**GeRDI** Generic Research Data Infrastructure

# Research Data Infrastructures – How Generic, How Specific? Overview of the GeRDI Project

Hans-Joachim Bungartz

Department of Informatics, Technical University of Munich

Leibniz Supercomputing Center (LRZ)

German Research & Educational Network (DFN)

# Personal Background

## Research

- Education in Mathematics & Informatics (sorry – no physics, no materials … ☹)
- Scientific Computing, High-Performance Computing
- Computational Science and Engineering
- Data analytics via "numerical machine learning"

## IT infrastructure

- Heading DFG's Commission for IT Infrastructure 2006-2013
- Heading DFN since 2011 (network, services – EduRoam, DFN-AAI, …)
- Task Force of the Wissenschaftsrat that led to NHR / National Supercomputing
- Various advisory boards

# Attention! Brief philosophical detour (concerning the 4 paradigms) ☺

# 3 Paradigms ... (my list counts only 3 ... ☺)

**#1: Experiment**          ... somehow data-based

(cf. statistics: conclude about reality, based on experiments)

**#2: Theory**          ... somehow model-based

(cf. stochastics: conclude about reality, based on models)

**#3: Computation**          ... bridging the gap

# A Brief History of "Computational"
## (Computational ←→ Computation-aided)

**1st generation of "computational": qualitative (forward) simulation**

**2nd generation of "computational": optimization, inverse problems**

**3rd generation of "computational": quantitative through analytics/ML**

**Ultimate goal of "computational": predictive science**

**The current AI/ML hype is also due to the fact that we now have the computational power (HW, algorithms) that was absent in the early days of AI/NN/…**

# End of detour – back on Earth again… ☺

# Topics

1. Motivation

2. Vision & Mission

3. GeRDI at a Glance

4. Important Aspects

5. Outlook

Funded by
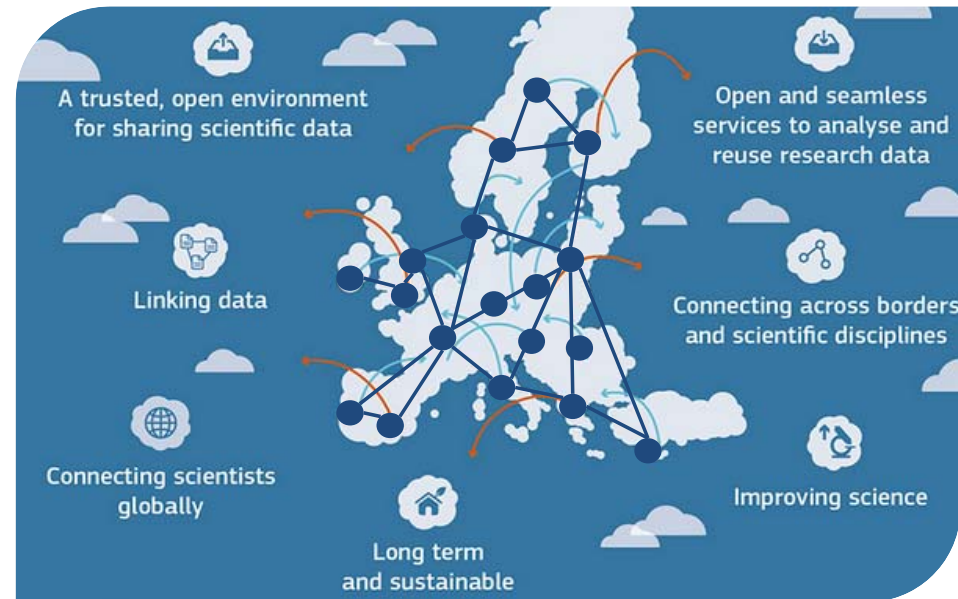
# 1. Motivation

GeRDI — Generic Research Data Infrastructure

What was the motivation to initiate GeRDI?

# European Developments: European Open Science Cloud

## Idea: European Research Area via a Network of Research Data Centers

# Nat'l Developments: Recommendations of the German Council for Research Information Infrastructures (RfII)

**Rat für Informations Infrastrukturen**

**LEISTUNG AUS VIELFALT**
Empfehlungen zu Strukturen, Prozessen und Finanzierung des Forschungsdatenmanagements in Deutschland

| Nr. im Text | Empfehlungen mit höchster Priorität |
|---|---|
| 4.1.1, 4.1.2 | Phasenmodell für die Entwicklung von Informationsinfrastrukturen – Planbarkeit und Mindeststandards sicherstellen – geordnete Übergänge in geeignete Trägerschaft über unabhängige Begutachtungen organisieren |
| 4.2.1 | Aufbau einer Nationalen Forschungsdateninfrastruktur (NFDI) – Kompetenzen bündeln und Grundversorgung mit Services für das Forschungsdatenmanagement schaffen |
| 4.2.3 | Arbeitsteilige Organisation von Services in Verbundstrukturen – übergreifende Infrastruktur- und Kompetenzzentren etablieren |

