# Towards a Seamlessly Interlinked Research Data and Software Ecosystem at HZDR

2nd Practice Forum Research Data Management, October 20, 2022
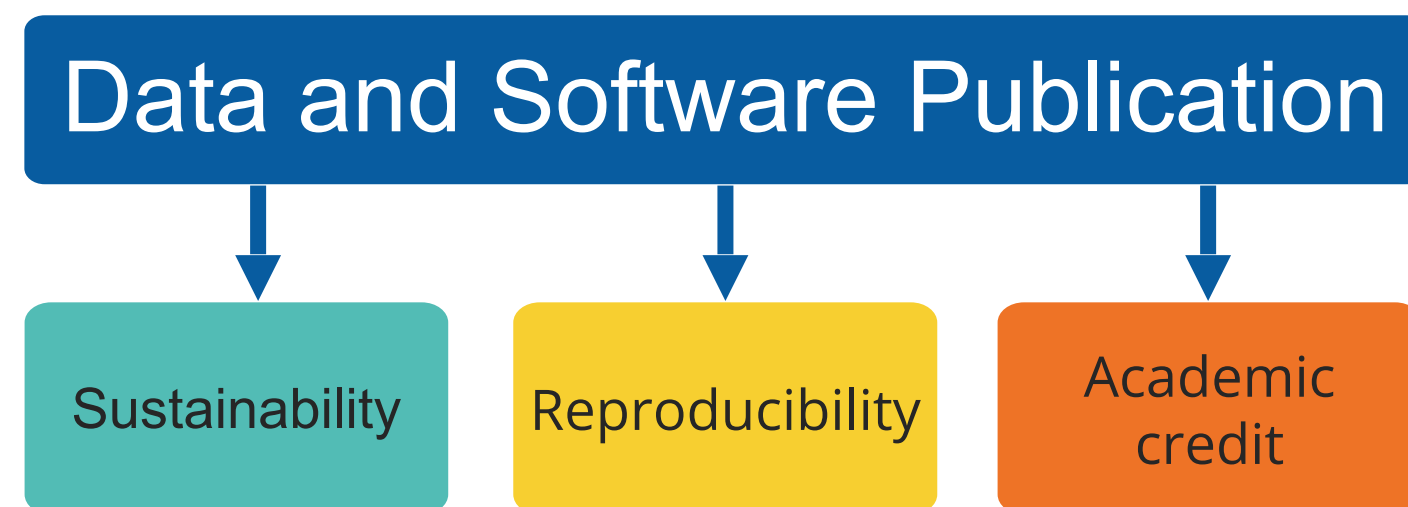
Oliver Knodel // contact: o.knodel@hzdr.de

# Motivation — Software, Data and Everything in Between

— Data and software are an important result of a scientific experiment.

— Scientific publication, research software and data must receive the same academic credit:

**Data and Software Publication**

| Sustainability | Reproducibility | Academic credit |
|---|---|---|

— FAIR principles also exist for research software and should be taken into account [1].

— In addition to the publication itself a seamless interlinking between all available data products is also necessary to improve findability.

**[1]** Barker, M., Chue Hong, N.P., Katz, D.S. *et al.* Introducing the FAIR Principles for research software. *Sci Data* **9**, 622 (2022). doi.org/10.1038/s41597-022-01710-x

...to avoid this:

