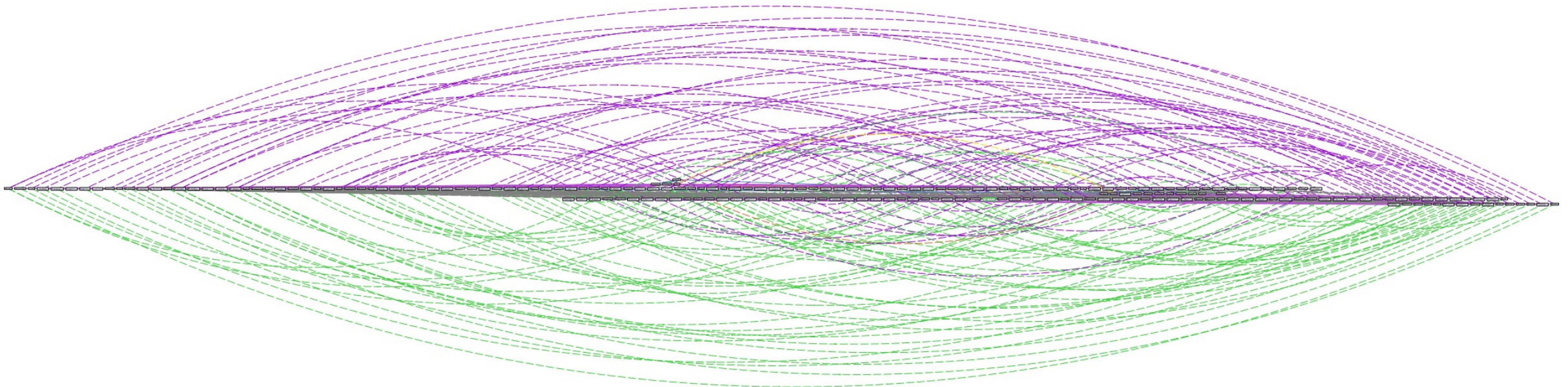# Data Collections Explorer

## An easy-to-use tool for sharing and discovering research data
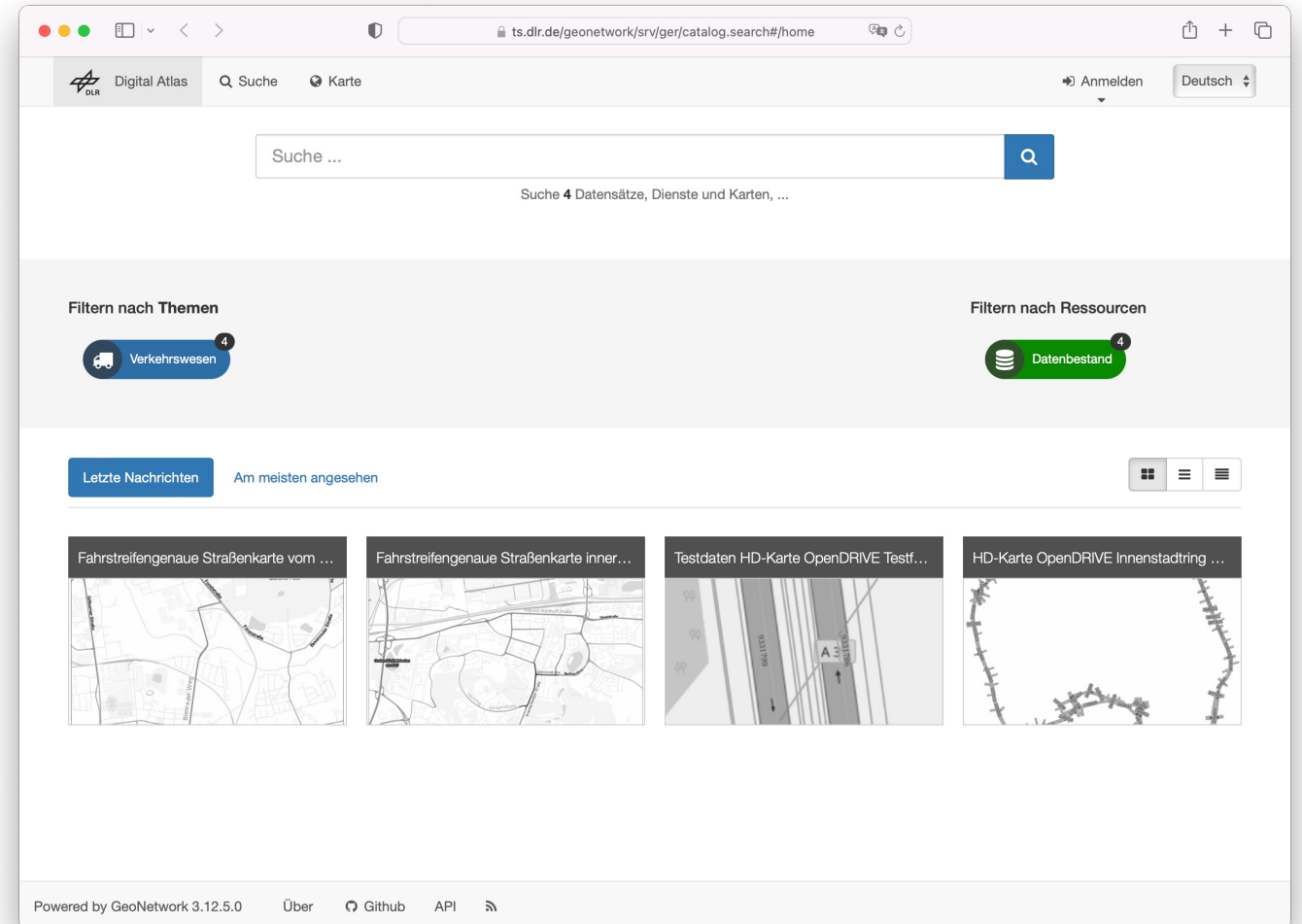