# Research Data Management Initiatives in Germany (excerpt)

## Hessen, Hamburg, Baden-Württemberg, Helmholtz, Leibniz, Fraunhofer, …

Industrial Data Space
FhG

Hamburg Open Science

LeibnizData

HeFDI – Hessische
Forschungsdateninfrastrukturen

Helmholtz Data
Federation

eScience
Baden-Württemberg

# On the one hand …

- A lot of political initiatives (open, data, FAIR, …)

- Increasing awareness of importance of research infrastructures

- Even more: RI as an essential part of science – not just "something that has to be there"

- Emerging funding frameworks – DFG-LIS/KfR, EOSC, NFDI, …

- Several institutional responses
  - Research Associations
  - States
  - Projects (CERN-LHC, Copernicus, SkA, …)
  - "well organized" communities

- Rather mature level for "Huge-Data communities" (Helmholtz, e.g.)

# On the other hand …

- Various isolated (community-driven) attempts

- In Germany often via DFG funding

- Creates increasing concern at DFG:

  – in their bodies (AWBI/LIS, KfR/WGI, …)

  – among reviewers: "*How many of such research data-related database systems are we expected to fund*?"

  – Lack of sustainability – standard issue in scientific software

  – Lack of interoperability (the "I" in FAIR …)

  – Too many similar developments, missing economy of scale

- Not very well developed situation for the "long tail" (universities, "Small-Data communities") ➔ primary target group for GeRDI
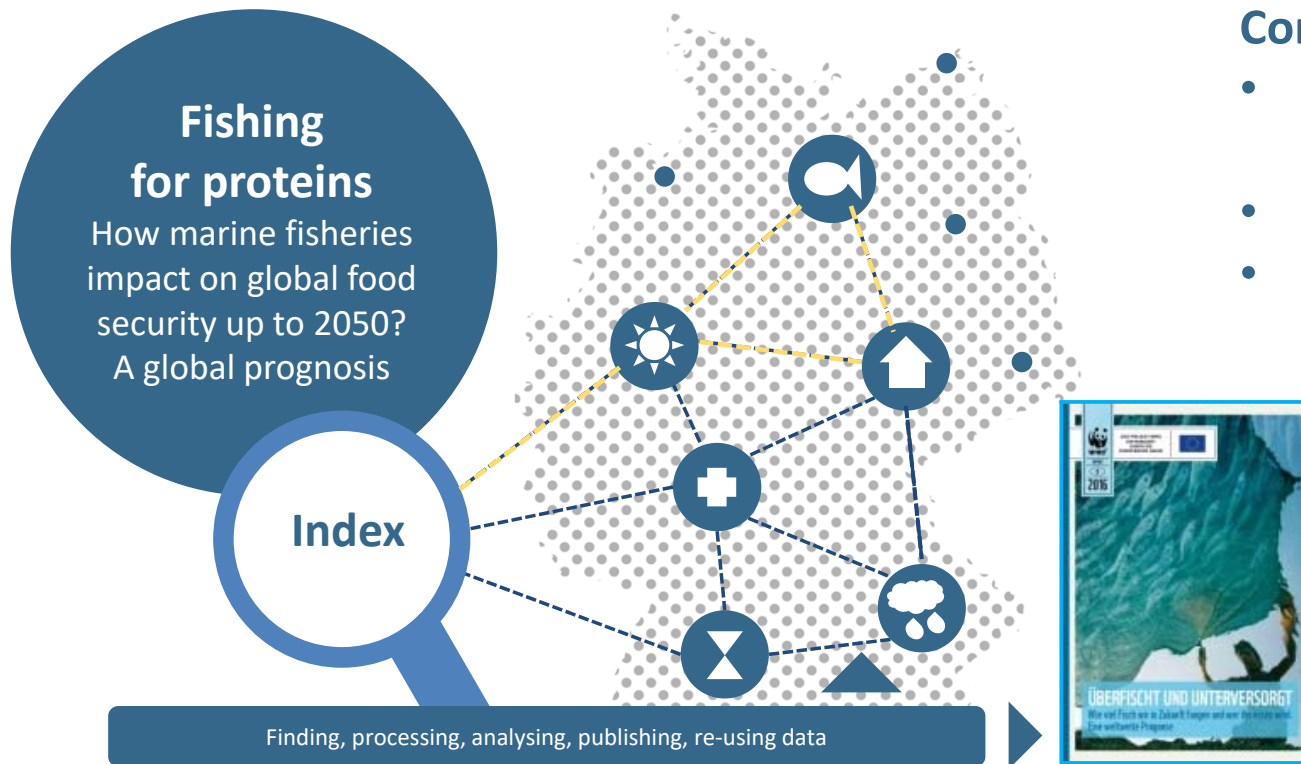
# Generic vs. Specific