**F**indable **A**ccessible **I**nteroperable **R**eusable

For Research Software

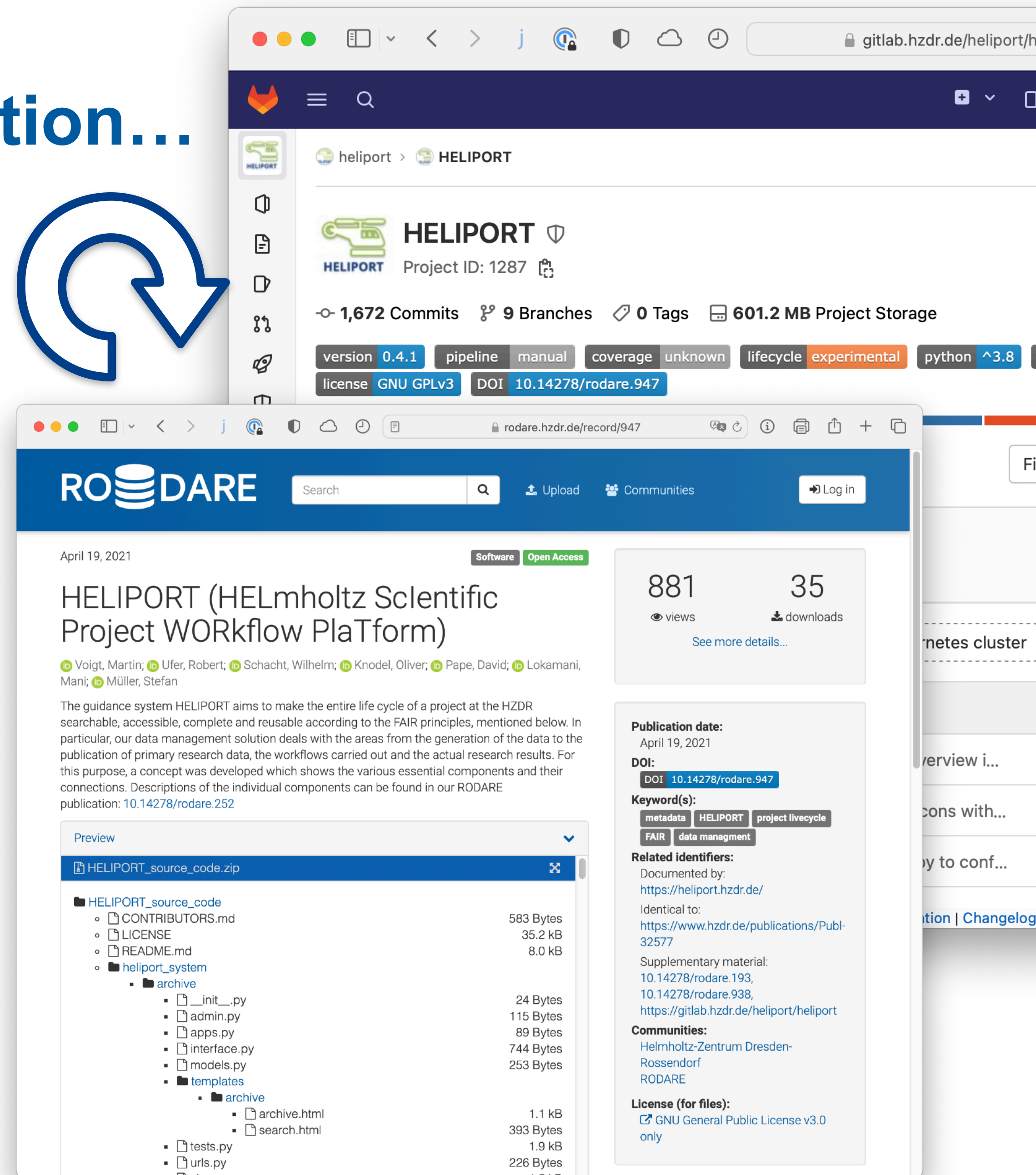# The Git Repository is not a Publication…

**Software Repository**

— Software is typically available (not *published*) in version control systems with open or restricted access:

— We need workflows or methods to publish software and data to ensure long-term *availability* and to meet the FAIR principles.

**Software Publication**

— Software must be cited in a similar way to scientific publications.

— Common data repositories (e.g. institutional, domain-specific, Zenodo) support typically the publication type *software*.

— With an additional *software publication* we can cite specific versions of a software including rich metadata:

  • Title, authors (including ORCIDs), Abstract, license, …

  • Related Identifiers to link additional resources.

  • Typically it would not be practical to link *all* scientific datasets produced or analyzed by the software to the software publication.
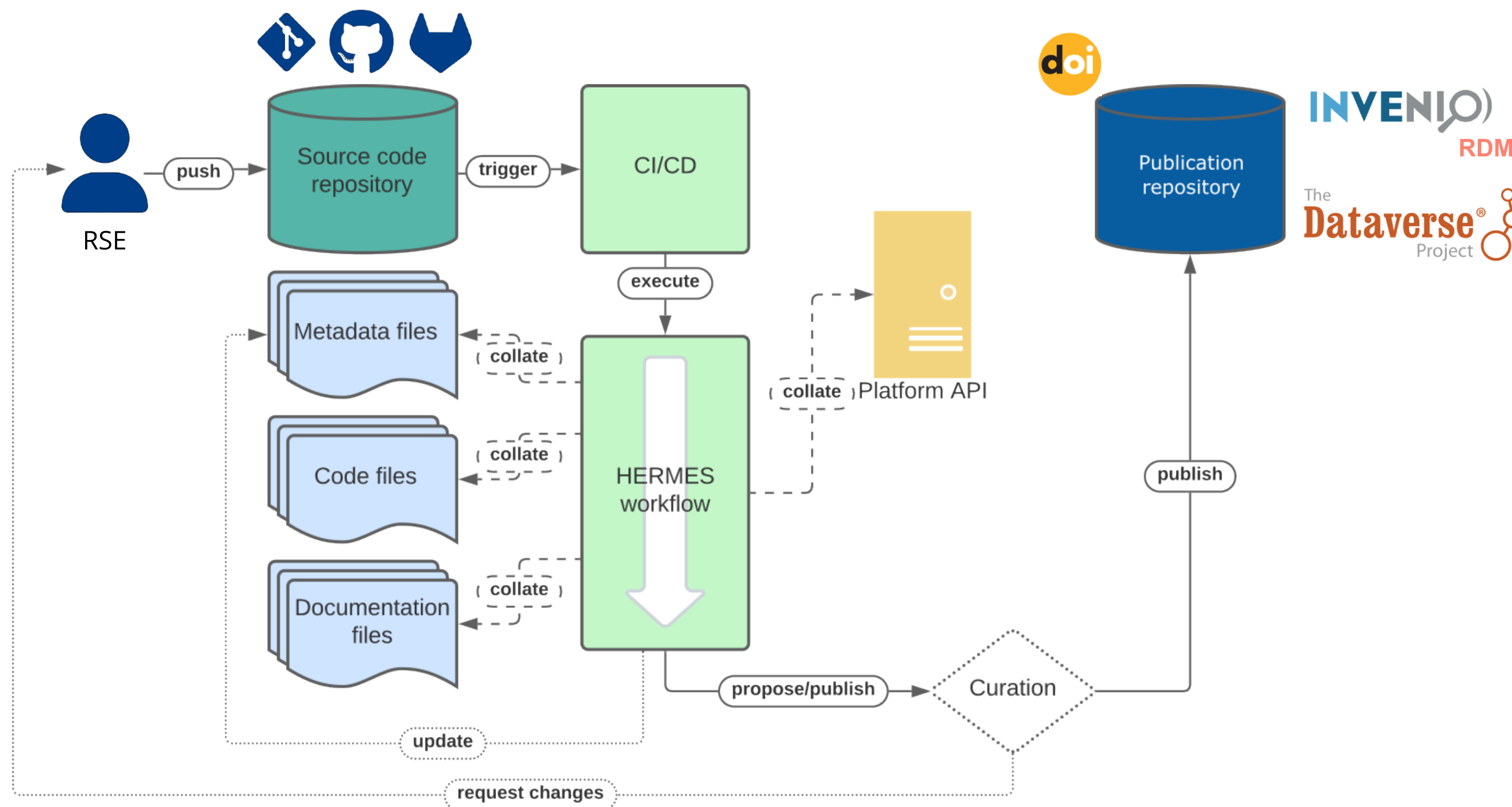
**Repository Structure**