Philipp Ost, Yusra Shakeel, Philipp Tögel – 1st CoRDI, Karlsruhe, 12.–14. September 2023
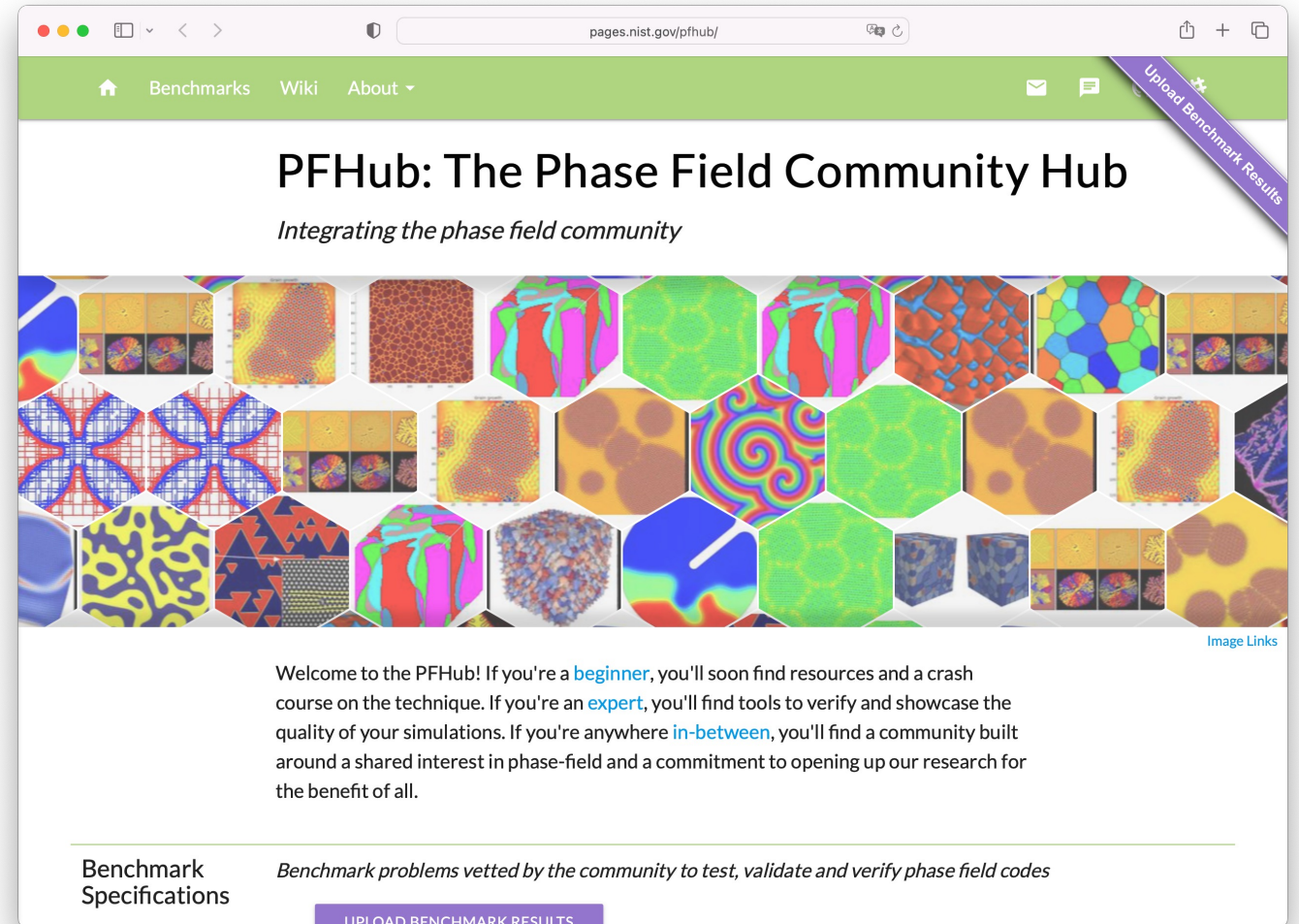
# Motivation – How to share datasets?

