# HZDR Data Management Strategy

## Meeting at Leibniz Institute of Polymer Research Dresden (IPF)

November 2019

Oliver Knodel, Thomas Gruber and Stefan Müller // contact: o.knodel@hzdr.de

# HZDR – Facts and Figures

— **Member of the** HELMHOLTZ **Association**
RESEARCH FOR GRAND CHALLENGES

— **Employees  approx. 1,200**
   including about 350 scientists
   + 150 doctoral students
   as well as employees and guest
   scientists from more than **50** countries

— **Research Sites**
   **DRESDEN**
   Leipzig, Freiberg, Schenefeld near Hamburg (XFEL), Grenoble (FR)

Credits: Killig, DESY, ESRF/Ginter
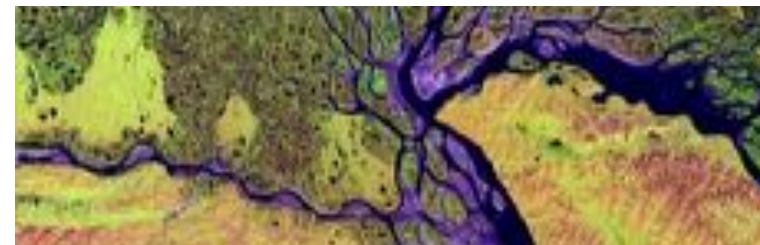
# HZDR – Location in Germany



GEOMAR Helmholtz-Zentrum für Ozeanforschung Kiel

Deutsches Elektronen-Synchrotron DESY

Helmholtz-Zentrum Geesthacht Zentrum für Material- und Küstenforschung (HZG)

Max-Delbrück-Centrum für Molekulare Medizin in der Helmholtz-Gemeinschaft (MDC)

Alfred-Wegener-Institut Helmholtz-Zentrum für Polar- und Meeresforschung (AWI)

Helmholtz-Zentrum für Infektionsforschung (HZI)

Helmholtz-Zentrum Berlin für Materialien und Energie (HZB)

Helmholtz-Zentrum Deutsches GeoForschungs-Zentrum GFZ

Helmholtz-Zentrum Dresden-Rossendorf (HZDR)

Forschungszentrum Jülich

Deutsches Zentrum für Luft- und Raumfahrt (DLR)

Helmholtz-Zentrum für Umweltforschung - UFZ

Deutsches Zentrum für Neurodegenerative Erkrankungen (DZNE)

GSI Helmholtzzentrum für Schwerionenforschung

Helmholtz-Zentrum für Informationssicherheit – CISPA

Deutsches Krebsforschungs-zentrum (DKFZ)

Karlsruher Institut für Technologie (KIT)

Max-Planck-Institut für Plasma-physik (IPP) (assoziiertes Mitglied)

Helmholtz-Zentrum München - Deutsches Forschungszentrum für Gesundheit und Umwelt

HELMHOLTZ ZENTRUM DRESDEN ROSSENDORF

HELMHOLTZ
RESEARCH FOR GRAND CHALLENGES

DRESDEN concept

# The six research fields of the Helmholtz Association

**ENERGY**

EARTH AND ENVIRONMENT

**HEALTH**

AERONAUTICS, SPACE AND TRANSPORT

KEY TECHNOLOGIES

**MATTER**

Dr.-Ing. Oliver Knodel | Department of Information Services and Computing | Computational Science Group | www.hzdr.de

# Large Research Infrastructures



## ELBE – Center for High-Power Radiation Sources

— Electron accelerator ELBE feeds free-electron lasers FELBE & THz source TELBE.

— Generates positrons, protons and neutrons as well as X-ray and gamma radiation.

— High-intensity lasers (1 Petawatt) **DRACO** and **PENELOPE** (under construction).

## Dresden High Magnetic Field Laboratory (HLD)

— Nanoscale Producing Europe's highest pulsed magnetic fields.

## Ion Beam Center (IBC)

— Nanoscale surface analysis and modification.



Credits: Bierstedt, Killig (2 x)

Dr.-Ing. Oliver Knodel | Department of Information Services and Computing | Computational Science Group | www.hzdr.de

# One of Our Scientific Computing Challenges

**Main Challenge: Pre-/post-processing and archiving of research data**

— Filter and compress measured or simulated data at the *edge* of the datacenter.

— Accelerate compute-intensive tasks with dedicated low-latency (e.g. FPGAs), high-bandwidth (e.g. GPUs) hardware.

**Heterogenous Data Center**

— Many research questions require compute intensive deep learning approaches suitable for our HPC Cluster with GPUs (and FPGAs).

— In the End the research data is located in the data centre anyway.