- The specialist's / domain scientist's view:
  - Starting point: "My science is special, my data are special, my requirements are special … hence my Research Data Infrastructure has to be special, too."
  - Consequence: "Let's get some help from experts and develop something for us."

- The generalist's / IT person's view:
  - Starting point: "Data are data, we can work with abstractions of concrete things; data need data base / data management systems. That's it."
  - Consequence: "Let's build a generic RDI useful for all fields."

- Truth: somewhere in-between
  - At the top: domain-dependent data ➜ specific top
  - At the bottom: bits & bytes ➜ generic bottom
  - Explore the frontier!

**Generic Research Data Infrastructure**

# 2. Vision & Mission

# GeRDI – The basic idea

## Fishing for proteins
How marine fisheries impact on global food security up to 2050? A global prognosis

**Index**

Finding, processing, analysing, publishing, re-using data

ÜBERFISCHT UND UNTERVERSORGT

## Connect repositories across disciplines

- Collaboration with various research communities
- Focus on the data life cycle
- Contribution to the European Open Science Cloud and national initiatives

CAU Kiel University Christian-Albrechts-Universität zu Kiel

DFN DEUTSCHES FORSCHUNGSNETZ

lrz Leibniz-Rechenzentrum der Bayerischen Akademie der Wissenschaften

ZBW Leibniz-Informationszentrum Wirtschaft Leibniz Information Centre for Economics

ZIH

Funded by DFG

# Vision

Significant increase of interoperability.

Principle "as generic as possible, as specific as necessary" allows for effective and efficient usage of research data cross disciplines.

# Mission

GeRDI will provide **generic**, **sustainable** and **open** software connecting research data repositories to enable multidisciplinary and FAIR research data management.

This **software** will be based on common standards and be developed in close collaboration with several research communities to ensure a best match to the requirements of different disciplines.

GeRDI will promote a wide usage of its software and thus contribute to establishing an active GeRDI community, which will continue to flourish beyond the life span of the project.

All project results – in particular software, training support and business model – will form a German contribution to the European Open Science Cloud.

# 3. Project at a Glance

## Generic Research Data Infrastructure

How GeRDI is organised

# Facts and Figures

- Initiated by funding agencies (BMBF, later DFG) – "pre-NFDI" initiative

- Funded by DFG within its programme „Wissenschaftliche Literaturversorgungs- und Informationssysteme" (DFG-LIS)

- Outline:
  – Phase 1 (2016-2019): concepts, pilot with selected communities
  – Phase 2 (2019-2022): production, roll-out
  – Budget: ca. 3 Mio Euro (Phase 1)

- 5 partners + associated community partners

- Phase 2 delayed – revision of proposal required

# Project Consortium

**Plus Advisory Board gathering other key players in Germany**

## Infrastructures



## Communities

Funded by

# GeRDI Perspective & Embedding

|  | Phase I 2016-2019 | Phase II 2019-2022 | from 2022 |
|---|---|---|---|
| **GeRDI Project** Concept and Evaluation of a Research Date Infrastructure | Services, Pilot, Operation Model, Training | Transfer of Phase I: Put GeRDI into Operation, Roll-Out, New Communities |  |
| **BMBF/DFG** Setup, Operation, Organization of a National Research Data Infrastructure | Expert discussions about Research Data Management | Tendering of a National Research Data Infrastructure, Roll-Out, Sustainable Organization | Self-sustained and self-organized<br>• Operation<br>• Maintenance<br>• Development |

Image: www.digitalbevaring.dk

CAU Kiel University Christian-Albrechts-Universität zu Kiel · DFN DEUTSCHES FORSCHUNGSNETZ · lrz Leibniz-Rechenzentrum der Bayerischen Akademie der Wissenschaften · ZBW Leibniz-Informationszentrum Wirtschaft Leibniz Information Centre for Economics · ZIH

Funded by

DFG

CC BY NC SA

# Work Packages

# Aspects of Sustainability (cf. NFDI discussions)

- Software
  - Maintenance?
  - Training?
  - Modifications, extensions, variants, …? Open source?