- **DLR Digital Atlas**
  - Hosted and curated by DLR
  - Easy sharing of data with collaboration partners

# Motivation – How to serve a community?

- PFHub: The Phase Field Community Hub
  - Hosted by NIST
  - Curated by the phase field community

# Motivation – How to share a dataset?

- KITTI Vision Benchmark Suite
  - Hosted by responsible research group
  - Data access requires an account
  - Usage restrictions apply

# Easy Sharing of Research Data

- What do these examples have in common?
  - Allow relatively easy sharing of research data
  - Allow working with/on the data as part of a community
  - None of them are registered in re3data (as of 30.08.'23)
    - How does one find them?
    - There is a certain amount of "insider knowledge" required to know where to look

- How to share
  - … interim results before publication, possibly with a wider audience?
  - … data and related results with project partners?
  - … data jointly produced/curated/… by a group of researchers?

# Current State of the Art

- re3data – Repository of Research Data Repositories
  - "Gold standard" for research data repositories
  - Admits only quality controlled repositories

- Community specific initiative: Data Repository Finder
  - Geared towards the Life Sciences
  - Overview of repositories
  - "…helps researchers find data repositories where they can share data…"[1]
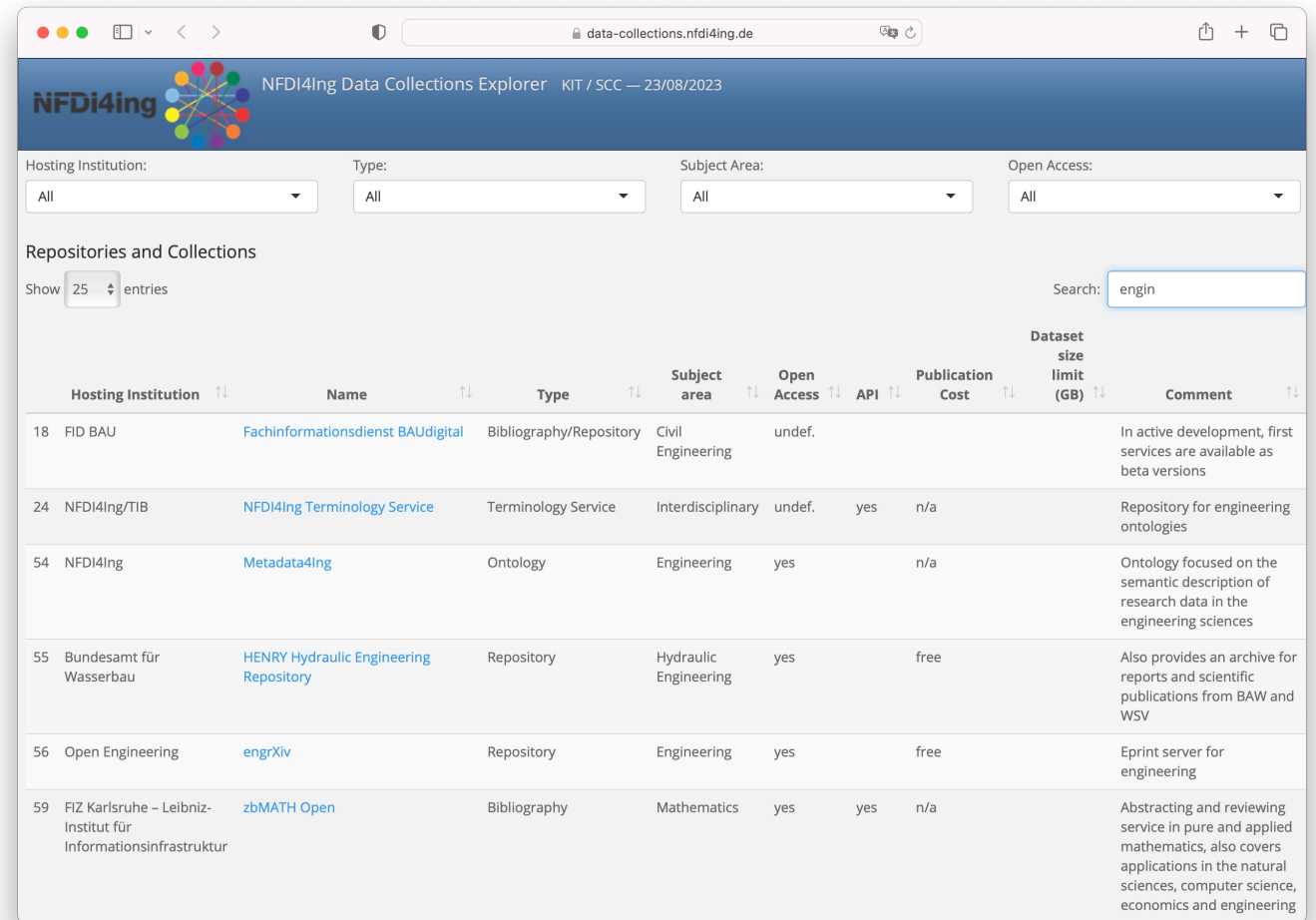
[1] https://data-repository-finder.ll.mit.edu

# Data Collections Explorer – Use Cases

- Typical use cases considered during design and implementation:
  - Scientists looking for repositories to publish their research data
    - Example: PhD students create data sets as part of their research
      - Where to publish them?
      - Are there size limits or costs involved?

  - Scientists searching for data sets
    - Example: An engineer is interested in material properties
      - Are there other materials that might fit the requirements?

September 11, 2023  P. Ost – philipp.ost@kit.edu                                                                      SCC

# Data Collections Explorer – Current Version

- Successful proof-of-concept
  - Started out with 38 entries
  - Continuously updated
    - Currently at 87 entries
- Access it at
  [data-collections.nfdi4ing.de](data-collections.nfdi4ing.de)

