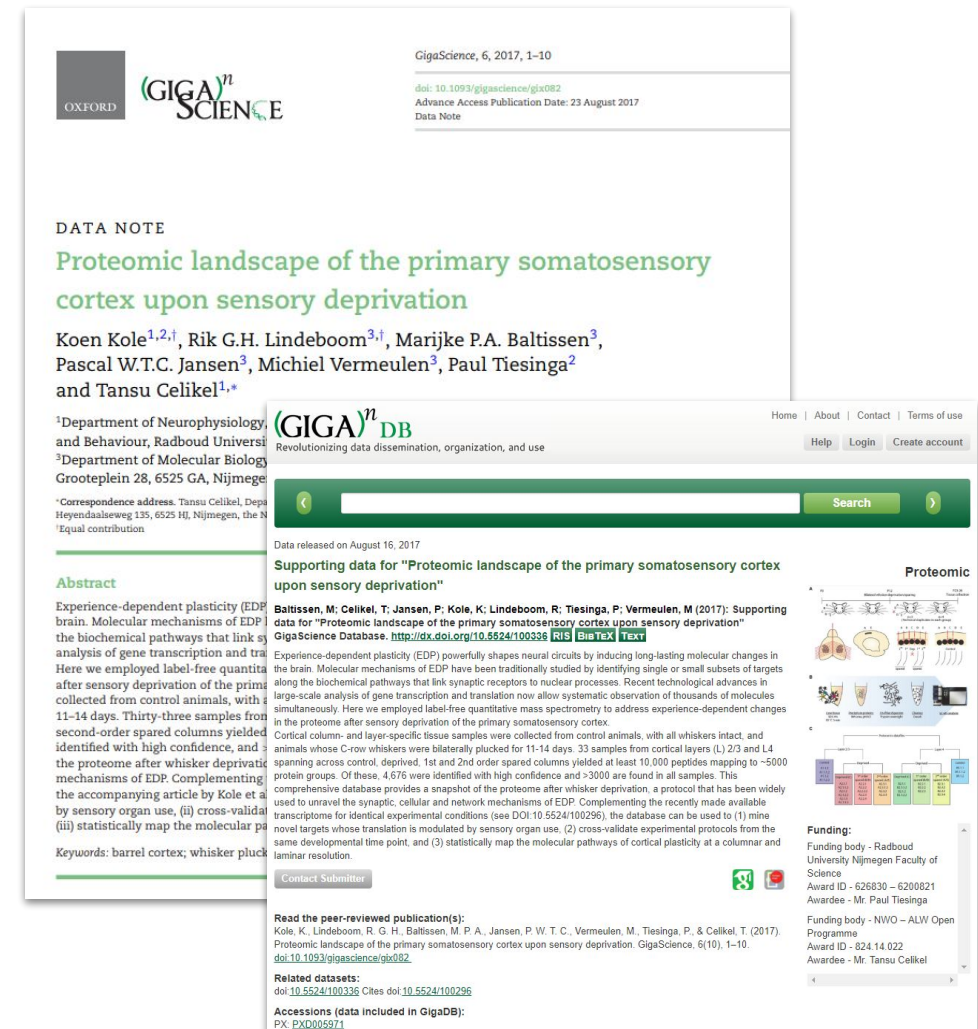**"Availability of the supporting data**
Data supporting this work are available in the *GigaScience* repository, *Giga*DB [14].
The raw mass spectrometry proteomics data have been deposited in the ProteomeXchange Consortium via the PRIDE partner repository [15] with the dataset identifier PXD005971"

## Reference:

[14] Kole K, Baltissen M, Lindeboom R et al.   Supporting data for "Proteomic landscape of the primary somatosensory cortex upon sensory deprivation." GigaScience Database  2017. http://doi.org/10.5524/100336

# Findability Agenda

1. **Persistent Identifiers (PIDs)**
   - **Which ones are there? How should they be used?**
   - **DOIs for research data → minimum criteria for a good data repo**

2. **Choosing FAIR repositories**

3. **Setup help**

4. **Welcome reception**

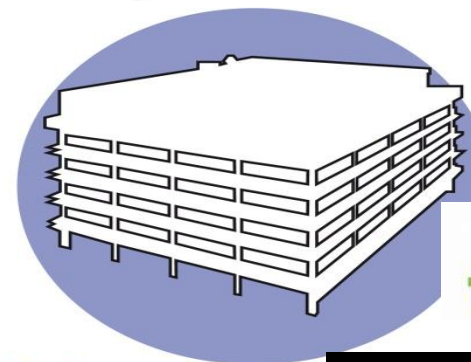# PIDs are everywhere:

