

How to **UNHIDE** and improve the metadata landscape of research software in Helmholtz.

Jens Bröder¹, F. D'Mello¹, G. Preuß², S. Fathalla¹, P. Buttigieg³,
O. Mannix², V. Hofmann¹, S. Sandfeld¹

¹Institute for Advanced Simulation, Materials Data Science and Informatics (IAS-9), Forschungszentrum Jülich GmbH

²Helmholtz-Zentrum Berlin für Materialien und Energie GmbH (HZB)

³Alfred Wegner Institute for Polar and Marine Research
deRSE 2024 | 07.03.2024

www.helmholtz-metadaten.de

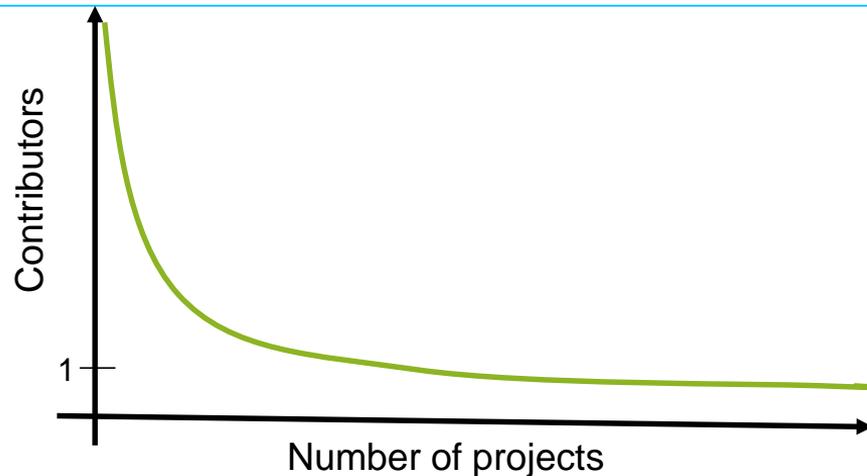


Better metadata for research software now!

Problems with RSE:

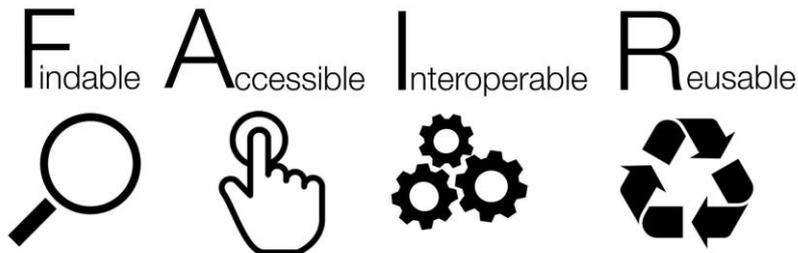
- Funding: Software is usually a by product
- Experience/Quality: Developer on the side
- Persistence: usually single developer
- Wrong incentive: I do my own software
- Hidden: uploaded somewhere

➔ Most projects die over a few years.



Better metadata helps:

- Find existing software for your problem.
- Understanding software.
- Foster reuse of software and collaboration.



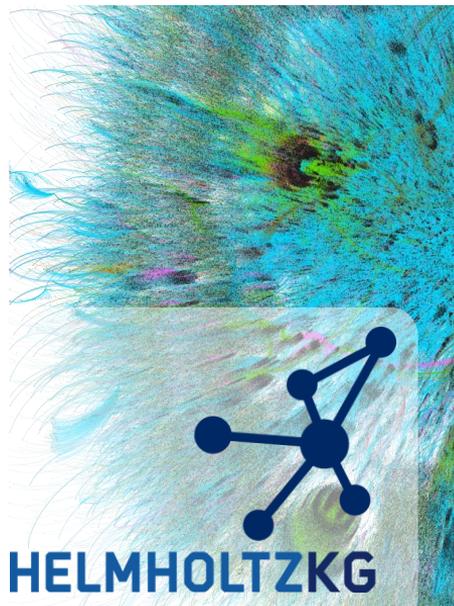
CC-BY-SA-4.0,
wikimedia commons

[1] M. D. Wilkinson. Scientific Data 3 (2016), pp. 1–9., also see www.go-fair.org

Outline

1. UnHIDE initiative and the Helmholtz knowledge graph:

- Why, what, where?

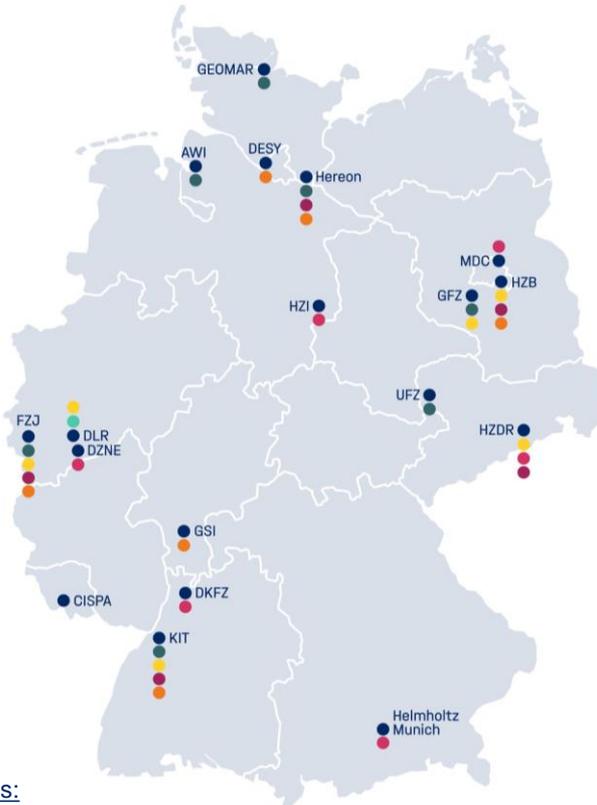
The logo for UnHIDE, featuring the word "UNHIDE" in blue capital letters. Above the letters "H" and "I" are two blue arrows pointing towards each other, with a small blue circle positioned between them.

2. Software metadata in detail:

- General problems & status
- Metadata in the graph
- How to improve?



What is the (meta)data reality for the Helmholtz association?



Research Fields:

- Energy, Earth and Environment
- Health, Information
- Aeronautics, Space and Transport, Matter

18 independent research centres in 6 different research fields all of which host digital infrastructure



libraries

metadata on published research & data sets



data repositories

research data (cold, medium-hot or hot)



code repositories, also external

code & software

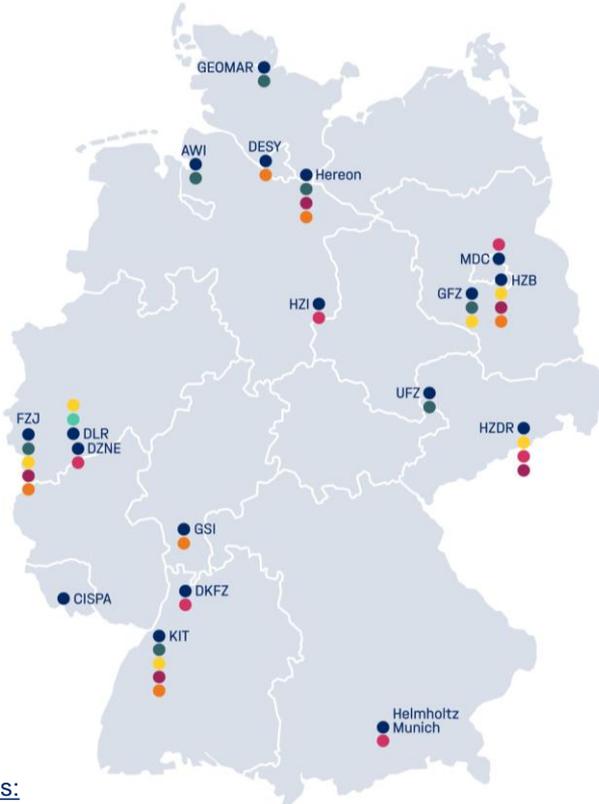


research infrastructures

heterogenous methods and approaches



- unified Helmholtz information and data exchange



Research Fields:

- Energy, ● Earth and Environment
- Health, ● Information
- Aeronautics, Space and Transport, ● Matter

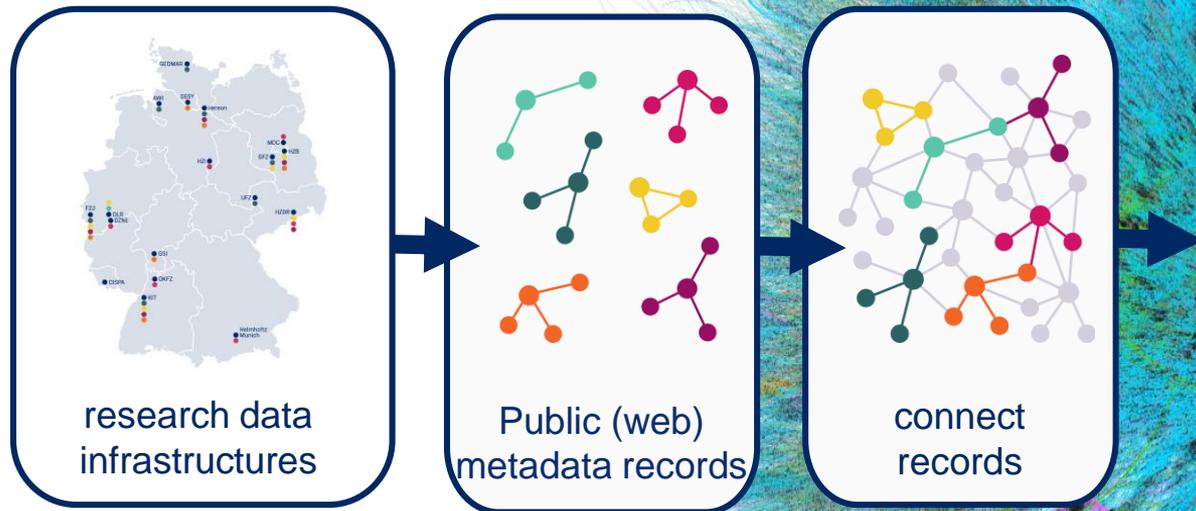
UNHIDE Scope

1. Create a lightweight interoperability layer by consolidating (meta)data in an association wide knowledge graph
2. Increase visibility of Helmholtz digital infrastructures
3. Improve quality and interoperability of Helmholtz metadata
4. Make digital assets easily findable
5. Assess the status quo of Helmholtz Metadata



HELMHOLTZKG

Helmholtz Knowledge Graph



conceptual & code co-development
with <https://oceaninfohub.org/>

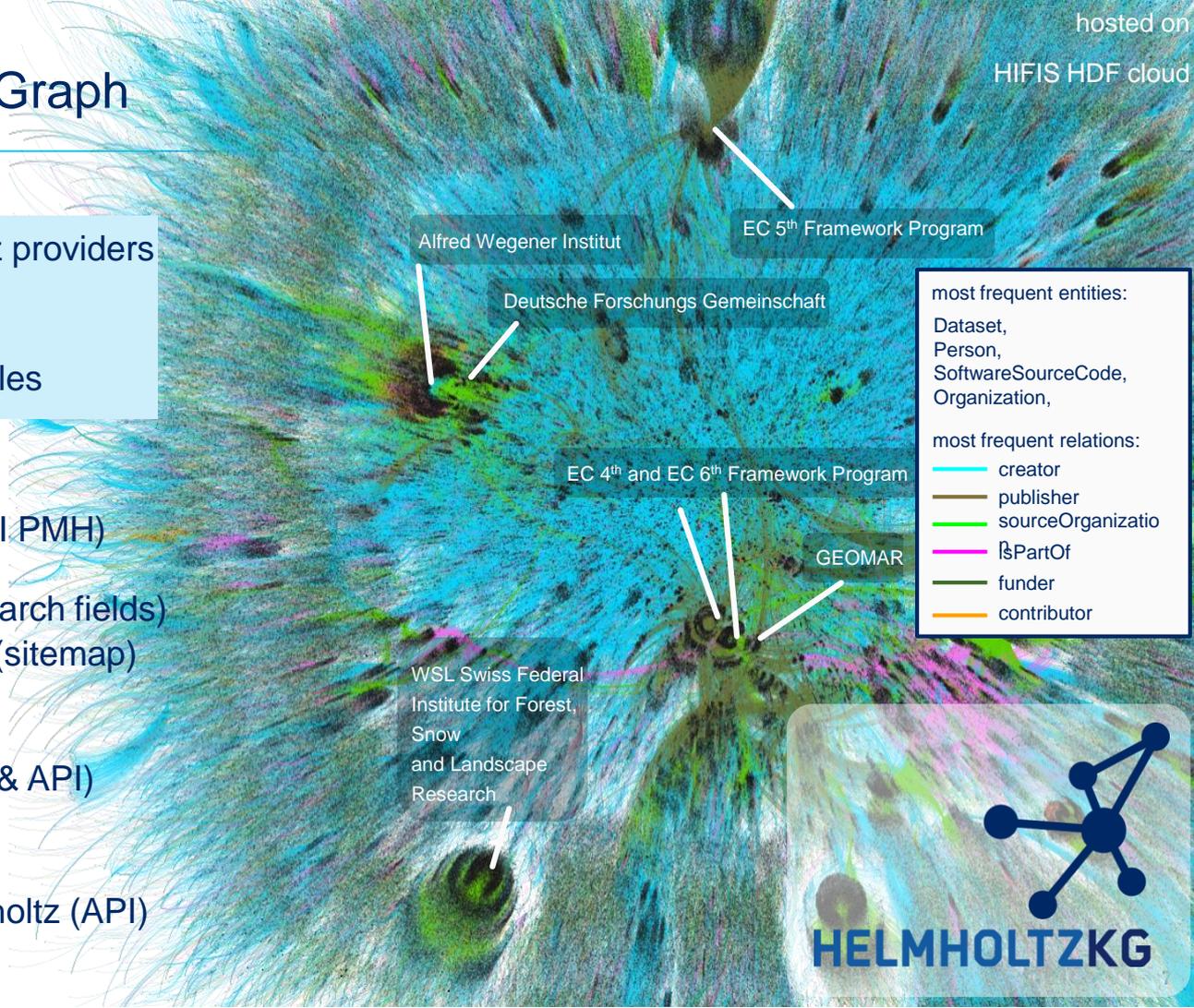


Helmholtz Knowledge Graph

data from 32 different Helmholtz providers
in one place!