— Software and date mixen in one repository,

— Separation between data and software repositories,

— Something in between…

# The HERMES Project: Automated Software Publication Workflow

— A simple and transparent software publication workflow for open and closed access software can be a platform for an understandable science.

— The metadata harvesting is essential to create a findable software publication.



HERMES

project.software-metadata.pub

github.com/hermes-hmc

team@software-metadata.pub

<HMC> HELMHOLTZ METADATA COLLABORATION

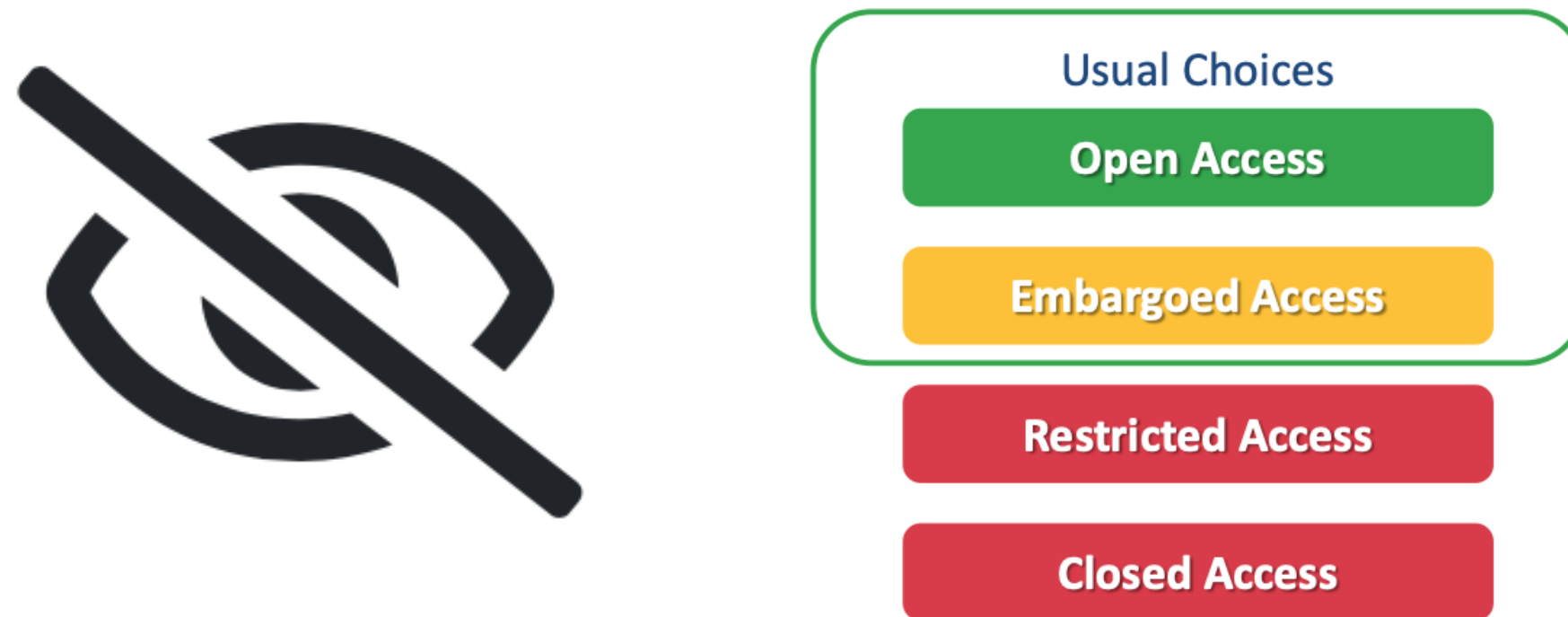DLR   JÜLICH Forschungszentrum   HZDR HELMHOLTZ ZENTRUM DRESDEN ROSSENDORF