## Researcher IDs

ORCID  Scopus®

isni

RESEARCHERID

THOMSON REUTERS

## Organisation IDs, Funder IDs

fund ref

Ringgold
Identify

GRID

## Resource IDs (articles, data, software, …)

doi®  ePIC
Persistent Identifiers for eResearch

Handle.Net®

ARK (Archival Resource Key)

URN-SERVICE

PICHE – Persistent Identifiers for Cultural Heritage Entities

# And even more new PIDs…

- Projects IDs
- Instrument IDs
- Ship cruises IDs
- Physical sample IDs,
- DMP IDs…

**Answer: Do researchers need to care about PIDs? → YES!**

**But what do they need to KNOW about PIDs?**

→ Remember: For a scientist, it is about the project, equipment, DMP, researcher, funder, resource …

→ It is not about the PID. PIDs are infrastructure.
→ In order to use PIDs, scientists do not need to know all about their whereabouts. A basic knowledge should be enough.

# TIB Survey 2017

**1400 scientists in the natural sciences & engineering** (across Germany)

→ 70% of the researchers are using DOIs for journal publications

 → less than 10% use DOIs for research data
(have a look: data available at doi.org/10.22000/54)

**Why?**
- 56% answered that they don't know about the option to use DOIs for other publications (datasets, conference papers etc.)
- 57% stated no need for DOI counselling services
- 40% of the questioned researchers need more information
- 30% cannot see a benefit from a DOI

With the new digital age: Possibilities & struggles!

Have a look: https://www.re3data.org/search
Out of more than 2115 repository systems listed in re3data.org in July 2018, only 809 (less than 39 %!) state to provide a PID service, with 524 of them using the DOI system

# A PID is

- Provenance
- Metadata
- Policies & Guarantees
- Machine readability
- Metrics



**Researchers should know that…**

Provenance means validation & credibility – a researcher should comply to good scientific practices and be sure about what should get a PID (and what not).

Metadata is central to visibility and citability – metadata behind a PID should be provided with consideration.

Policies behind a PID system ensure persistence in the WWW - point. At least metadata will be available for a long time.

Machine readability will be an essential part of future discoverability – resources should be checked and formats should be adjusted (as far possible).

Metrics (e.g. altmetrics) are supported by PID systems.

# PIDs provide interoperable Metadata

- Example:
→ Automatic ORCID profile update when DOI is minted

DataCite – CrossRef – ORCID
  collaboration

  → PID of choice for RDM:
  Here: The Digital Object Identifier (DOI)



If you authorize Crossref and DataCite to update your ORCID record

and you add your ORCID to your paper or dataset submission

when your publication gets a DOI, your ORCID record will get updated

AUTOMATICALLY!

# Digital Object Identifier (DOI)

- Persistent and unique identifier for objects in the digital environment

- DOIs refer to the objects not the location → remain valid

- DOIs are minted for research data, software and code, physical objects, grey literature

- DOI-System is an internationally recognised and supported standard

https://doi.org/10.15468/dl.n1glrt

Proxy    Prefix    Suffix

Proxy          Prefix          Suffix

https://doi.org/10.15468/dl.n1glrt

# The DOI Service

- National mandate to support German academic institutions and the research sector publishing and citing research output, including:

  - research data
  - software
  - videos
  - images
  - 3D models
  - grey literature

- With focus on Science and Technology
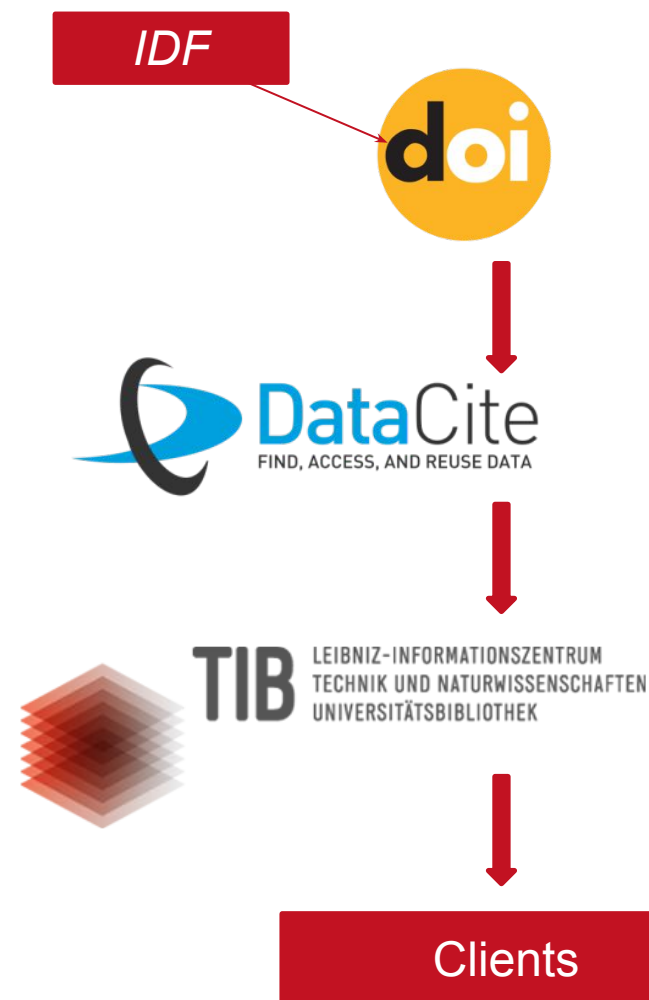- TIB is founding member of DataCite, a global DOI agency for research data