**Prototyping and Continuous Integration (CI)**

— We support scientific applications and workflows to improve data processing and even the full software lifecycle.

— Custom CPU, GPU and FPGA applications have to be tested and verified with every development cycle using CI.

**Generate Data**



T-Elbe: ~35 GB/s (24/7)

**datacenter**

**Store Data**
Filter & Compression

PCI-Express: ~6 GByte/s

**Post-Process Data**
Deep learning & Analyse

I/O per node: ~30 MByte/s

**Development**
Prototyping & CI

Dr.-Ing. Oliver Knodel | Department of Information Services and Computing | Computational Science Group | www.hzdr.de

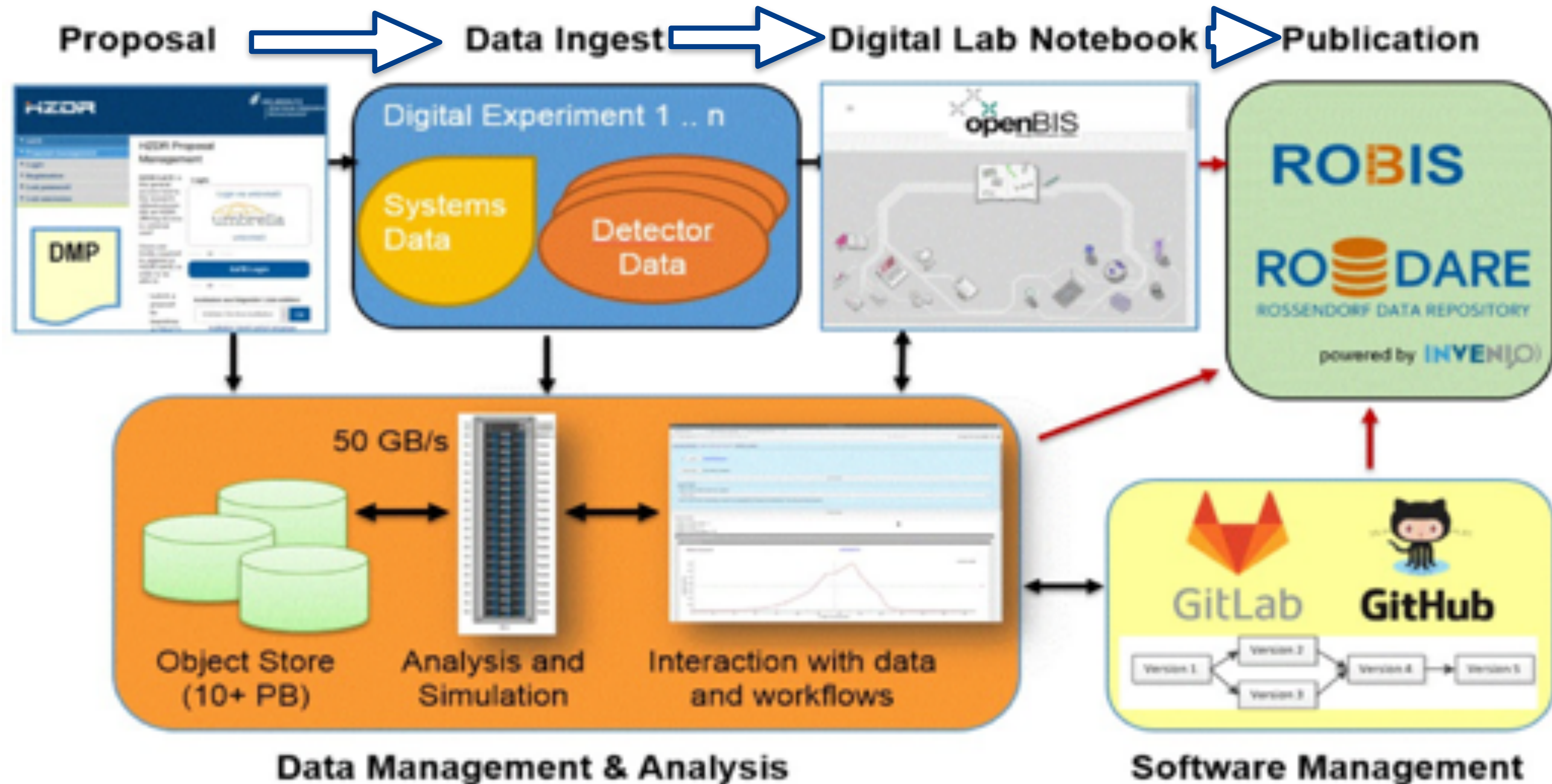# The other Challenge: An End-to-End Digital Data Lifecycle

— Data/Meta-Data standards are the key

— Support all stages with tools:
  • electronic lab books,
  • interactive analysis,
  • automated publication,
  • workflow management.