- 07/2021 – 06/2023
- Aim: Support RSEs in automatedly publishing their software with rich metadata

# Data (and Software) Publication Repositories

— The data can be published in the same or a different repository as the software (possibly there is a domain-specific repository for the scientific data).

— A dataset (at least the metadata) should always be published to provide sustainable scientific evidence.

— The data itself can also be published under restricted access.*

*Nevertheless it fulfills the FAIR principles, because the steps to access the data are documented.



— The data publication should reference:

• The software repository used to create or analyse the dataset.

• The scientific publication based on the datasets.

• The instrument or facility where the data was generated...

# Instrument DOIs and Landing Page

— For data publications we have the field *related identifiers*, were we can refer research facilities and instruments.

— Therefore, we plan to assign DOIs to instruments and provide DataCite records [2] and additional metadata on public landing pages.

— Components of the landing pages:

  • Mandatory: DOI, name, description, contacts, scope, location, ROR, device type.

  • Optional: Image, layout, sub-facilities, additional resources (JLSRF publication, internal website, ...) and the latest publications.

  • Citation export to BibTex, JSON, ...

[2] Bunakov, Vasily, Krahl, Rolf, Matthews, Brian, Vizcaino, Noeland Vukolov, Andrey, "Advanced infrastructure for PIDs in Photon and Neutron RIs", ExPaNDS project deliverable D2.5, Zenodo, Mar. 2022. doi: 10.5281/zenodo.5905351.

# Digital Objects and Handles Enable Long-term Sustainability

— At the HZDR, we use DOIs for resources containing the whole set of bibliographic metadata:

- Scientific articles,
- Published datasets and software,
- Instruments.

— Other identifiers included in the metadata are ORCID and ROR available in our internal databases for almost every scientist at HZDR.

— Further digital objects can request PIDs from our Handle server.

— The digital objects in our ecosystem can be correlated with each to create a comprehensible experiment providing data provenance.
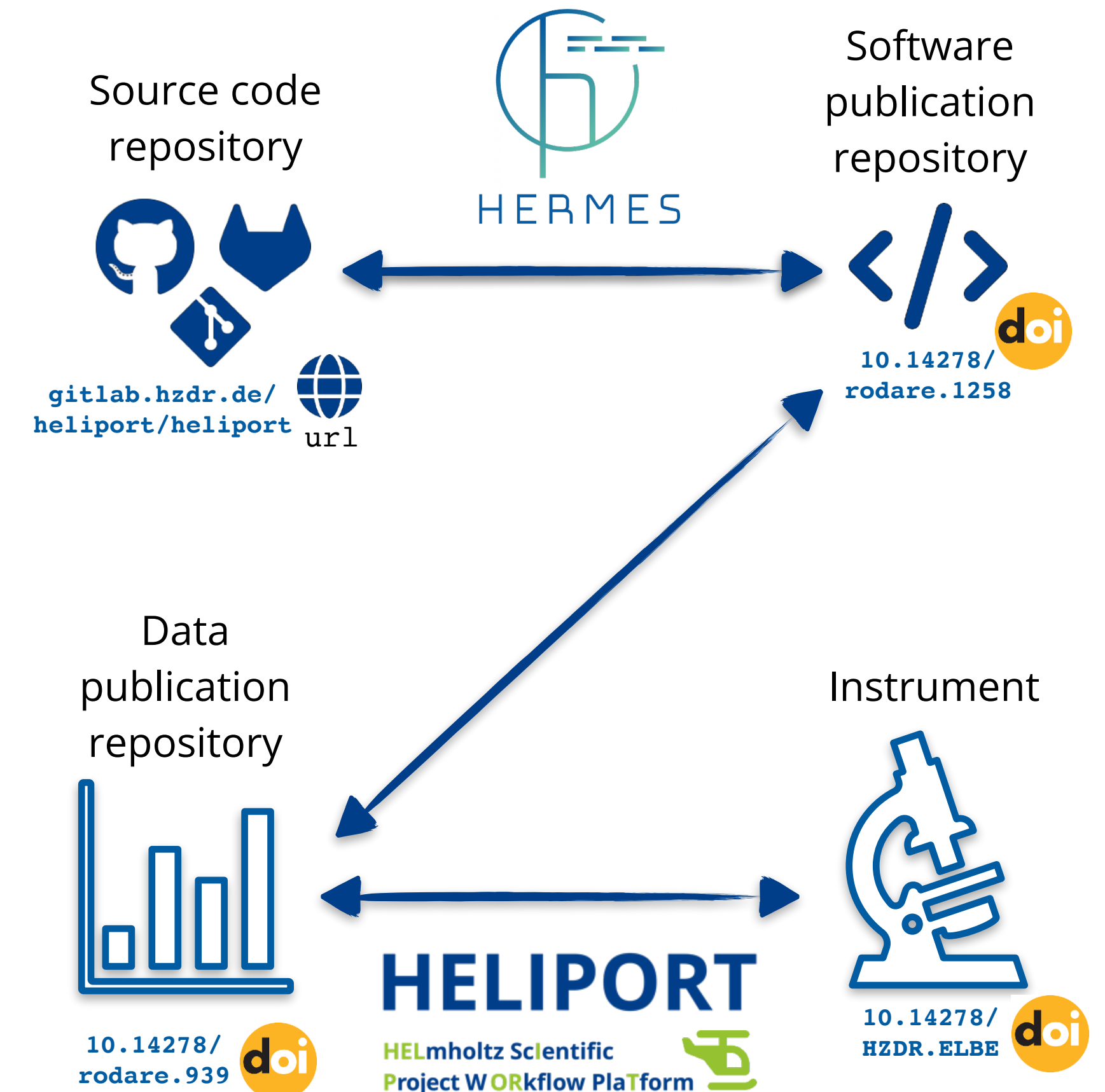
# Linking Data, Text and Research Software Together