# TIB Services

- Managing prefix assignment

- First-Level-Support

- Training, counseling and hand-outs

- Providing access to the registration plattform (Metadata Store and DOI Fabrica)

- Support in all aspects of good research data management practice

# Service Structure

- International DOI Foundation (IDF) manages DOI-System

- DataCite e.V. maintains and operates the infrastructure for DOI registration

- TIB as DOI Provider grants access to the infrastructure and provides services

- Client registers DOIs via TIB

# DOI at TIB: Facts

**Registered DOI**

- **Total 1,293,389 DOIs** (June 2018)

  - 65 % Research data
  - 25 % Grey literature
  - 10 % Images
  - 0,4 % AV media

**External clients of the DOI-Service**

- **Total 165 data centers**

  - Major research centers i.e. Pangaea, WDCC and ESO
  - 68 universities/university libraries
  - 12 Leibniz Institute /11 Helmholtz Zentren

**TIB user of the DOI-Service**

  - catalogue: Team DTF, retro-digitisation
  - TIB-Portal: AWI reports, AV-Media

# PID 101 for Researchers (*or: Resolving some PID myths*)

1. A PID is a „long lasting reference to a digital resource"

2. There are different sorts of PIDs & different uses,
   (e.g. for articles, data, persons, organizations, …)

3. PIDs are offered by organizations - Ask your institute/library

4. You do NOT have to pay for PIDs (by yourself)!

5. PIDs are mostly used for (persistent) citation – All published resources should have one

6. A correct citation always includes a PID → look in your citation manager

7. Metadata behind a PID are most important – please take care when providing them

8. PIDs are not perfect (they are issued by organizations, aka humans!)

9. PIDs are really useful & fun – they make yourself & your work more visible!

## And what researchers do NOT need to know…
### (*although some may want to know*)

- Total number of PIDs registered
- Names of the agencies
- Names & schemes of identifiers
- How persistence works
- How PID providers struggle
- How (and why) PID providers fight each other
- How perfect a PID (system) is (it is certainly not)

→ Researchers care about their research (= their passion)

→ As long as their not in information science themselves, PID providers should focus on communicating the practical points

→ Citeability & visibility; the benefit for the researcher should be crystal clear