— Get the data as early as possible into the data center.



Research Data Management Lifecycle taken from:
https://guides.library.ucsc.edu/datamanagement

# HZDR Data Management Strategy

# Electronic Lab Books for Better Meta-Data Management

— Long Evaluation Phase:

— Result:

  • **OpenBis** for structured Lab-Data.



  • **MediaWiki** for more free-form Documentation.



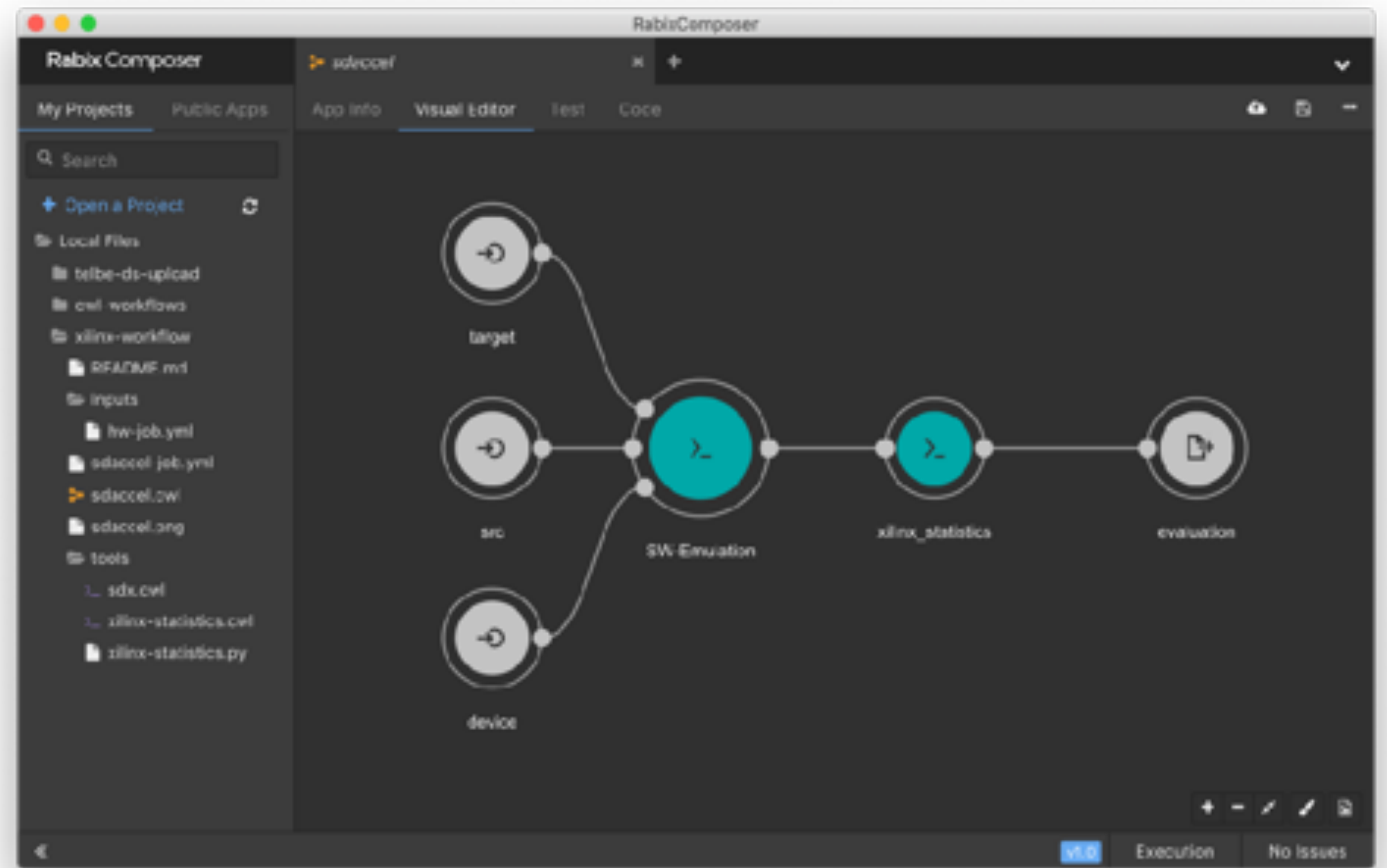— Both are necessary to meet the requirements of the experiments at HZDR.

