# Connecting Processes to Data via Meta-Data
## Perspectives of Data collection and Exchange
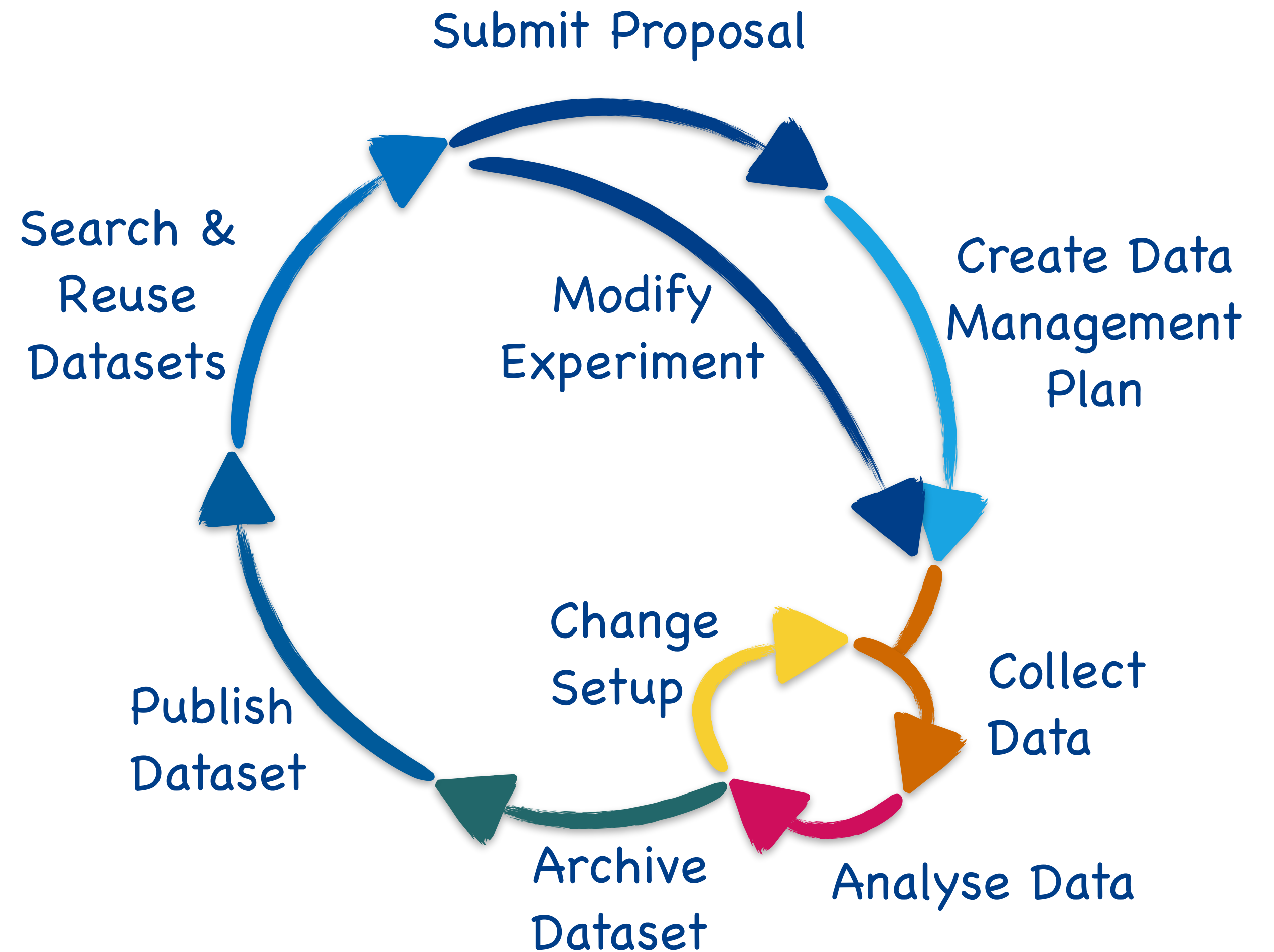
Thrill WP5 ML Workshop, HZDR, February 2024