**Software:**

I. HERMES can extract the metadata provided by Github or GitLab.

II. A software release can trigger a pipeline that initialises a publication with DOI based on the available (and third-party) metadata.

III. In a subsequent step, the DOI is added to the Readme file in the Git repository and the cross-linking is completed.
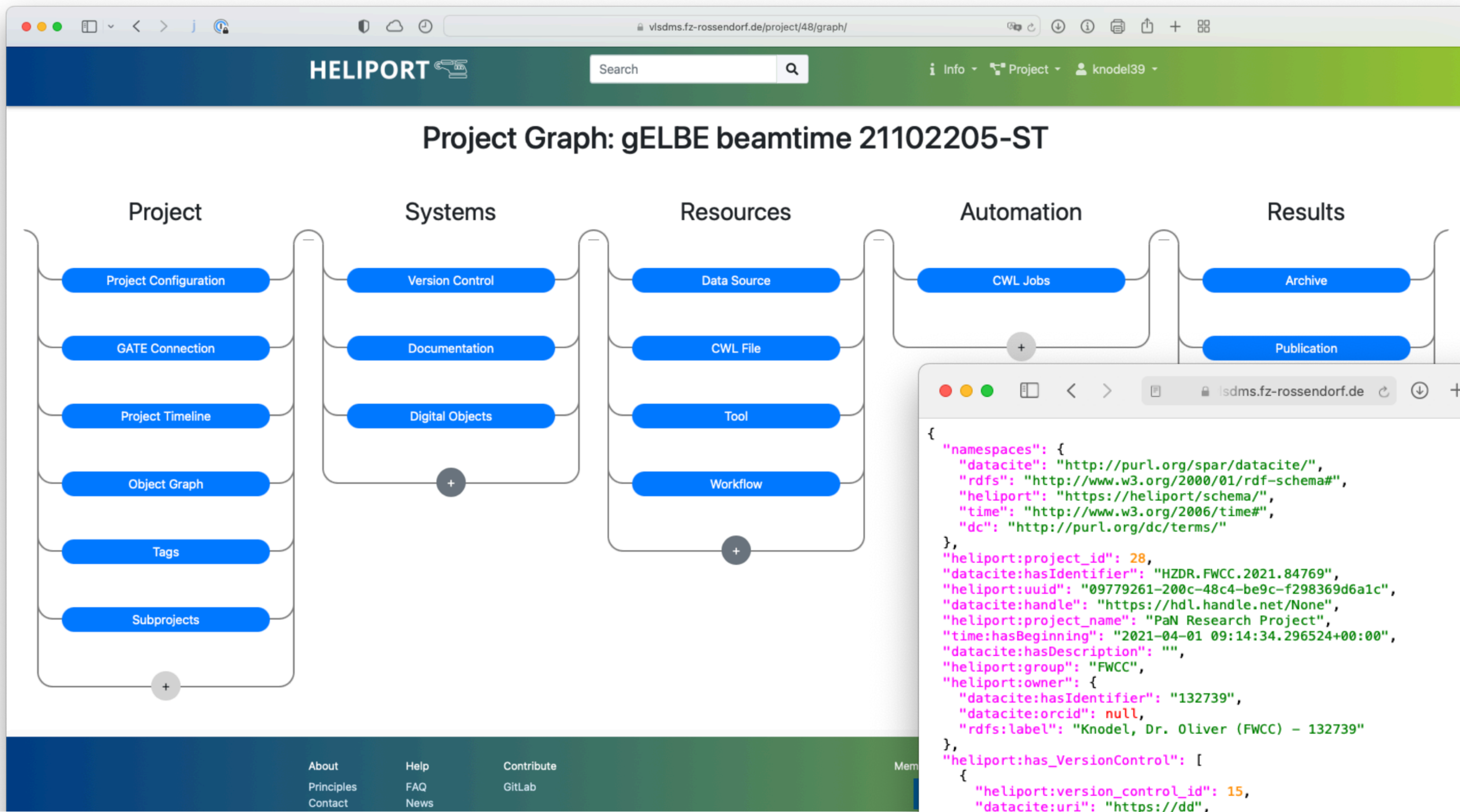
**Data:**