Dr.-Ing. Oliver Knodel | Department of Information Services and Computing | Computational Science Group | www.hzdr.de

# Workflow Engine for Scientific Workflows

— The execution of scientific workflows must be:

- Comprehensibly
- Archivable,
- Reusable,
- Reproducible,
- Publishable.

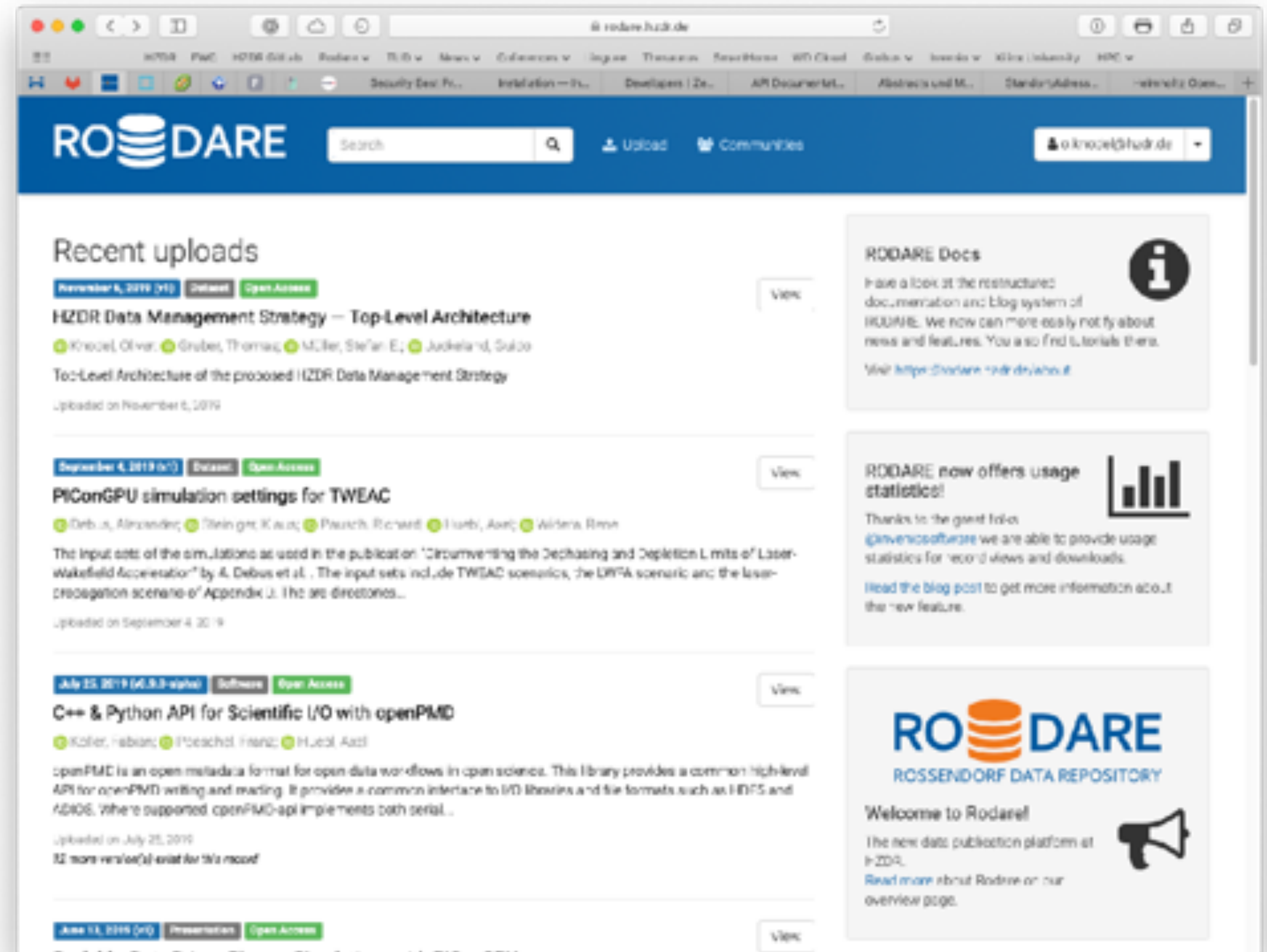— Based on an open standard: **OWL** or **WDL**

— Our evaluation is still in progress:

- **Reana** from Cern,
- Model-based approach from University **Turin**,
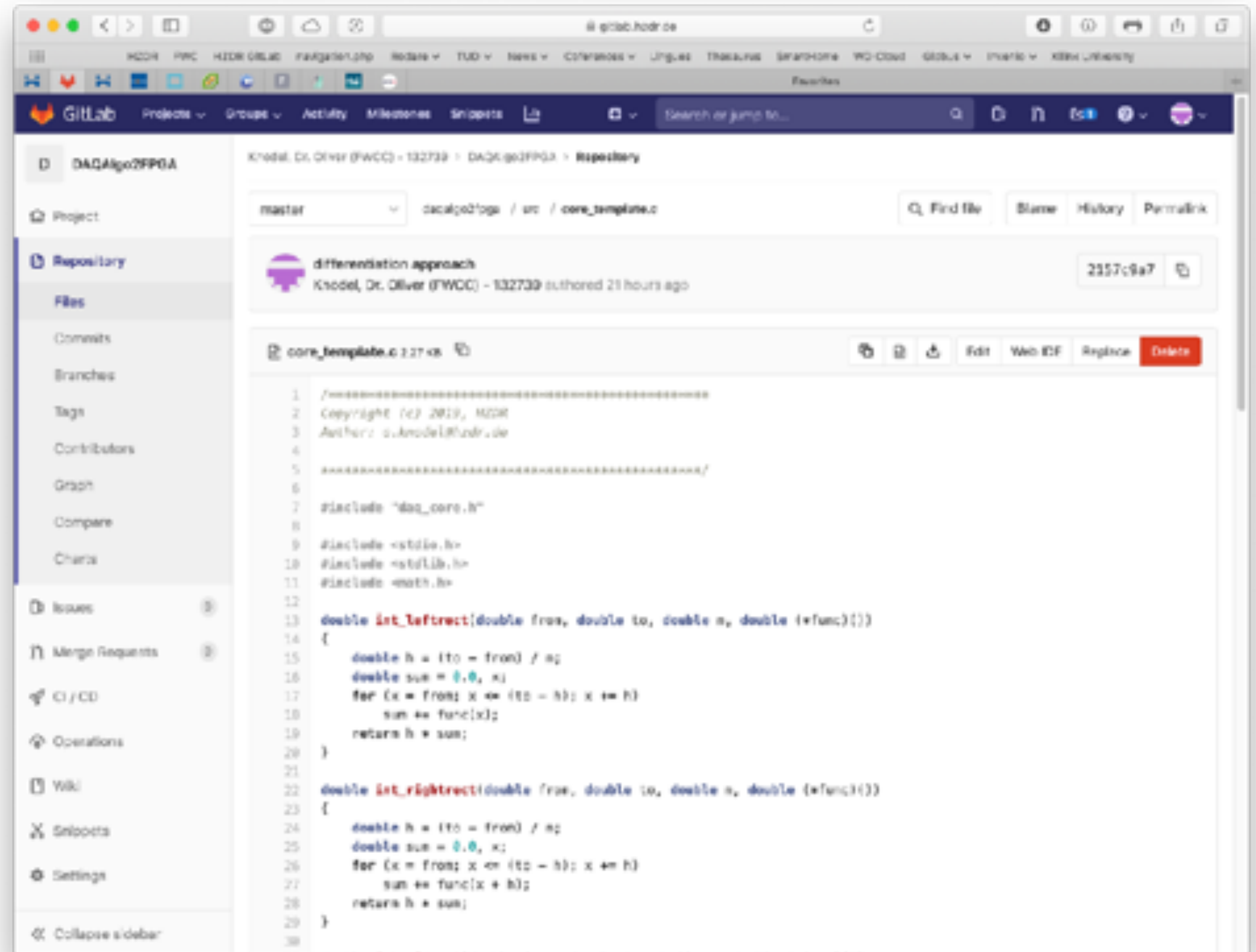- **Knime** as stand-alone Product.

Dr.-Ing. Oliver Knodel | Department of Information Services and Computing | Computational Science Group | www.hzdr.de

# RODARE - The HZDR Publication Platform

— Based on **Invenio** from CERN

— Highly modular and proven

- Own contributions:
- Shibboleth authenticator
- SFTP File Browser/Uploader
- Bittorrent Downloader
- GitLab-Integration