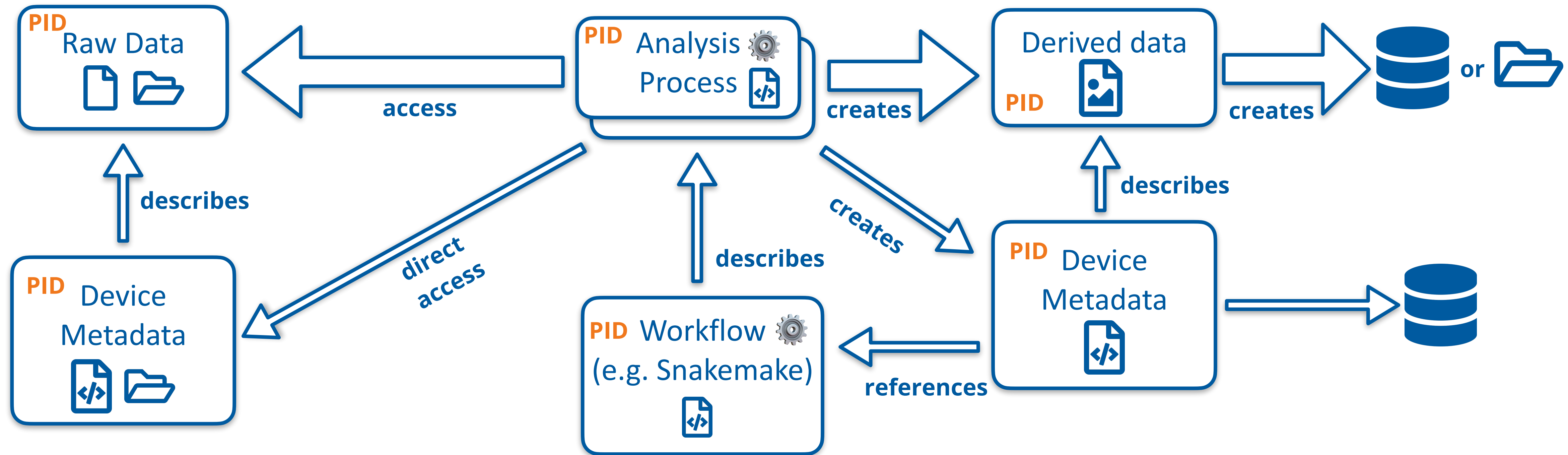**Oliver Knodel** // contact: o.knodel@hzdr.de

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

— We support many steps of our different research experiment (matter, energy and health) with tools:

- Electronic lab notebook (**E-Logbook**),
- Interactive analysis,
- Publication of datasets,
- Scientific workflow management,
- Handle generation and management.

— A uniform and smooth access to and between all services and **processes** in a digital ecosystem is necessary.

— The description and interconnection between all linked resources through metadata is essential to create a **comprehensible** and **FAIR** experiment.
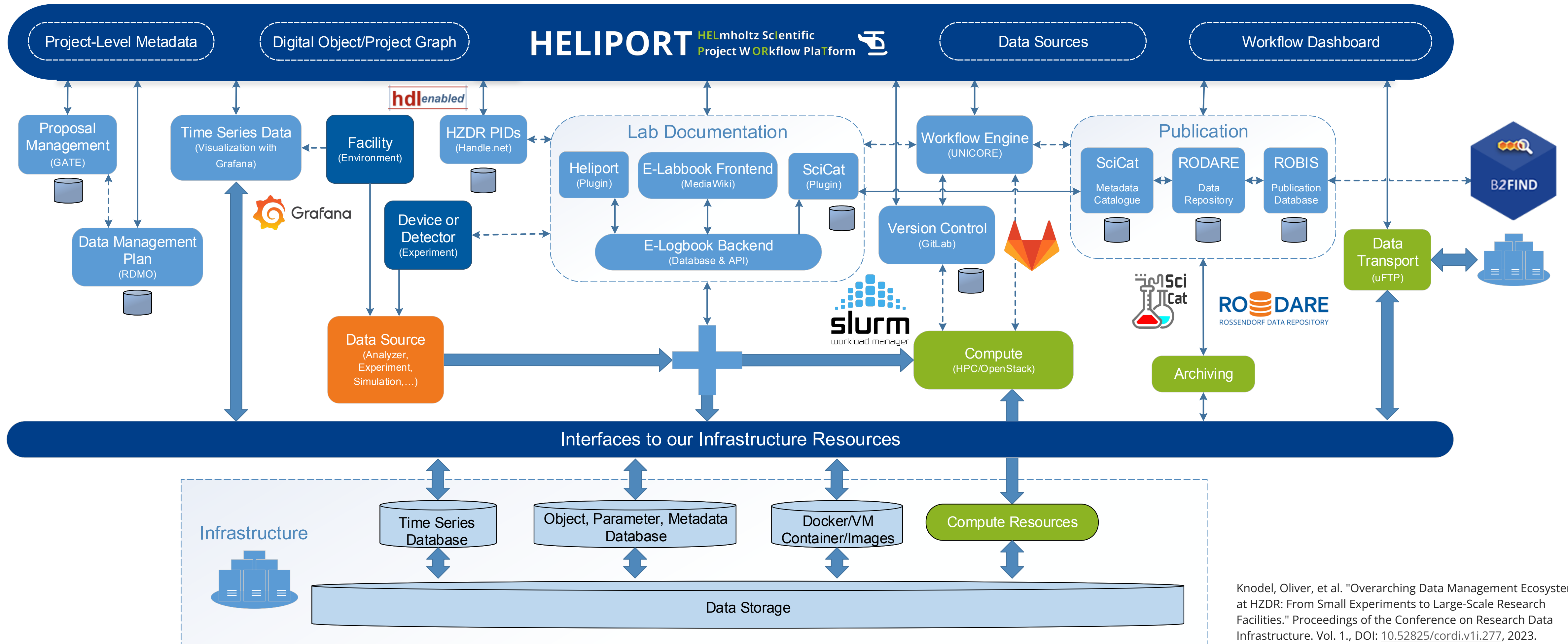
HZDR