> 2.15 mio records | 7.1 mio triples

- **Libraries**
16/18 Helmholtz libraries (OAI PMH)
- **data repositories** (from 3 research fields)
Rodare, Pangea, Jülich data (sitemap)
- **code repositories**
12 (all) GitLab Instances (Git & API)
- **global resources**
DataCite sub-graph for Helmholtz (API)



search.unhide.helmholtz-metadaten.de

UNHIDE

Helmholtz Association (HGF) Update

Try: GEOMAR -- Jülich Data -- DESY -- Rodare --

"Meaningfully combining data from heterogeneous sources is a knowledge graph's main value proposition."
Andrew Senior, The Knowledge Graph Cookbook

The Unified Helmholtz Information and Data Exchange (UNHIDE) Project aims to build a sustainable, interoperable, and inclusive digital ecosystem for all stakeholders. Existing and emerging data systems are linked via the Helmholtz Knowledge Graph (see right), with the ultimate goal of coordinating action and capacity to improve access to scientific publications, software, data and knowledge. The Project is funded by the centers of the Helmholtz Association and implemented by the Helmholtz Metadata Collaboration (HMC).

[To learn more about this project - click here](#)

search.unhide.helmholtz-metadaten.de

Categories

- Experts - 2794
- Documents - 1744
- Trainings - 5
- Datasets - 4134
- Software - 3.5k
- Projects - 5
- Institutions - 644
- Instruments - 0

HELMHOLTZ
Research for grand challenges.

HELMHOLTZ
Metadata
Collaboration

© 2023 Helmholtz Metadata Collaboration

Conductor Permalink

Extensions: cxml save to dav sponge User: SPARQL

OPENLINK®
VIRTUOSO
UNIVERSAL SERVER

Execute Query Reset

Execution timeout milliseconds

Options

- Strict checking of void variables
- Log debug info at the end of output (has no effect on some queries and output formats)
- Generate SPARQL compilation report (instead of executing the query)

Copyright © 2023 [Openlink Software](#)
 Virtuoso version 07.20.3238 (d89671fa1) on Linux (x86_64-ubuntu_bionic-linux-gnu) Single Server Edition (8 GB total memory, 633 MB memory in use)

sparql.unhide.helmholtz-metadaten.de/sparql

docs.unhide.helmholtz-metadaten.de

2. Software metadata in detail

- General problems & status
- Metadata in the graph
- How to improve?



Metadata in usage for (research) software

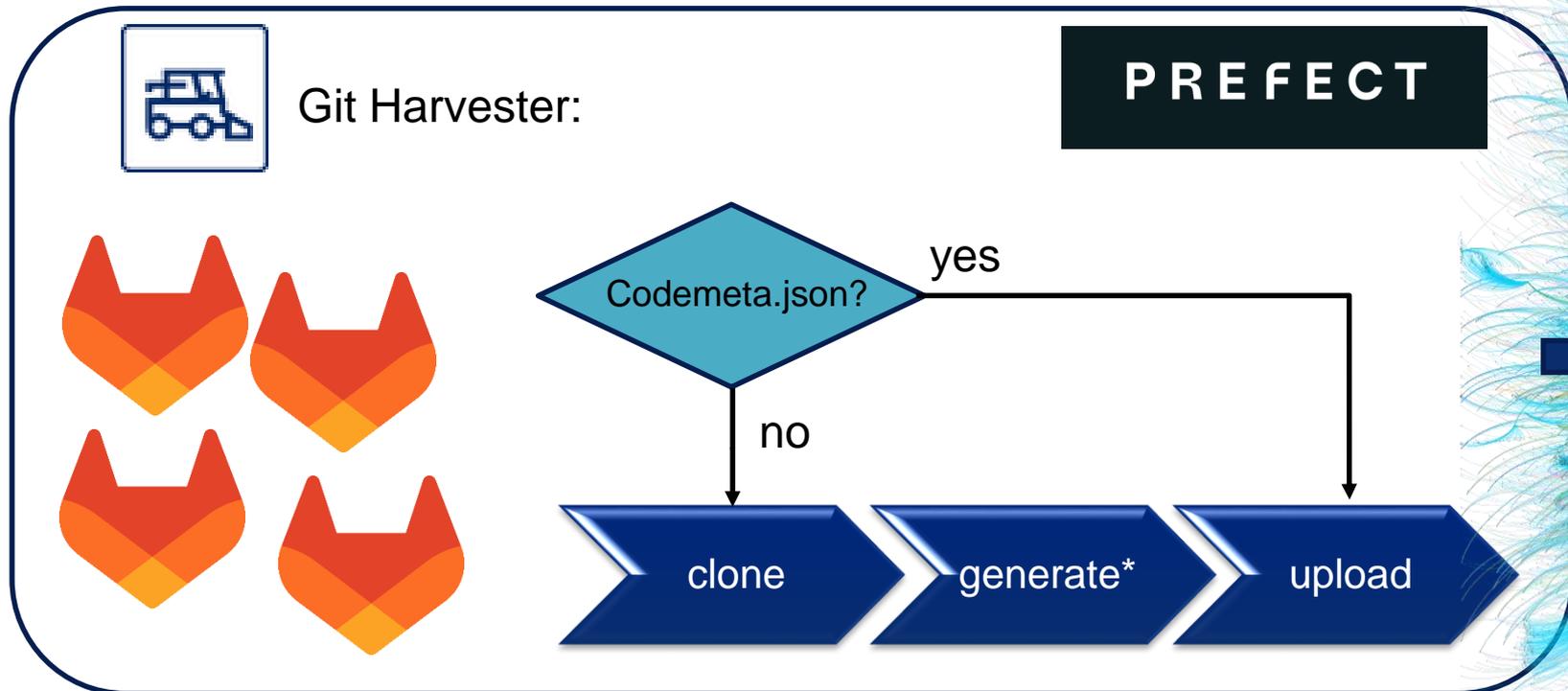
| What? | For Software |
|-----------------------------------|--|
| Standards | Git, language specific, citation.cff, codemeta, ... |
| Size | small |
| Metadata enrichment/ Platforms | <ul style="list-style-type: none">- Manually- Tracks the whole thing- Simple type of change (text)  |
| Re-use case | <ul style="list-style-type: none">- Similar purpose (usually)- As a whole- Install whole, use part (library) |

For datasets:

Different story!

Much, much harder

Data pipeline for software in detail:



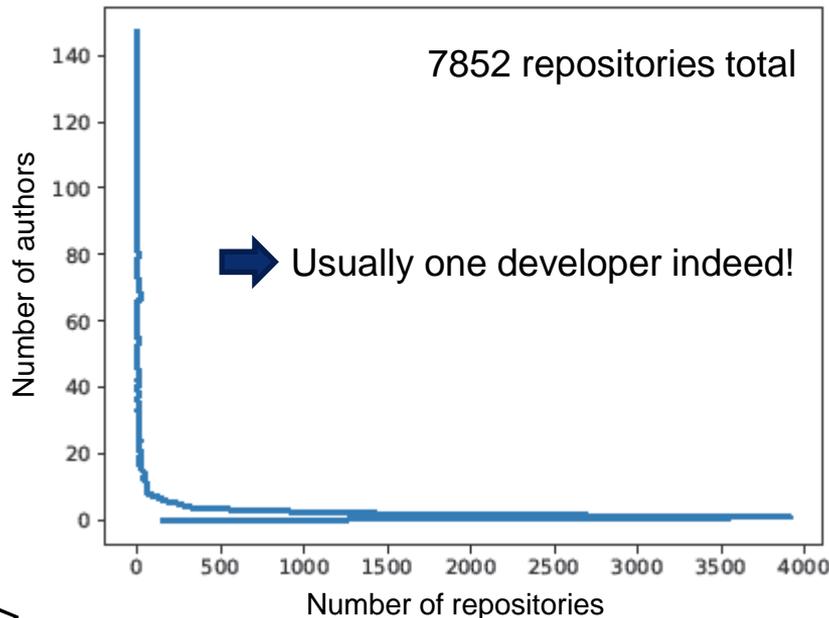
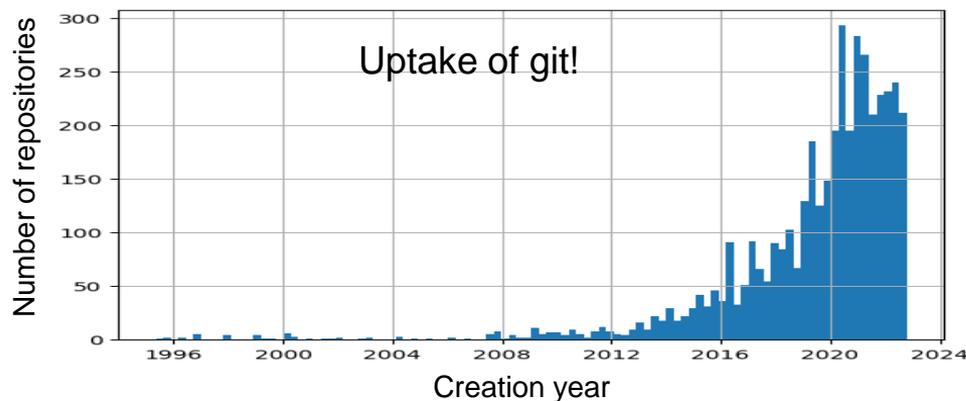
[1] <https://github.com/proycon/codemeta-harvester>

[2] <https://github.com/proycon/codemetapy>

[3] <https://github.com/github/linguist/>

*Run Codemeta-harvester¹ (codemetapy²,
github-linguist³, git commands)

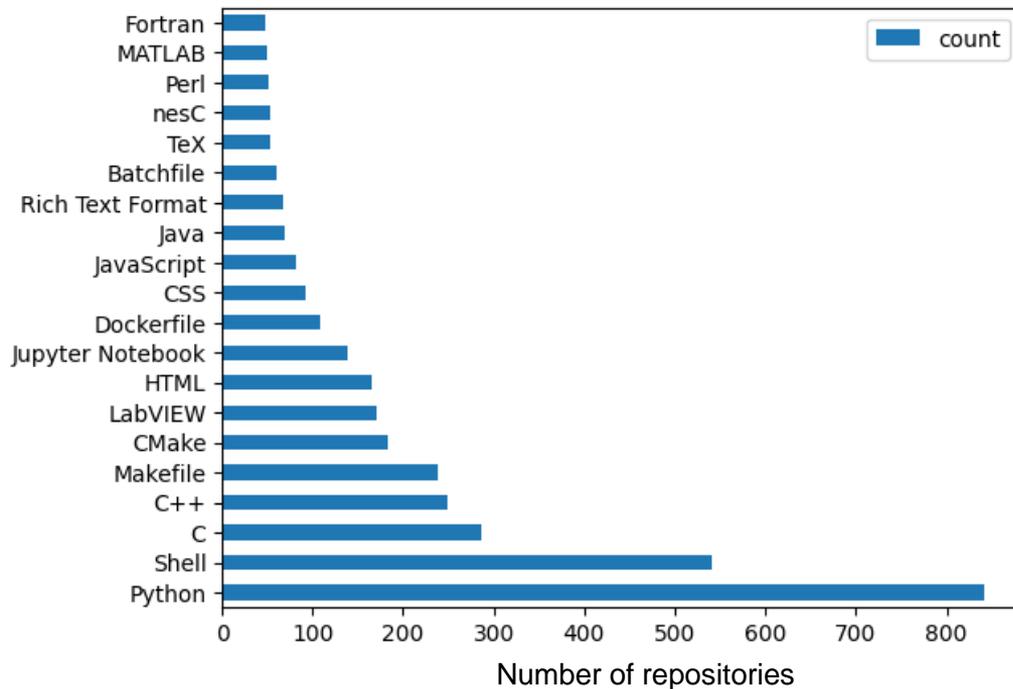
Overview of graph content with respect to RSE software



Statistics:

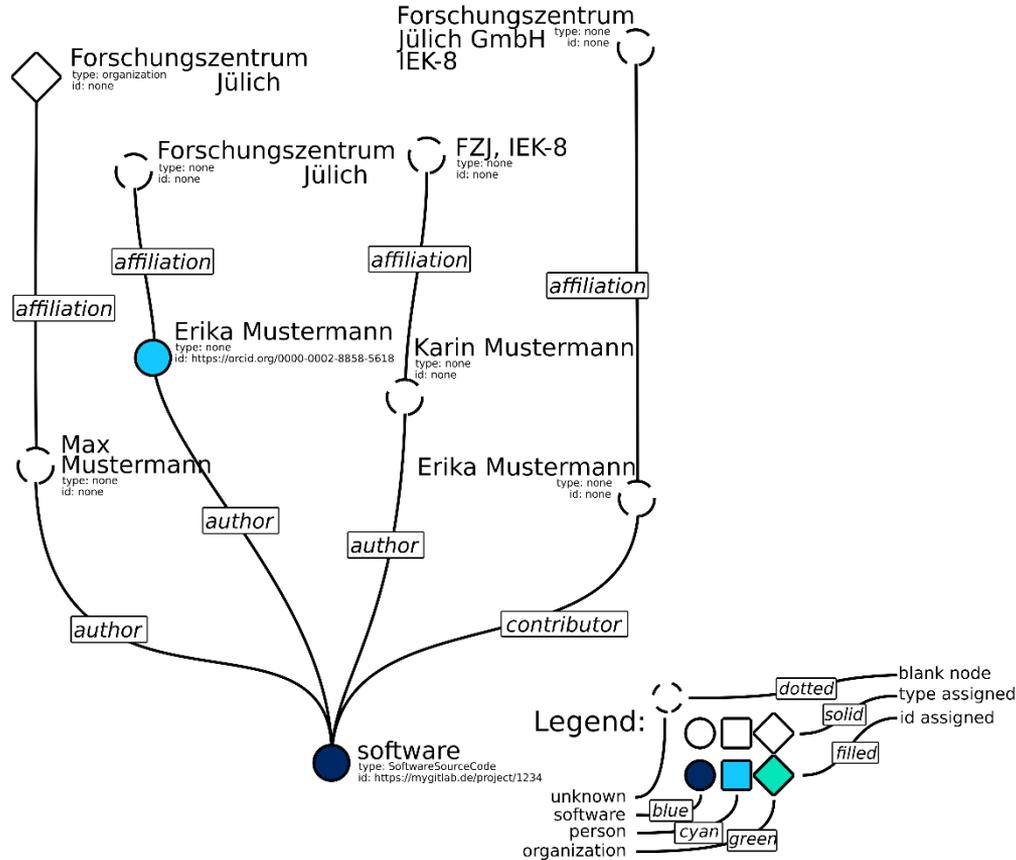
- 600+ active, 800+ WIP (no release)
- 2500+ licenses (MIT > GPL-3 > Apache > BSD > other GPL > others)
- 4500+ tags (possible releases)