- Hardware
  - Centralized at "hubs" (cf. HPC in Germany – PRACE, Gauß Center, Gauß Alliance, NHR)?
  - Semi-centralized at community hubs?
  - Decentralized with users / institutions / universities?

- Governance
  - Decision-making?
  - Integrated in NFDI governance vs. independent?

# 4a. Communities & Scenarios

Who are the (initial) users of GeRDI's services?

# Our Research Communities

**Various disciplines are involved**

- Alpine Environmental Data Analysis Centre

- Digital Humanities

- Environmental Resource and Ecological Economics

- Hydrology and River Basin Management

- Molecular Cell Biology and Genetics

- Paleoceanography

- Socio economics

- Tumor Diseases

Long tail – no existing solutions, but openness

Existing links to partners

**Why these?**
Long tail – no existing solutions
Needs & openness
Existing links to partners
Certain breadth of topics

Community Partner

Funded by

# Research Case 1: Fishery Management

Fishing for proteins: How marine fisheries impact on global food security up to 2050?

**Research**

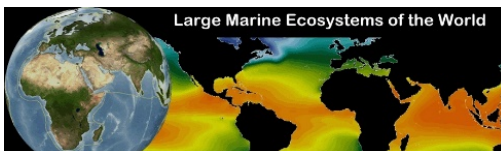Predict the amount of fish available for worldwide nutrition



Fig. Z3
Amount of per capita fish consumption, fish catches and population size on an LME basis in 2010. Data: Sea Around Us database/own maps

**Activities**
- Collect data from heterogenous sources (FAOStat, SeaAroundUs, …)
- Preprocess, analyze, and predict
- Publish
  - WWF-Report
  - Research paper, doi: 10.1111/gcb.13060
  - Research data, doi: 10.1594/PANGAEA.856741

# Research Case 1: Fishery Management

Fishing for proteins: How marine fisheries impact on global food security up to 2050?



| | |
|---|---|
| SEA AROUND US — FISHERIES, ECOSYSTEMS & BIODIVERSITY | large marine ecosystem based catch data |
| International Institute for Applied System Analysis (IIASA) | socio-economic data of 5 future scenarios |
| FAO STAT | prices of substitution goods |
| FishStatJ | trade statistics |
| Large Marine Ecosystems of the World | large marine ecosystem shape files |

# Research Case 2: Hydrology

Hydrology: How can we simulate & predict flash flood events?

**Research**

Modeling of flash flood events and assessment of impact factors



Photograph by Timothy Swinson, CC-BY 2.0 license

**Tasks**

- Aggregate data from various sources
- Preprocess and format data
- Provisioning of HPC systems
- Run hydrological and hydrodynamic simulations
- Analysis of the simulations
- Learning from federated data (simulations and measurements)
- Publish results on web portals

# Research Case 2: Hydrology

## Hydrology: How can we simulate flash flood events?

Vermessungsämter /
Survey bureaus


Landesamt für Digitalisierung,
Breitband und Vermessung

Maps and shape files

Hydrological and
hydrodynamical models/codes

Wasserwirtschaftsämter /
Water management

Measurement sensors

Water outflow rates

 Hydro_AS-2D

 Storm Water
Management Model

Wetterstationen /
Weather stations


DWD

Meterological data

 LARSIM

# Research Case 3: Microscopy

## Cell Imaging: How cells form tissue – Understand biological mechanisms

**Research**

- Development of novel microscopes
- Advanced imaging for biological samples (cells, tissues)

**Tasks**

- Navigate microscopy data using preview images
- Filter for image properties (magnification ratio, time resolution, size)
- Select images manually based on visual properties
- Preview of selected images as animated image sequence
- Interface to analysis service (e.g. cell tracking)
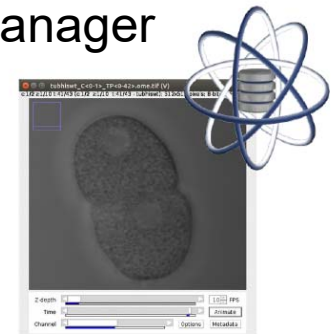
# Research Case 3: Microscopy

## Cell Imaging: How cells form tissue – Understand biological mechanisms

**Digital light sheet microscopy (DLSM)**

- Produces image data or image sequences
- 3D images of molecules
- 10GB to 100GB per image sequence
- Metadata encapsulated in images (OME-TIFF)
- Ingest of microscopy data into KIT Data Manager
- Setup KIT Data Manager as community-specific repository solution