— HZDR is Part of CERN Community Collaboration Project

Dr.-Ing. Oliver Knodel | Department of Information Services and Computing | Computational Science Group | www.hzdr.de
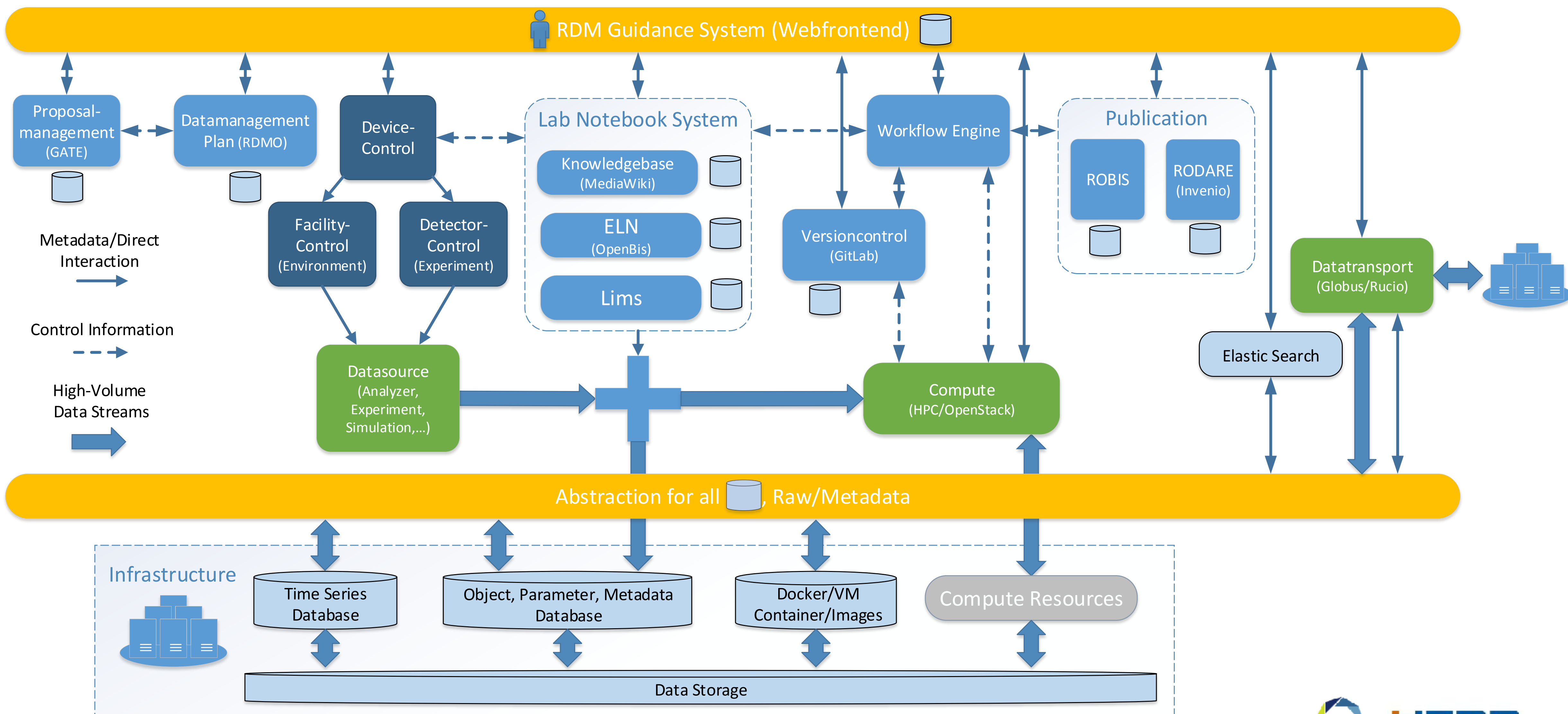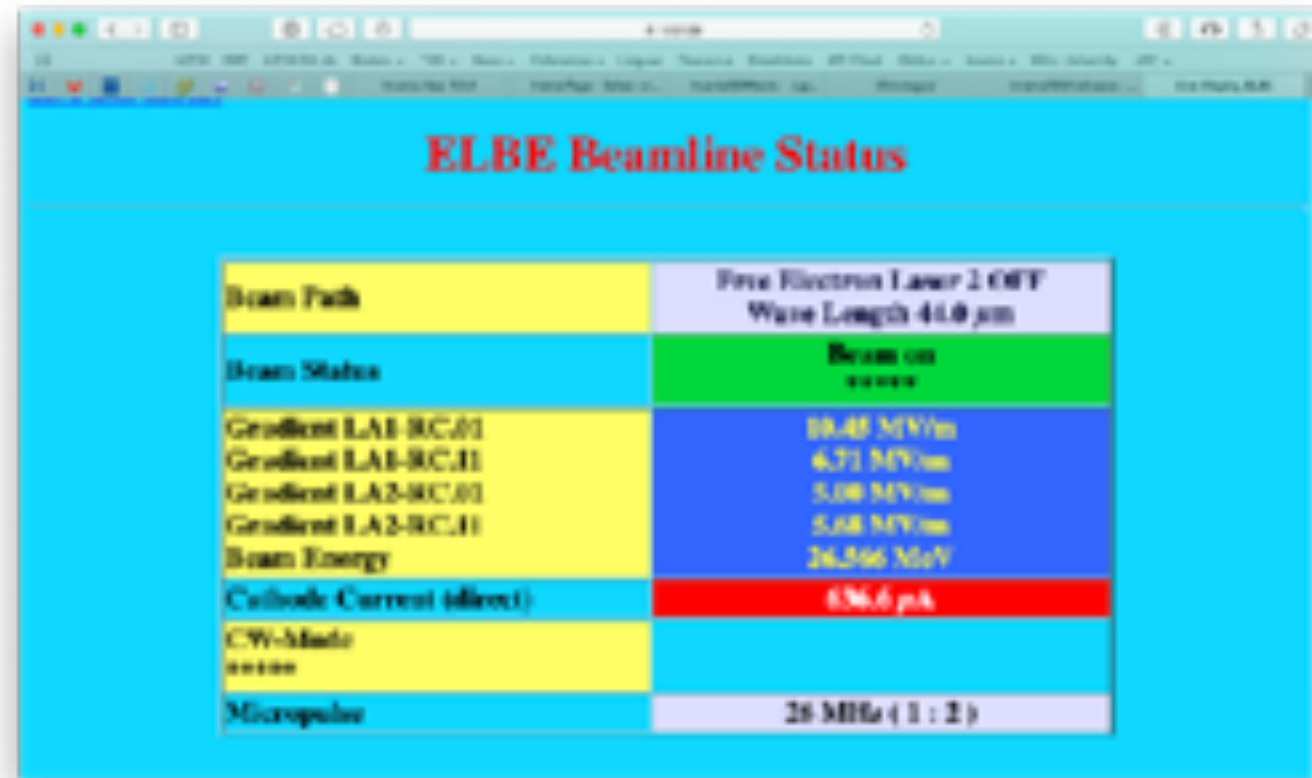
# IT-Hosted Services for Collaborative Work