# Connecting Processes to Data via Meta-Data



— The entry point is usually the raw data set that is generated from the test setup (e.g. detector, camera)

— Each data product should be described by a (standardised) metadata schema

— Process (or workflows) can access the RAW data via the corresponding metadata (with associated identifier)

— Workflows themselves should be described regarding the FAIR principles

— Derived data should contain additional descriptive metadata and an exchangeable data format (e.g. HELPMI)

— **PIDs** (e.g. Handles, DOIS) are essential to provide persistent identifier for data and metadata
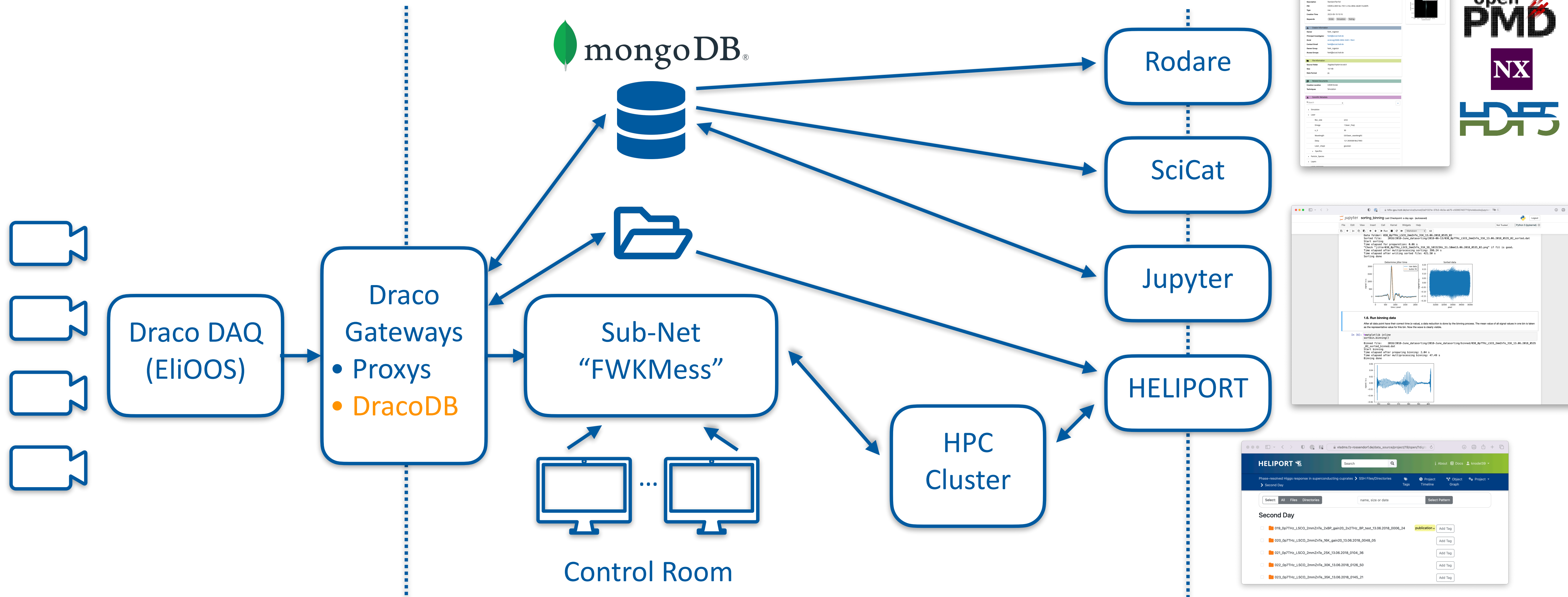
# Digital Research Landscapes at HZDR