I. After data collection or processing, a pipeline can start collecting metadata from a proposal system or other related services.

II. The metadata and information from a computational workflow can be used to create a data publication with references to a specific software version (DOI) and the instrument where the data was taken.



Source code repository

gitlab.hzdr.de/
heliport/heliport url

HERMES

Software publication repository

10.14278/
rodare.1258

Data publication repository

10.14278/
rodare.939

HELIPORT

HELmholtz ScIentific
Project WORkflow PlaTform

Instrument

10.14278/
HZDR.ELBE

DRESDEN concept    HZDR

# Overview of the Project Resources from a Higher Level

> The HELIPORT project aims at developing a platform which accommodates the **complete life cycle** of a scientific project and links all corresponding programs, systems and workflows to create a more **FAIR** and comprehensible project description.



**HELIPORT**
**HEL**mholtz **ScI**entific
**P**roject **W**ORkflow **PlaT**form

🌐 heliport.hzdr.de

 codebase.helmholtz.de/heliport

✉ heliport@hzdr.de

<HMC> HELMHOLTZ METADATA COLLABORATION

HZDR HELMHOLTZ ZENTRUM DRESDEN ROSSENDORF
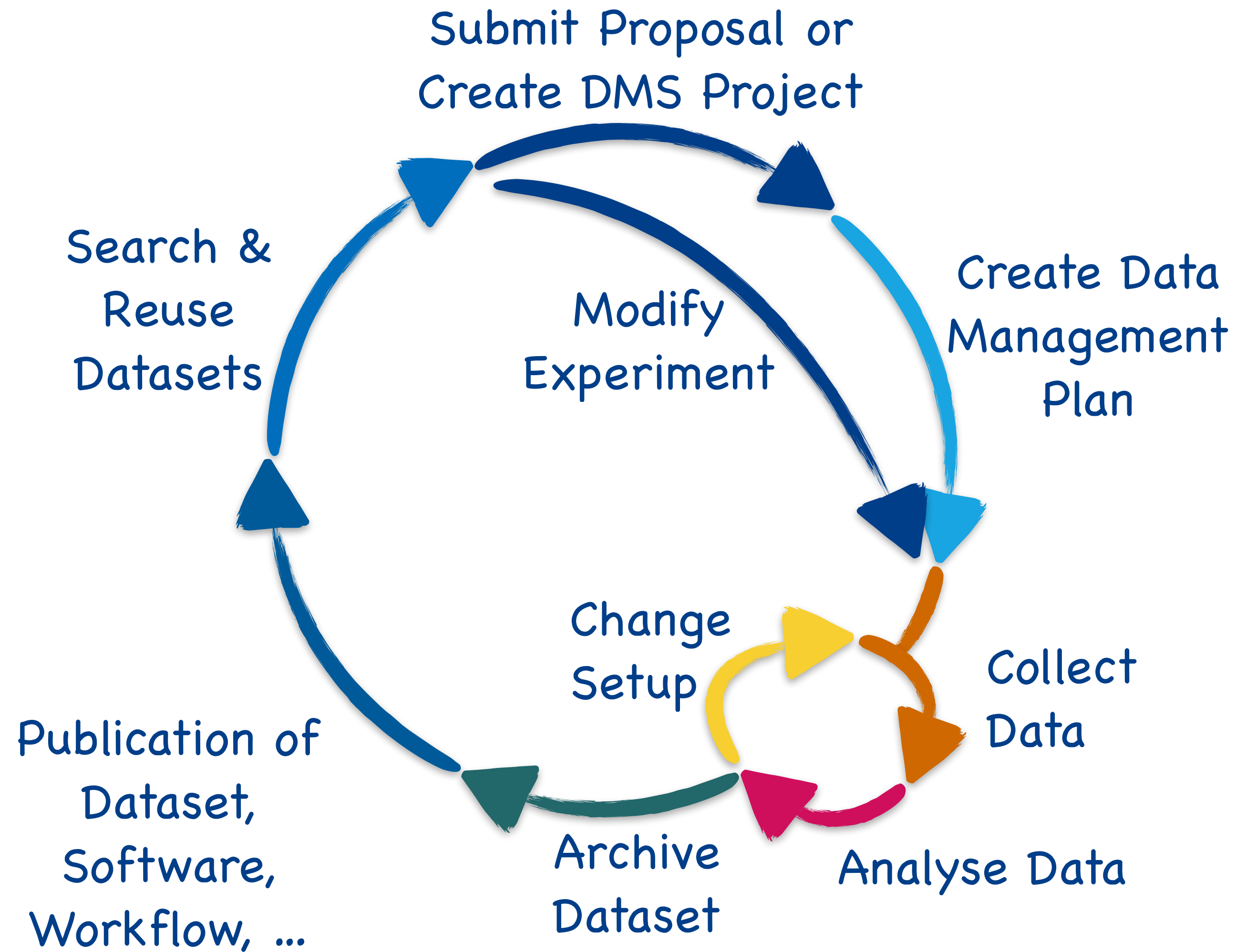
HI JENA Helmholtz Institute Jena