# Summary: PIDs (DOIs) = the glue!

**Digital CV / CRIS**

**Data & AV Repository**

**Portal**

DOI

DOI

DOI
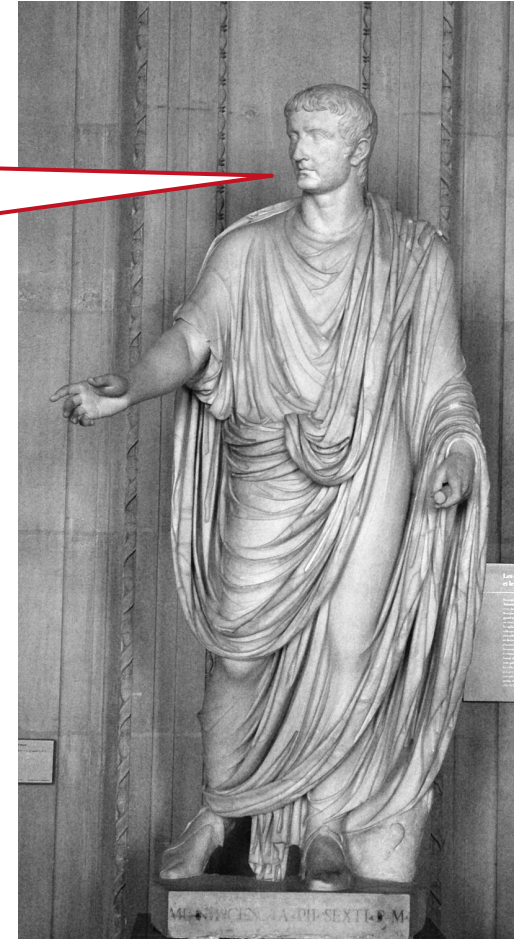
DOI

# Data Journals

Articles in Data Journals
- describing research data
- Containing information about data creation and collection, data qualities, functionalities
- Showing possibilities of re-using data
- Recommendation or provision of data deposit
- Discipline specific or generic data journals
- Examples
  - Nature Scientific data
  - Biomedical Data Journal (BMDJ)

# GitHub + Zenodo = DOI (guides.GitHub.com/activities/citable-code)

- official integration thanks to science.Mozilla.org/projects/codemeta

- intrinsic IDs (Git's SHA1 hashes) vs. "minted" PIDs

  - technical vs. procedural persistence

- Zenodo: file backup & persistent landing page for each release version

  - SoftwareHeritage.org ingests automatically, but not in real-time (source)

- Read more Software.ac.uk/tags/software-preservation

- demo: GitHub.com/TIBHannover/BacDiveR/issues/14

When in Rome, dress like the Romans.

Public Domain by Marie-Lan Nguyen, via
commons.Wikimedia.org/w/index.php?curid=549920

# Findability Agenda

abstraction layers that direct you to the actual object locations

1. **Persistent Identifiers (PIDs)**
   - **Which ones are there? How should they be used?**
   - **DOIs for research data → minimum criteria for a good data repo**

2. **Choosing FAIR repositories**

3. **Setup help**

4. **Welcome reception**

# The right repository?

*There is no right or wrong*

Decision for or against a specific repository depends on various criteria, e.g.
- Data quality
- Discipline
- Institutional requirements
- Reputation (researcher and/or repository)
- Visibility of research
- Legal terms and conditions
- Data value (FAIR Principles)
- Exit strategy (tested?)
- Certificate (based only on documents?)

*Decision has to be taken by the researcher him-/herself*

Help: <u>Checklist – where to keep research data</u> (DCC)

# Benefits of data sharing/publication in (good) data repositories

- data are kept safe in a secure environment
- data are regularly backed up and preserved (long-term) for future use
- data can be easily discovered by search engines and included in online catalogues
- intellectual property rights and licencing of data are managed
- access to data can be administered and usage monitored
- the visibility of data can be enhanced
- enables more use and citation
- citation of data increases researchers scientific reputation
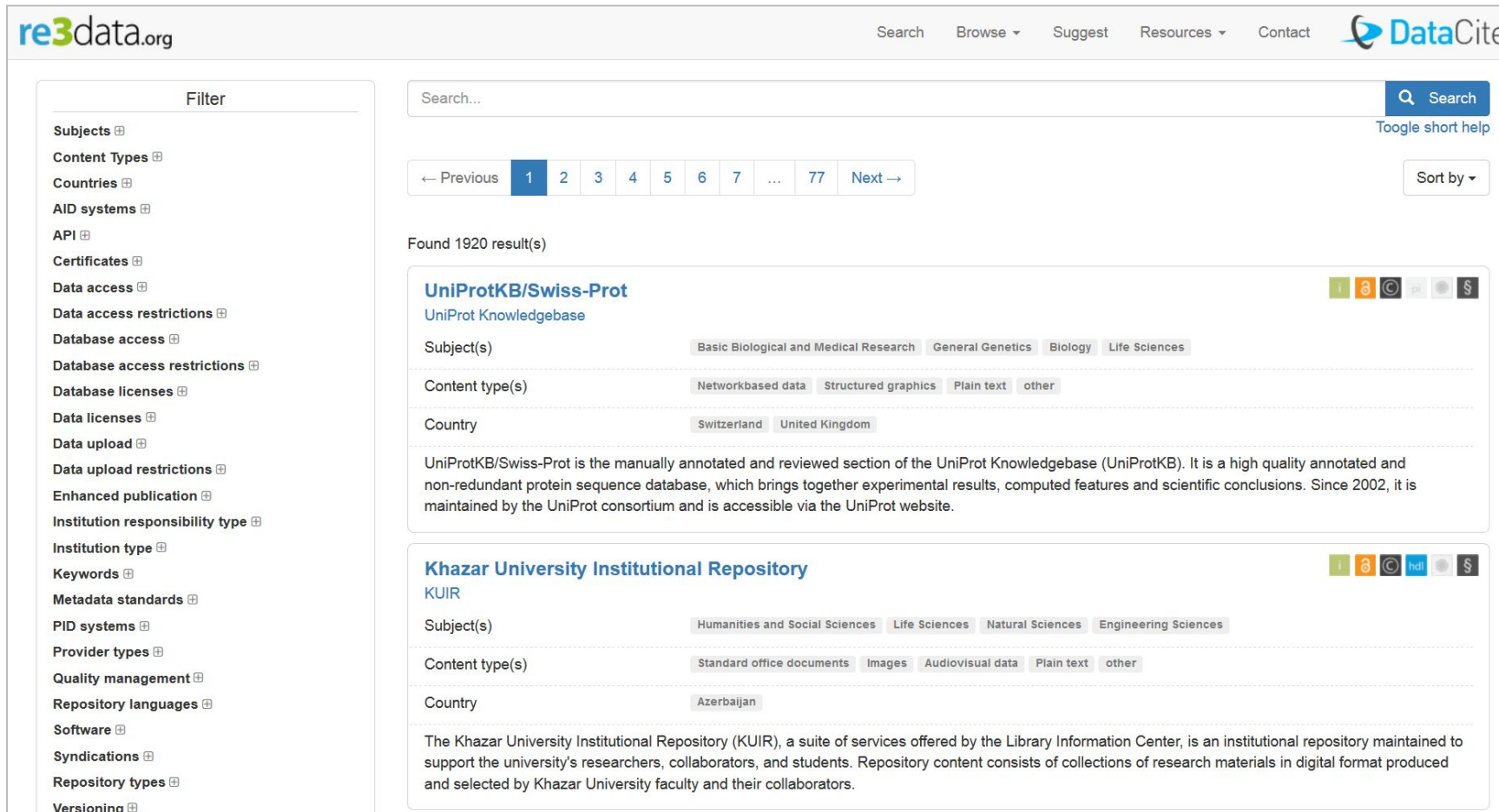
Some recommendations:

→ look for the usage of PIDs

→ look for the usage of standards (DataCite, Dublin Core, discipline-specific metadata

→ look for licences offered

→ look for certifications (DSA / Core Trust Seal, DINI/nestor, WDS, …)

# FAIRsharing.org vs. re3data.org (June 2018)

| | FAIRsharing.org<br>standards, databases, policies | re3data.org<br>REGISTRY OF RESEARCH DATA REPOSITORIES |
|---|---|---|
| Mode of use, operator | Online platform<br>University of Oxford | Online platform<br>DataCite (KIT) |
| Type of data | Metadata records of "*Standards, databases, policies, and Collections /Recommendations*" | Metadata records of<br>*data repositories* |
| No. of records | Standards: 1169<br>Databases (=Repositories): 1064<br>Polices: 112<br>Collections/Recommendations: 44<br>Identifier schema: 7 | Metadata standards: 76<br>Repositories: 2096<br><br>Identifier schema (PID Systems): 5 +others |
| Findability | URL; recommendations (EU H2020, journal guidelines) | URL; recommendations (EU H2020, journal guidelines) |
| Accessibility (timeline) | Since 2016 (started with biosharing.org) | Since 2012 (merged with Databib 2015) |
| Interoperability | Mark up with schema.org & BioSchema project planned, API | XML templates, API |
| Re-useability | API (read-only, pre-contact required)<br>Content licensed via CC-BY-SA 4.0 | API<br>Content licensed via CC-0 1.0 |
| Data policy (who can upload information) | Registered users who are maintainers, ORCiD preferred, "sanity checks" by curators, approval/refusal | Repository operators, application form, analyzes of the repository website by curators, metadata checked, approval/refusal |

# How to find a repository

- Ask your colleagues & collaborators
- Look for institutional repository at your own institution
- Search at re3data.org – Registry of Research Data Repositories

# Findability Agenda

*abstraction layers that direct you to the actual object locations*

1. **Persistent Identifiers (PIDs)**
   - Which ones are there? How should they be used?
   - DOIs for research data → minimum criteria for a good data repo

2. **Choosing FAIR repositories**

3. **Setup help needed? Red stickies, please!**

4. **Welcome reception at Waterloostraße 1; directions:
   TIBHannover.GitHub.io/2018-07-09-FAIR-Data-and-Software/#schedule**

LEIBNIZ INFORMATION CENTRE
FOR SCIENCE AND TECHNOLOGY
UNIVERSITY LIBRARY

**TIB**

# Which questions do you have for us?

**Contact information:**

Katrin.Leinweber@TIB.eu & Angelina.Kraft@TIB.eu

T +49 511 762-14693 & -14238

Leibniz
Association