HELIPORT — HELmholtz Scientific Project WORkflow PlaTform

Project-Level Metadata · Digital Object/Project Graph · Data Sources · Workflow Dashboard

Proposal Management (GATE) · Time Series Data (Visualization with Grafana) · Facility (Environment) · HZDR PIDs (Handle.net) · Lab Documentation · Workflow Engine (UNICORE) · Publication

Heliport (Plugin) · E-Labbook Frontend (MediaWiki) · SciCat (Plugin) · SciCat (Metadata Catalogue) · RODARE (Data Repository) · ROBIS (Publication Database)

Data Management Plan (RDMO) · Device or Detector (Experiment) · E-Logbook Backend (Database & API) · Version Control (GitLab)

Data Source (Analyzer, Experiment, Simulation, …) · Compute (HPC/OpenStack) · Archiving · Data Transport (uFTP)

Interfaces to our Infrastructure Resources

Infrastructure — Time Series Database · Object, Parameter, Metadata Database · Docker/VM Container/Images · Compute Resources

Data Storage

Knodel, Oliver, et al. "Overarching Data Management Ecosystem at HZDR: From Small Experiments to Large-Scale Research Facilities." Proceedings of the Conference on Research Data Infrastructure. Vol. 1., DOI: 10.52825/cordi.v1i.277, 2023.
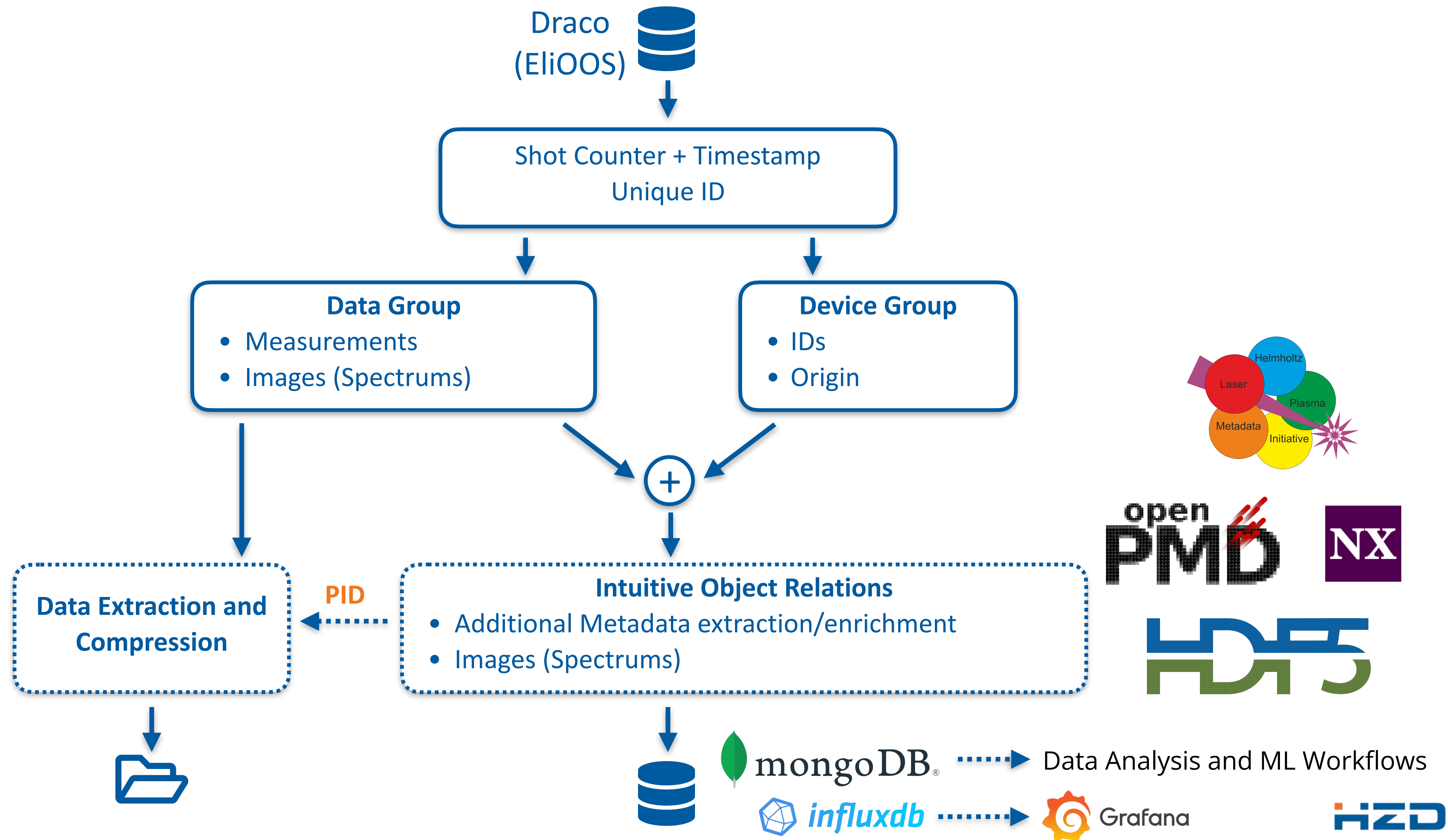
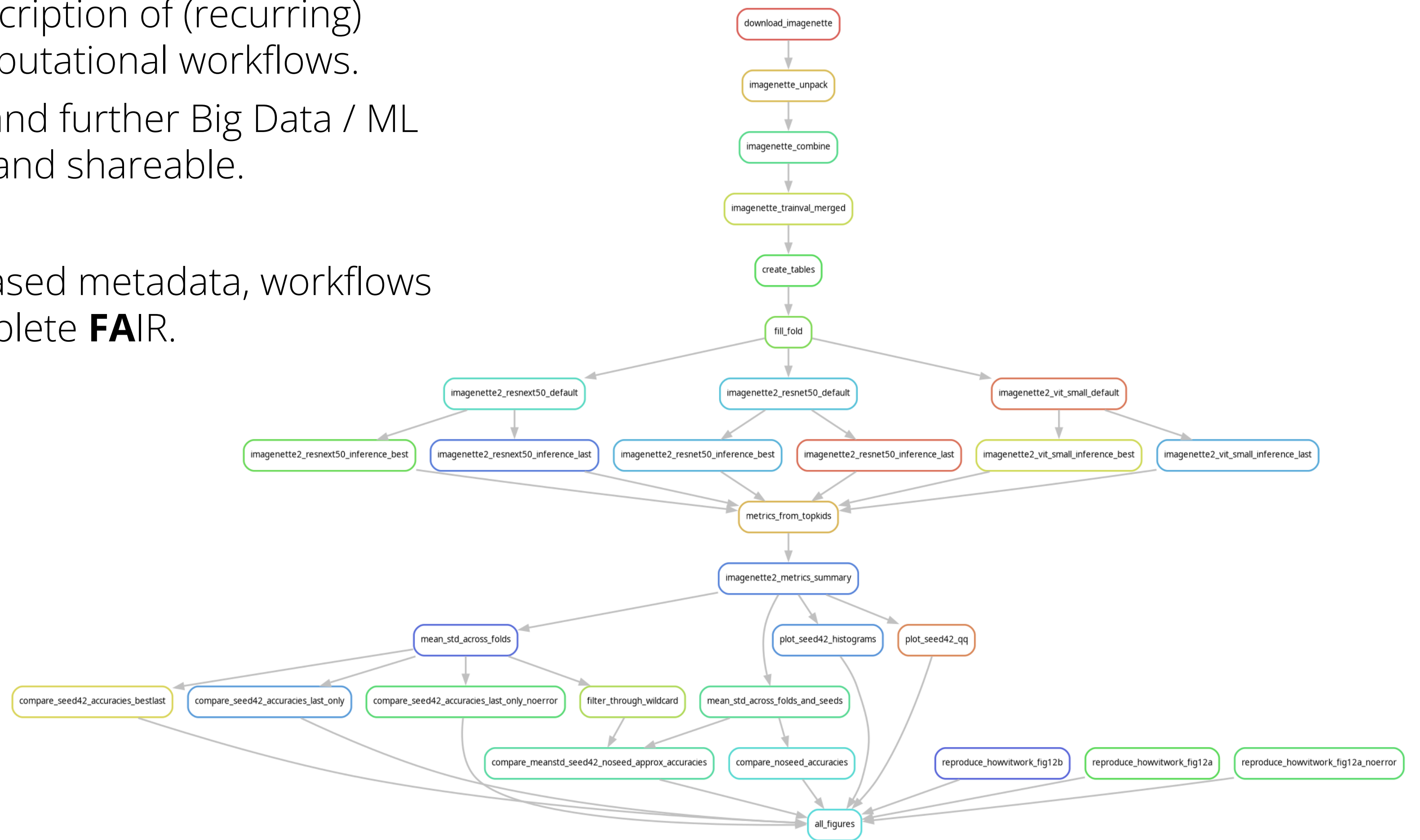# Draco Data Flow for Advanced Data Management