# Data Collections Explorer

- Basic implementation completed and available
  - CSV table to store entries
  - Served as HTML table using R Markdown and shiny
- Code is available on GitHub:
  [github.com/kit-data-manager/Data-Collections-Explorer](github.com/kit-data-manager/Data-Collections-Explorer)
  - Adapt it, use it for your own community
- Constant updates
  - Mostly triggered by input from scientists
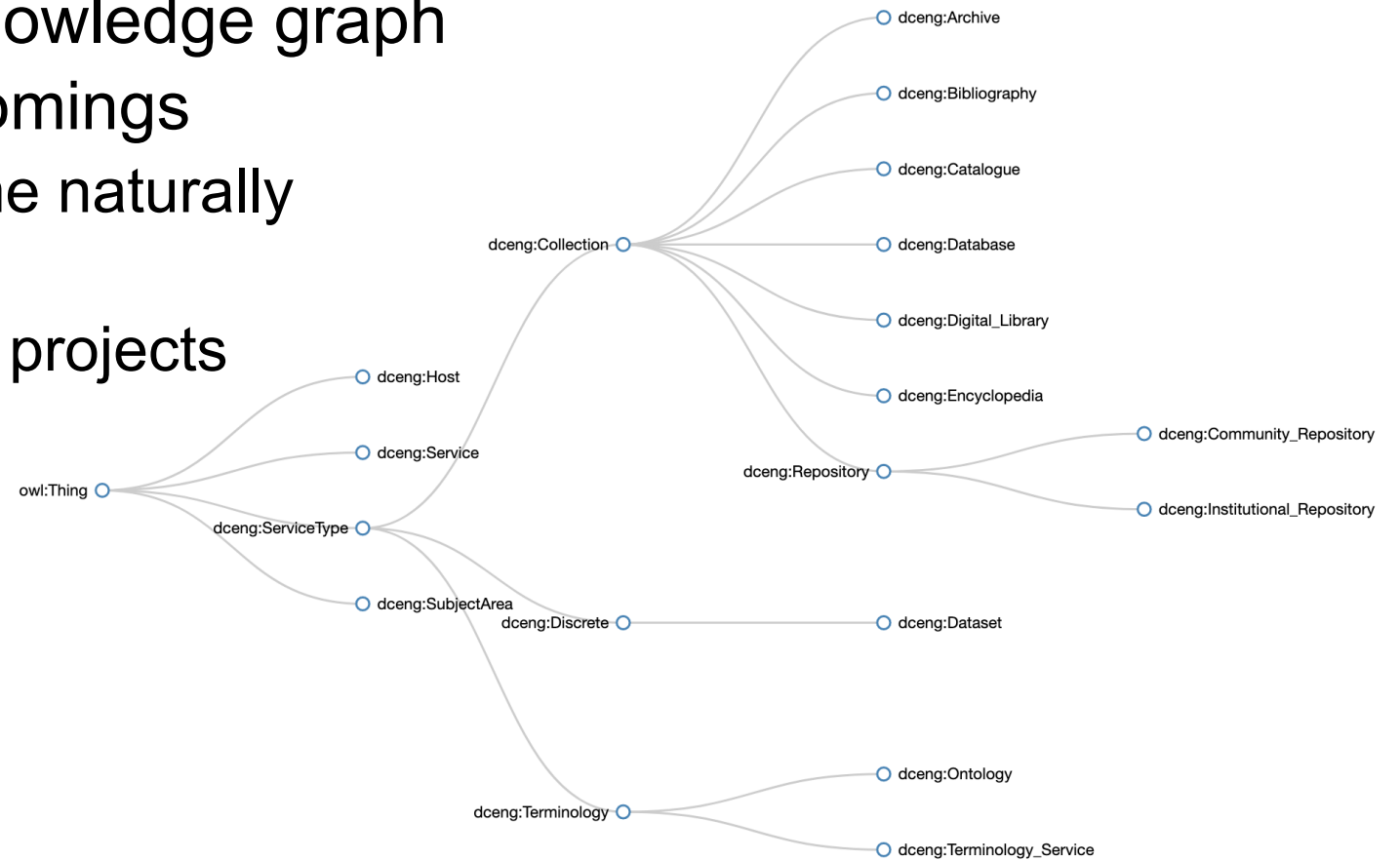- Nonetheless, some caveats apply…

# Collecting Scientists' Feedback

- Feedback collected from all NFDI4Ing Task Areas
  - Questions of interest:
    - What are scientists requirements and wishes of scientists regarding the Data Collections Explorer?
    - How would they like to work with the Data Collections Explorer?
    - Which repositories or data sets are important for their daily work?
  - Aggregated results:
    - Easy-to-use interfaces, including submission of new entries
    - API access
    - Flexible search
    - Use of vocabularies
- Discussions take time, but provide valuable input

September 11, 2023  P. Ost – philipp.ost@kit.edu                                    SCC

# Data Collections Explorer – Knowledge Graph

- Replace CSV table with a knowledge graph
- Addresses almost all shortcomings
  - One-to-many mappings come naturally
  - Access via SPARQL
  - Easier integration with other projects
  - More flexibility
  - Easy machine accessibility

# Conclusions

- The Data Collections Explorer is an easy-to-use flexible tool
  - Share and find data focussed on a dedicated community
  - No big entry hurdles compared to re3data
  - Built for NFDI4Ing, easily adaptable for other communities

- New graph based version is work in progress; what to look forward to:
  - API access
  - New user interface
  - Vocabulary integration
  - And more… ☺