— The **Department of Information Services and Computing** supports all the institutes as well as external users with a wide range of IT-services:
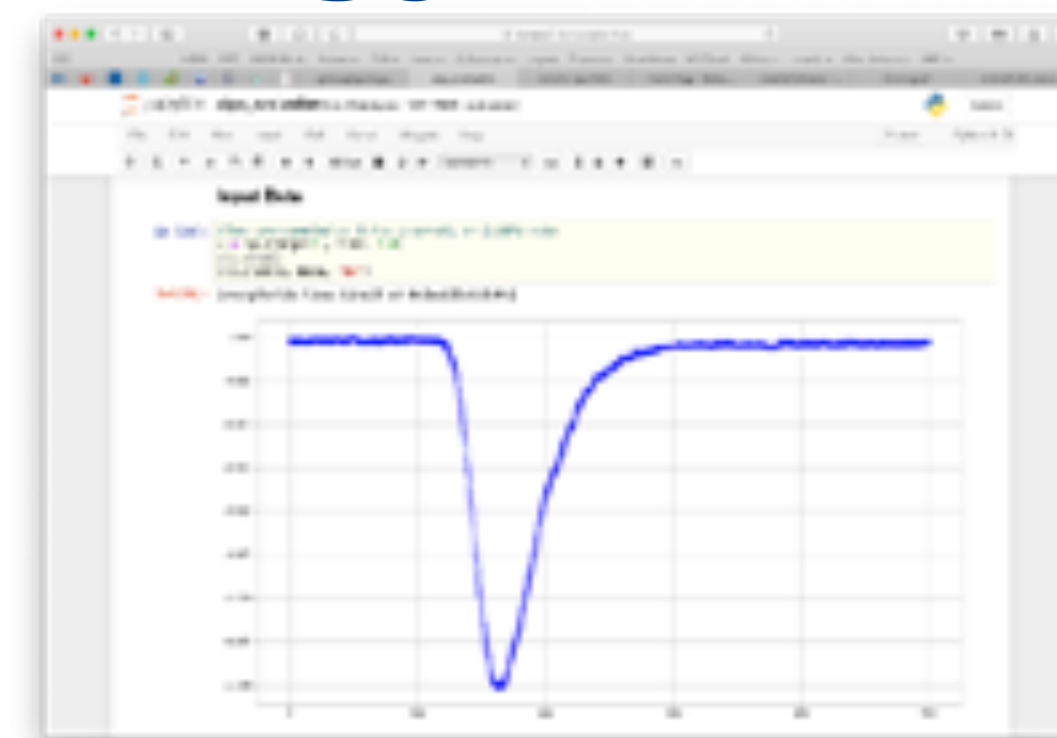
Member of the Helmholtz Association