# Extended Draco Data Pipeline (to fully take advantage of the HZDR infrastructure)



Draco
(EliOOS)

Shot Counter + Timestamp
Unique ID

**Data Group**
- Measurements
- Images (Spectrums)

**Device Group**
- IDs
- Origin

(+)

**Data Extraction and Compression**

PID

**Intuitive Object Relations**
- Additional Metadata extraction/enrichment
- Images (Spectrums)

mongoDB

influxdb

Grafana

Data Analysis and ML Workflows

# Computational Workflows (Processes) and Data Management (MetaData)

— In our HZDR infrastructure, the description of (recurring) work can be automatised with computational workflows.

— Workflows enable deeper insights and further Big Data / ML methods that are comprehensible and shareable.

— Workflows enable FA**IR**.

— With interchangeable, standards-based metadata, workflows can be used in different RIs to complete **FA**IR.



DAG of an ML Workflow in Snakemake based on "Machine Learning State-of-the-Art with Uncertainties"; DOI: 2204.05173

# Workflow Architecture at HZDR (in development)

— HELIPORT offers an infrastructure which permits the integration of various workflow languages and access modes to HPC infrastructures.

— The infrastructure keeps track of and collects the metadata and enables access to all resources involved.

— Next steps:

- Python library sending workflow information directly to HELIPORT,

- Provision of provenance information from Jupyter notebooks,

- Use case: **PIConGPU**

# HELIPORT
**HEL**mholtz Sc**I**entific
**P**roject W **OR**kflow Pla**T**form

> 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 using **APIs**.

**Project Members:**

HZDR
HELMHOLTZ ZENTRUM
DRESDEN ROSSENDORF

HI JENA
Helmholtz Institute Jena

JÜLICH
FORSCHUNGSZENTRUM

**Funded by:**

<HMC> HELMHOLTZ
Metadata
Collaboration





Knodel, Oliver, et al. "HELIPORT: A Portable Platform for {FAIR Workflow | Metadata | Scientific Project Lifecycle} Management and Everything" In Proceedings of the 4th International Workshop on Practical Reproducible Evaluation of Computer Systems (P-RECS '21). AAM. 2021. 10.1145/3456287.3465477.

ToDo: Metadata crosswalk to
**schema.org**
**ResearchProject**

# Metadata Catalogue SciCat and Data Repository RODARE (Draft)

**Curated Metadata Source**

**Public Metadata Catalogue**

**Subsequent Access to Data**

**ExperimentLogging app (ExL)**

**SciCat**

**RODARE**

**E-Logbook**

**HELIPORT**

**Metadata from Experiment/Simulation**

**Dataset**

**Filesystem**

**Tape Archive**

**Fully Automated Process for DRACO**

# Conclusions

— Access to data via metadata and PIDs is required to enable ML according to the FAIR criteria

— Connecting Systems, services and processes using APIs for Metadata exchange is essential.

— Metadata and data should follow data standards and schemas to allow exchange of research products and to provide **FAIR** and **comprehensible** research.

— The computational workflows are essential to automate recurring processes.