Overview of Language use



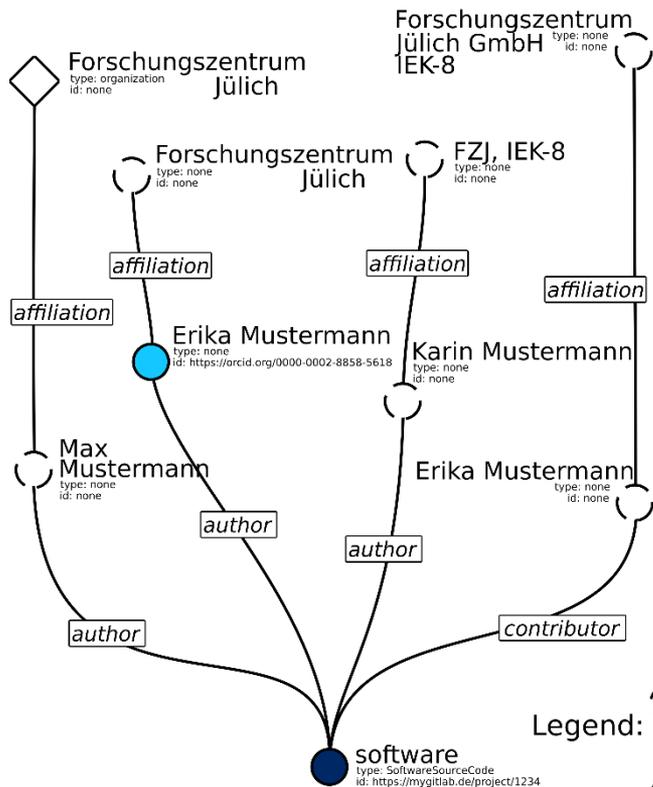
Future: Improving software metadata

Missing PIDS and semantics:

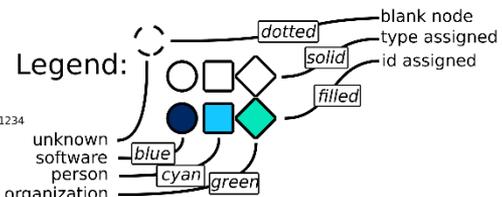
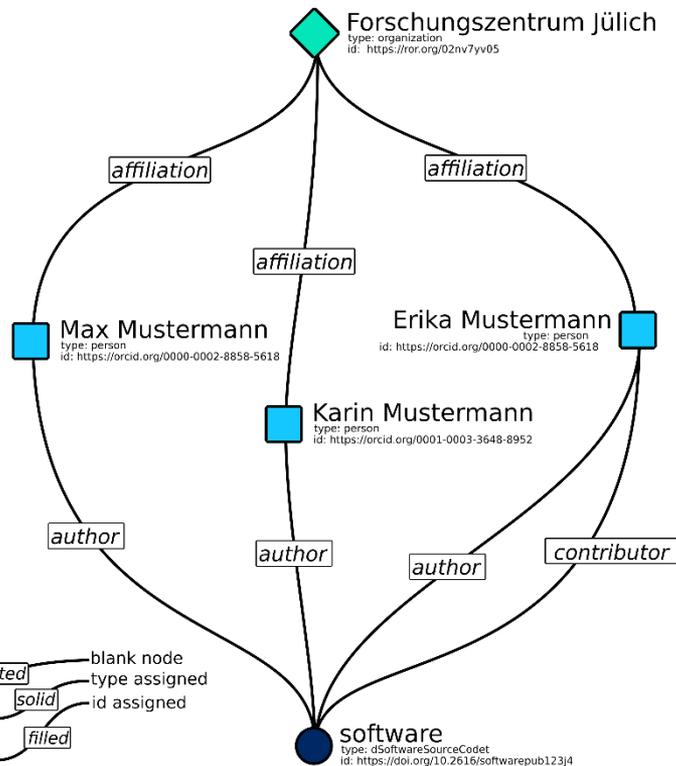


Future: Improving software metadata

Missing PIDS and semantics:



Complete PIDS and semantics:



Strategy



Volker Hofmann
Strategy Development
FZJ / Information



(Pier Luigi Buttigieg)
Former Member
Strategy Development
AWI

Development



Jens Bröder
Tech Lead
FZJ / Information



Gabriel Preuß
Full stack
HZB / Matter



Said Fathalla
Triple Store; SPARQL
FZJ / Information



Fiona D'Mello
Web FrontEnd
FZJ / FDC

Community



Oonagh Mannix
Outreach
HZB / Matter



(Nina Weisweiler)
Former Member
Infrastructures, GFZ

Acknowledgements:

Anton Pirogov: Technical design discussions,
Mustafa Soylyu: Deployment help,
Vivien Serve, Markus Kubin and
Özlem Özkan: UX-workshops,
Pedro V. Barranco: Data/endpoints investigations,
Annika Strupp: early discussion about sense
Stefan Sandfeld (director): Providing resources.
Other HMC staff contributing to visions.

Outlook:

- A metadata bot, with automatic suggestions, improve gitlab/github apis to better expose metadata, software publications
- Harvest also data from github. Problem: finding software by affiliation is hard.



Outlook:

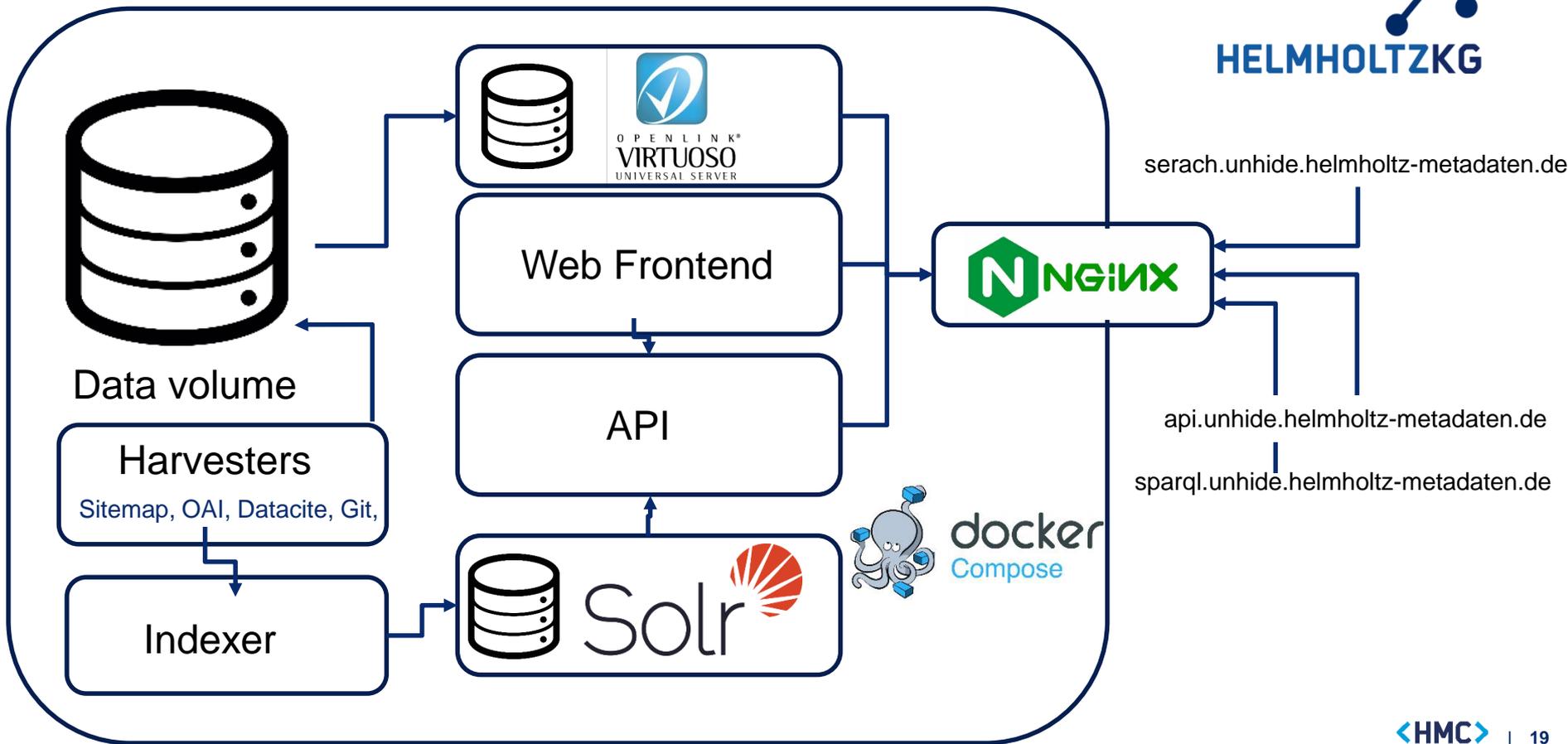
- A metadata bot, with automatic suggestions, improve gitlab/github apis to better expose metadata, software publications
- Harvest also data from github. Problem: finding software by affiliation is hard.