Dr.-Ing. Oliver Knodel | Department of Information Services and Computing | Computational Science Group | www.hzdr.de

# HZDR Data Management Strategy — Top Level Architecture

Member of the Helmholtz Association

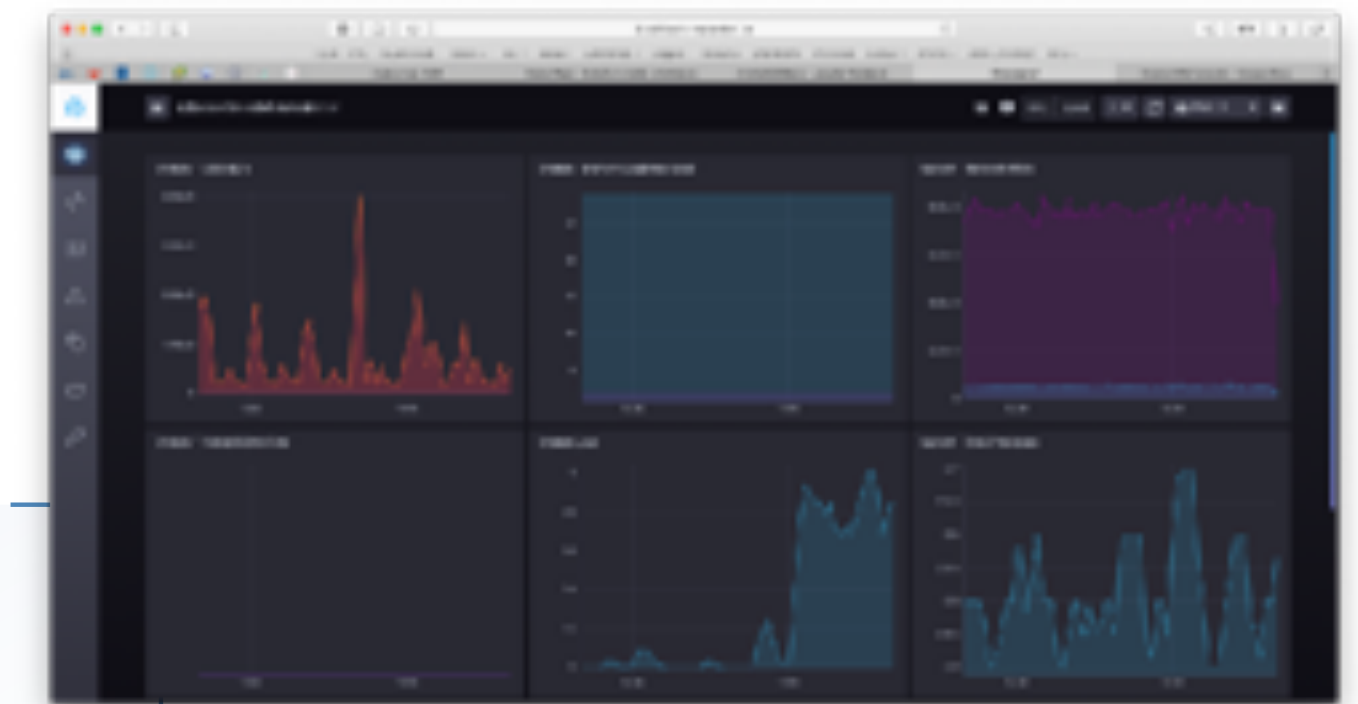Dr.-Ing. Oliver Knodel | Department of Information Services and Computing | Computational Science Group | www.hzdr.de

# HZDR Data Management Strategy — Possible Dataflow


Live Beamline Monitor


RAW Data


Time-Series Data


Structured (Meta)data

## ELBE Facility

Detector
(User Experiment)

OPC-UA
Endpoint

ELBE (Accelerator)
Facility Control System

ADC | FPGA | SFP

RAW Data

Facility/Environment Data

## Data Center

TSDB

SFP | FPGA

Meta

Virtual Server

Metadata

Lab Book

Dr.-Ing. Oliver Knodel | Department of Information Services and Computing | Computational Science Group | www.hzdr.de

# HZDR Data Management Strategy — Data Publication

— Data/Meta-Data standards are the key of a usable data publication (e.g. NeXus, CERIF, …).

— **Metadata:** the who, what, when, where, why, how of your research.

— All data generated during the experiment:
  • RAW and
  • Facility Data,
  • Results/Analysis,
  • Workflows and
  • Metadata from the Lab Book.

Dr.-Ing. Oliver Knodel | Department of Information Services and Computing | Computational Science Group | www.hzdr.de