JÜLICH FORSCHUNGSZENTRUM

- 07/2021 – 06/2023
- Aim: Collect every system, service or digital product of a research project in an uniform metadata package.

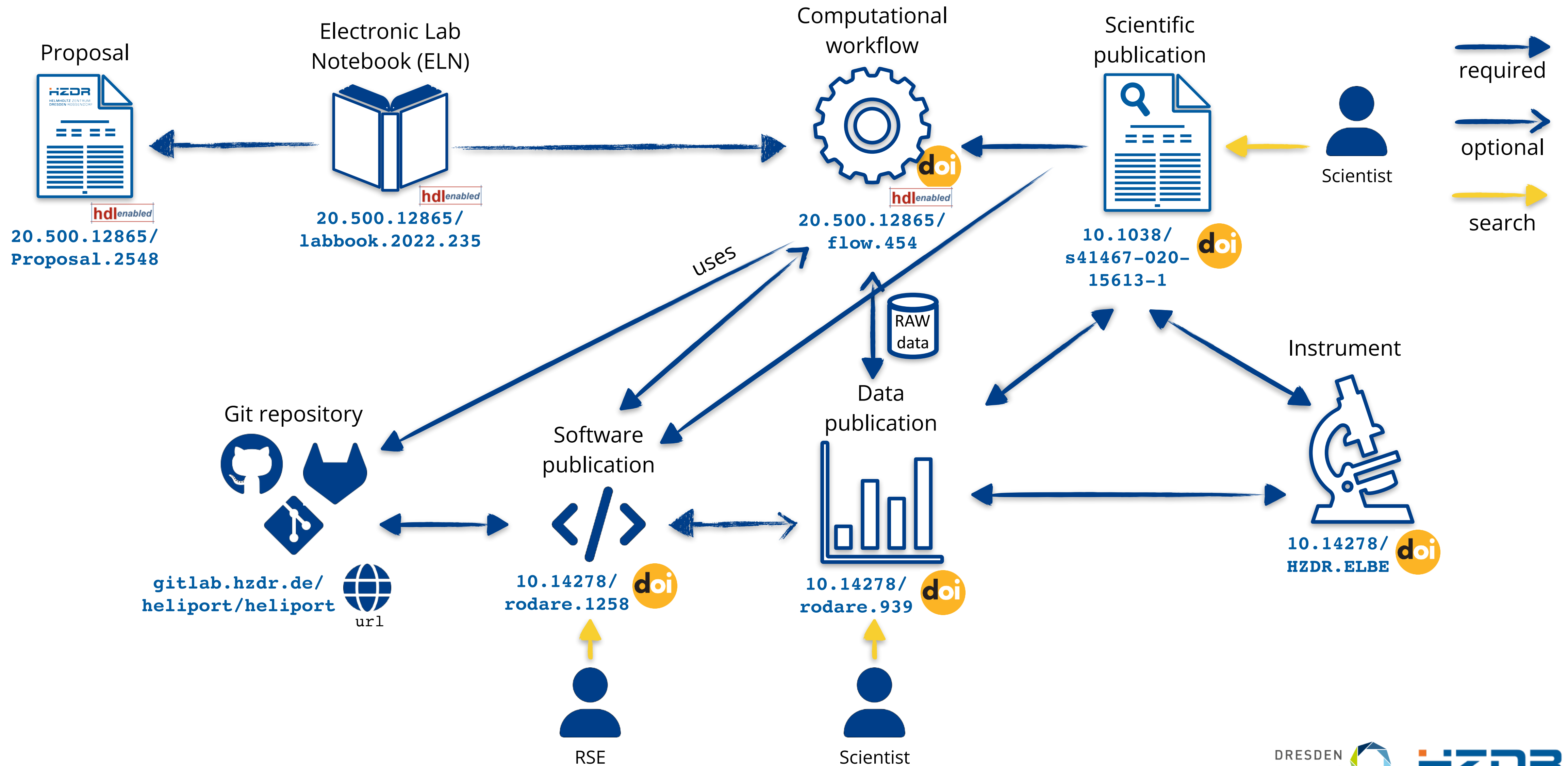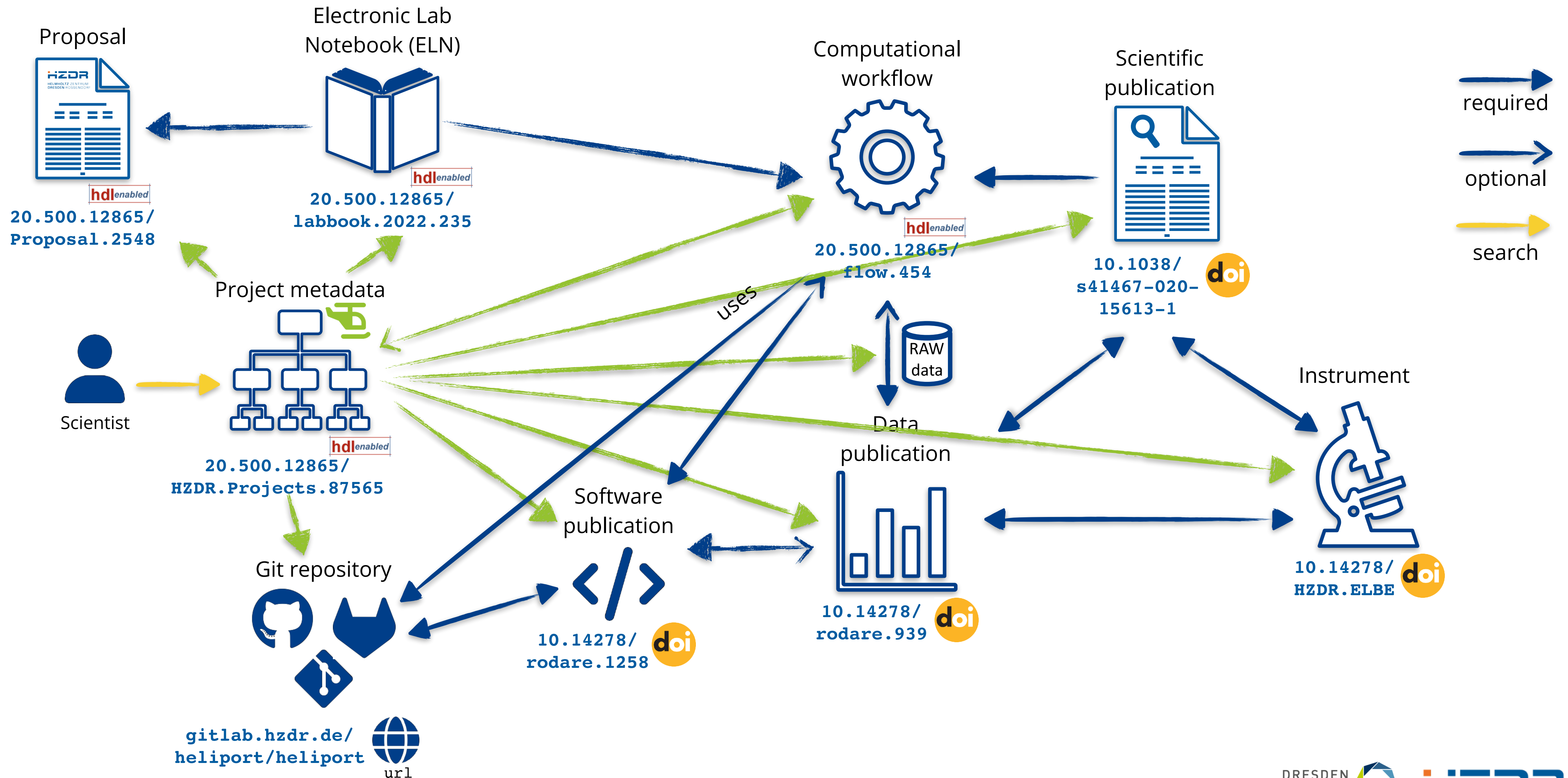DRESDEN concept

HZDR

# Our Challenge: An End-to-End Digital Data Lifecycle

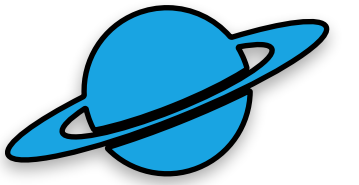# Top-Level View of the Interlinked Digital Objects of an Experiment at the HZDR

# Top-Level View of the Interlinked Digital Objects of an Experiment at the HZDR

# Conclusions and Outlook

**Conclusion:**

— For an interlinked ecosystem, it is necessary to consider different entry points for the provision of metadata.

— The cross-linking of the services and systems is unavoidable to enable comprehensible science.

➡ Automated pipelines and workflows are the key to exchange metadata and support scientists and RSEs.

**Status:**

— We provide DOIs for software and data (instrument DOIs are work in progress),

— Handles can be created for all types of digital objects.

➡ With HELIPORT and HERMES, we develop systems that automate the exchange of metadata between internal and external systems and services.