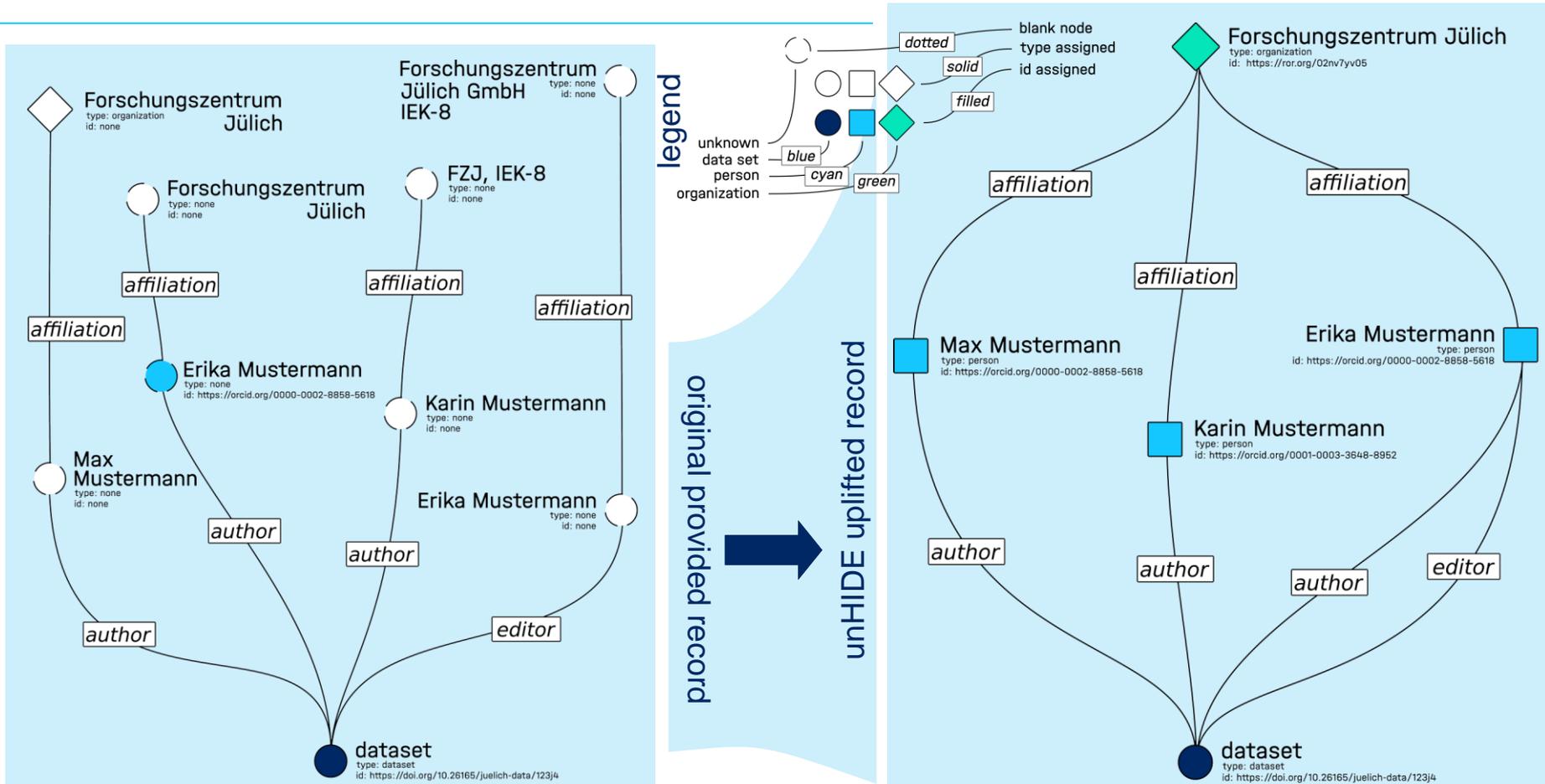
Sum up & take away:

- The goal of the Helmholtz knowledge graph is to make the metadata in Helmholtz visible and improve it ultimately at the source. https://helmholtz-metadaten.de/en/unhide_helmholtz-kg
- We provide federated search for software project through metadata of all Helmholtz gitlab instances

Technical view behind graph:

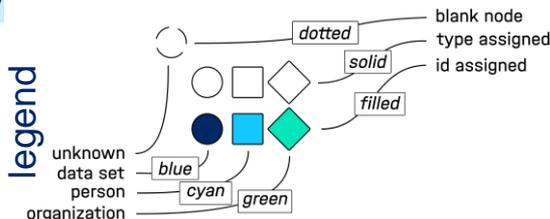


Uplifting increases structuredness of the graph



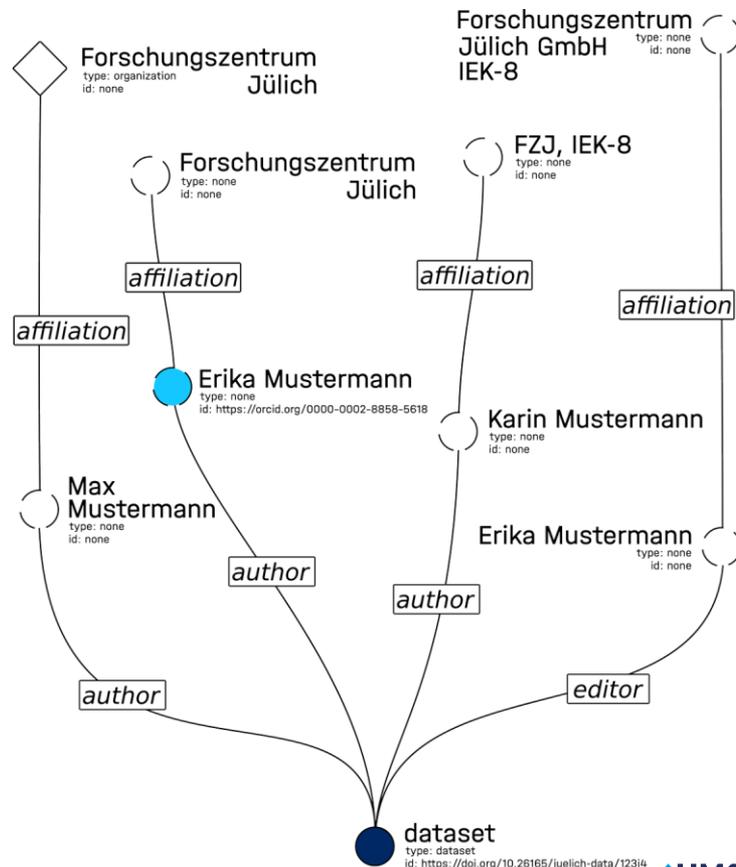
Improving metadata quality at the source

Provided data records are often **incomplete**, **heterogeneous** & **messy**



```
{
  "@context": "http://schema.org",
  "@id": "https://doi.org/10.26165/juelich-data/123j4",
  "@type": "dataset",
  "name": "Name of the dataset.",
  "author": [
    {
      "affiliation": {
        "@type": "Organization",
        "name": "Forschungszentrum Jülich",
        "name": "Max Mustermann"
      },
      "@id": "https://orcid.org/0000-0002-8858-5618",
      "affiliation": {
        "name": "Forschungszentrum Jülich",
        "name": "Erika Mustermann"
      },
      "affiliation": {
        "name": "FZJ, IEK-8",
        "name": "Karin Mustermann"
      }
    },
    {
      "affiliation": {
        "@type": "Organization",
        "name": "Forschungszentrum Jülich GmbH, IEK-8",
        "name": "Erika Mustermann"
      }
    }
  ],
  "editor": {
    "affiliation": {
      "@type": "Organization",
      "name": "Forschungszentrum Jülich GmbH, IEK-8",
      "name": "Erika Mustermann"
    }
  }
}
```

original provided record



Improving metadata quality at the source

The assembled data allows uplifting records through **type inference**, **general harmonization** which allows **resolving entities** and **assigning IDs**.

Uplifting is recorded in reversible patches and will be **fed back to data providers**

```
{ "@context": "http://schema.org",
  "@id": "https://doi.org/10.26165/juelich-data/123j4",
  "@type": "dataset",
  "name": "Name of the dataset.",
  "author": [ { "affiliation": { "@type": "Organization",
    "name": "Forschungszentrum Jülich" },
    "name": "Max Mustermann" },
    { "@id": "https://orcid.org/0000-0002-8858-5618"
    "affiliation": { "name": "Forschungszentrum Jülich" },
    "name": "Erika Mustermann" },
    { "affiliation": { "name": "FZJ, IEK-8" },
    "name": "Karin Mustermann" } ],
  "editor": { "affiliation": { "@type": "Organization",
    "name": "Forschungszentrum Jülich GmbH, IEK-8" },
    "name": "Erika Mustermann" }
```

original provided record



unHIDE uplifted record

```
{ "@context": "http://schema.org",
  "@id": "https://doi.org/10.26165/juelich-data/123j4",
  "@type": "dataset",
  "author": [ { "@id": "https://orcid.org/0000-0003-3648-8952",
    "@type": "Person",
    "affiliation": { "@id": "https://ror.org/02nv7yv05",
      "@type": "Organization",
      "name": "Forschungszentrum Jülich GmbH" },
    "name": "Max Mustermann" },
    { "@id": "https://orcid.org/0000-0002-8858-5618",
    "@type": "Person",
    "affiliation": { "@id": "https://ror.org/02nv7yv05",
      "@type": "Organization",
      "name": "Forschungszentrum Jülich GmbH" },
    "name": "Erika Mustermann" },
    { "@id": "https://orcid.org/0001-0003-3648-8952",
    "@type": "Person",
    "affiliation": { "@id": "https://ror.org/02nv7yv05",
      "@type": "Organization",
      "name": "Forschungszentrum Jülich GmbH" },
    "name": "Karin Mustermann" } ],
  "editor": { "@id": "https://orcid.org/0000-0002-8858-5618",
    "@type": "Person",
    "affiliation": { "@id": "https://ror.org/02nv7yv05",
      "@type": "Organization",
      "name": "Forschungszentrum Jülich GmbH, IEK-8" },
    "name": "Erika Mustermann" },
  "name": "Name of the dataset." }
```

Basics Linked Data in JSON-LD serialization example

```
{ "@context": {  
  "dc11": "http://purl.org/dc/elements/1.1/",  
  "ex": "http://example.org/vocab#",  
  "xsd": "http://www.w3.org/2001/XMLSchema#",  
  "ex:contains": {"@type": "@id"}  
},  
"@graph": [  
  {"@id": "http://example.org/library",  
   "@type": "ex:Library",  
   "ex:contains": "http://example.org/library/the-republic"},  
  {"@id": "http://example.org/library/the-republic",  
   "@type": "ex:Book",  
   "dc11:creator": "Plato",  
   "dc11:title": "The Republic",  
   "ex:contains": "http://example.org/library/the-republic#introduction"},  
  {"@id": "http://example.org/library/the-republic#introduction",  
   "@type": "ex:Chapter",  
   "dc11:description": "An introductory chapter on The Republic.",  
   "dc11:title": "The Introduction"}  
]  
}
```

Library:

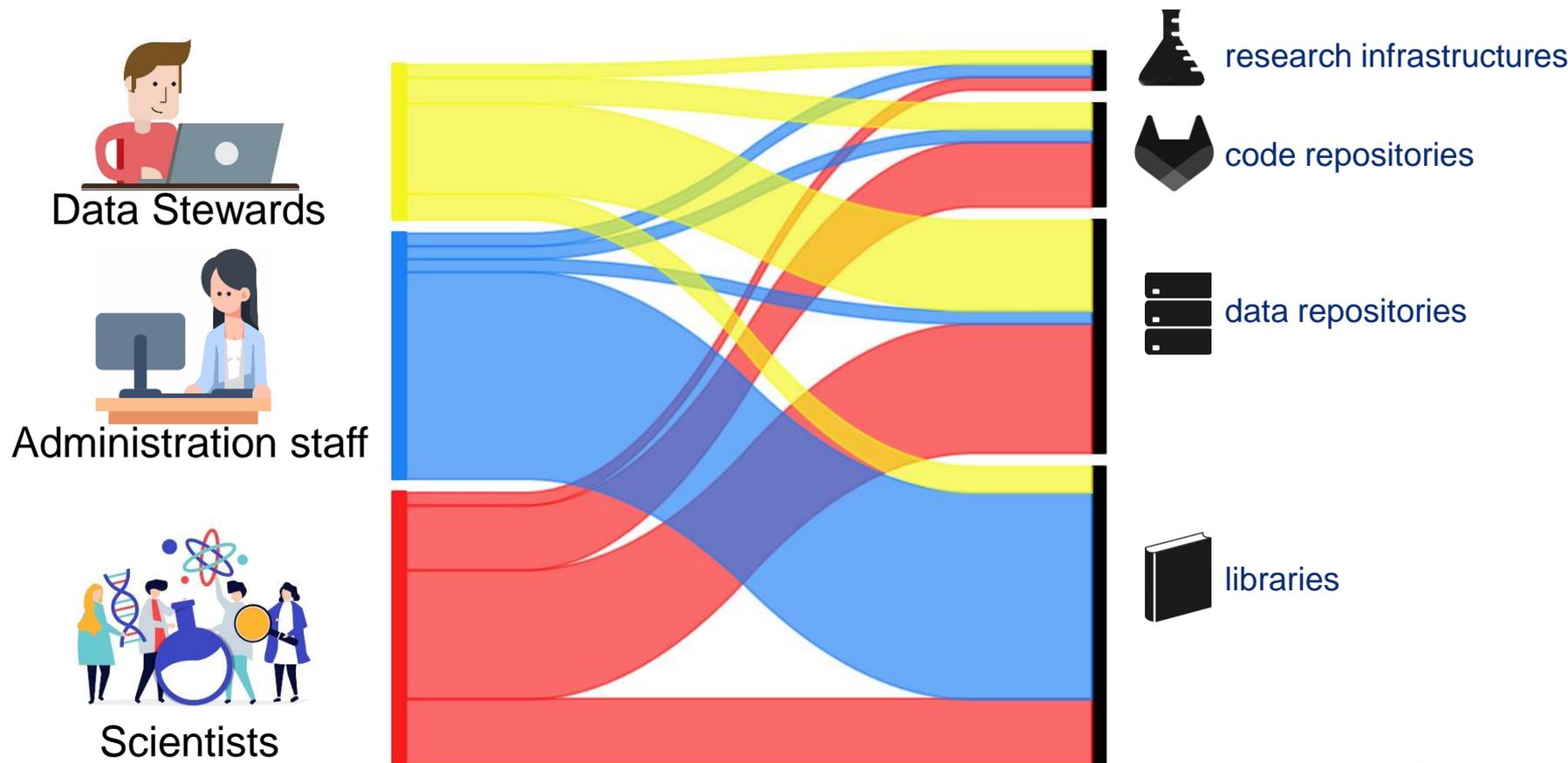
Semantic/RDF triple:



```
{  
  "@context": "http://schema.org/",  
  "@id": "http://orcid.org/0000-0002-1584-4316",  
  "@type": "Person",  
  "name": "Jane Doe",  
  "jobTitle": "Professor",  
  "telephone": "(425) 123-4567",  
  "url": "http://www.janedoe.com"  
}
```

Person:

The Helmholtz (Meta)data ecosystems

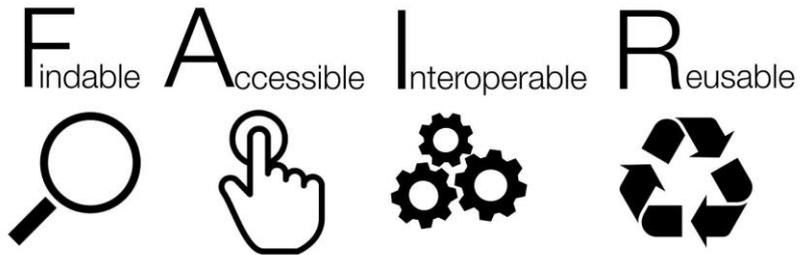
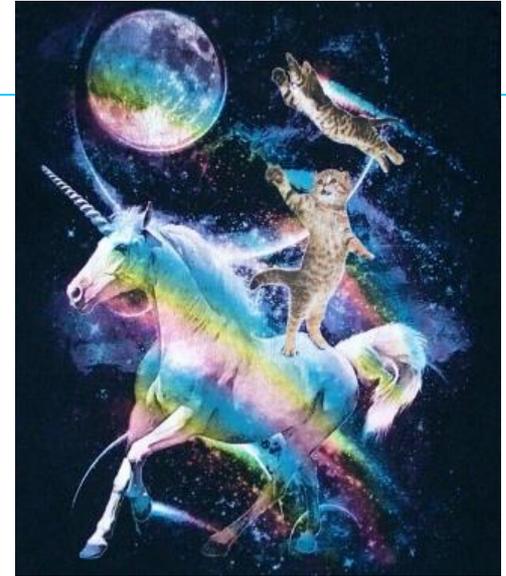


* Only abstract descriptive!

Motivation: sustainable science

In the ideal world:

- Human **knowledge accumulates** as best and sustainable as possible. (Knowledge transfer)
- Research is (easily) **reproducible**. (for simulations that should be doable right?)
- Data, Software and Results are **interoperable** and **re-usable**. (drives Progress)



CC-BY-SA-4.0,
wikimedia commons

How can you bring your research a bit closer to these ideals?

[1] M. D. Wilkinson. Scientific Data 3 (2016), pp. 1–9., also see www.go-fair.org



Introduction

- About
- Implementation overview
- Data Sources

Data in unHIDE

- Overview
- Dataset
- Documents
- Experts
- Institution
- Instruments
- Software
- Training

Interacting with unHIDE data

- Use case examples
- Web search
- SPARQL endpoint
- REST API

Related Knowledge

- Structured data on the web
- Tools around Linked Data
- Other Graphs

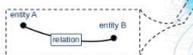
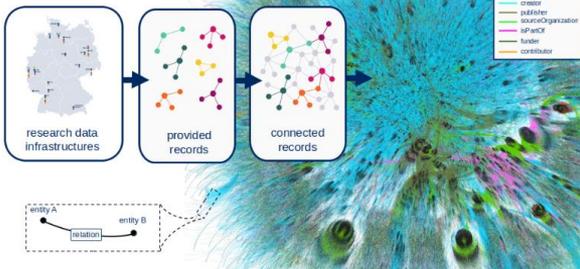
Technical implementation

- Data pipeline
- Architecture



unHIDE - unified Helmholtz Information and data exchange

Helmholtz Knowledge Graph



Introduction & Scope

Research across the Helmholtz Association (HGF) depends and thrives on a complex network of inter- and multidisciplinary collaborations which spans across its 18 Centres and beyond.

However, the (meta)data generated through the HGF's research and operations is typically siloed within institutional infrastructure and often within individual teams. The result is that the wealth of the HGF's (meta)data is stored and maintained in a scattered manner, and cannot be used to its full value to scientists, managers, strategists, and policy makers.

To address this challenge, the Helmholtz Metadata Collaboration (HMC) is launching the **unified Helmholtz Information and Data Exchange (unHIDE)**. This initiative seeks to create a lightweight and sustainable interoperability layer to interlink data infrastructures and provide greater, cross-organisational access to the HGF's (meta)data and information assets. Using proven and globally adopted knowledge graph technology (Box 1), unHIDE will develop a comprehensive association-wide Knowledge Graph (KG) the "Helmholtz-KG": a solution to connect (meta)data, information, and knowledge.

Box 1

What is a Knowledge Graph?

- A "graph", from graph theory, is a structure that models pairwise connections between objects using "nodes" connected by "edges".
- A "knowledge graph" uses such a graph structure to capture knowledge about how a collection of things (ranging from people to data) are connected to one another (for instance). This helps organisations learn to



Contents

- Introduction & Scope
- Table of contents
- Contributors and Partners
- Acknowledgements
- References

general: hmc@fz-juelich.de
 strategy: v.hofmann@fz-juelich.de
 tech: j.broeder@fz-juelich.de