KIT Data Manager

# FAIR principles

# GeRDI support

**F**indable

Harvesting data sources, Search Index

**A**ccessible

Data download support

**I**nteroperable
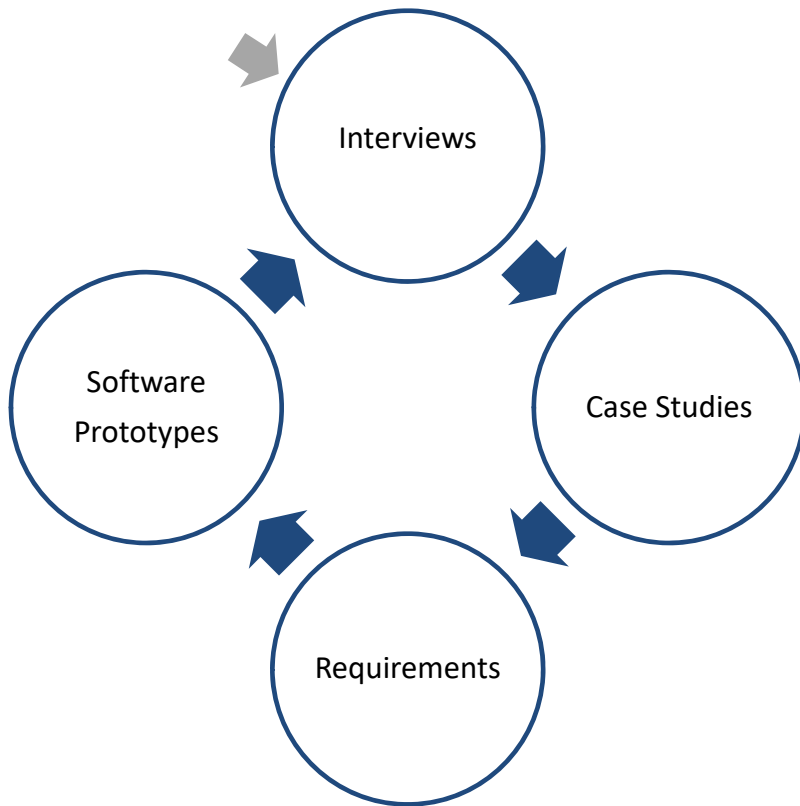
Mapping of Metadata to GeRDI-Metadata Schema

**R**eusable

Publishing results to a GeRDI-indexed repository

# GeRDI Generic Research Data Infrastructure

# 4b. Requirements Engineering

How do we define the requirements on GeRDI?

CAU Christian-Albrechts-Universität zu Kiel · DFN DEUTSCHES FORSCHUNGSNETZ · lrz Leibniz-Rechenzentrum der Bayerischen Akademie der Wissenschaften · ZBW Leibniz-Informationszentrum Wirtschaft Leibniz Information Centre for Economics · ZIH entrum für Informationsdienste und Hochleistungsrechnen

# Use (Research) Case Driven Approach

# Requirements Analysis

Interviews

Case Studies

Requirements

Software Prototypes

**Stakeholders**
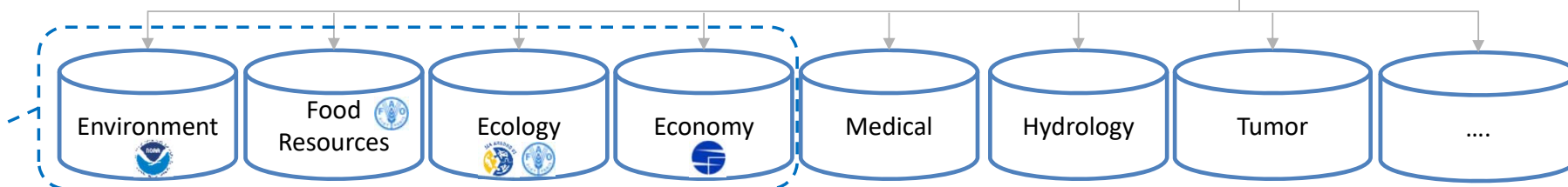
- Researchers
- Data Providers (repositories and archives)
- Developers (GeRDI project team)

**Continuously …**

… present state of the art

… derive new usage scenarios

… determine scenarios

… add to Backlog

… release new software versions

# Main Usage Scenarios

- Integrate and harvest existing repositories
- Integrated access to data of multiple disciplines
- Enable **new interdisciplinary research**
- Deploy new data repositories
- Data processing, analysis, and publication
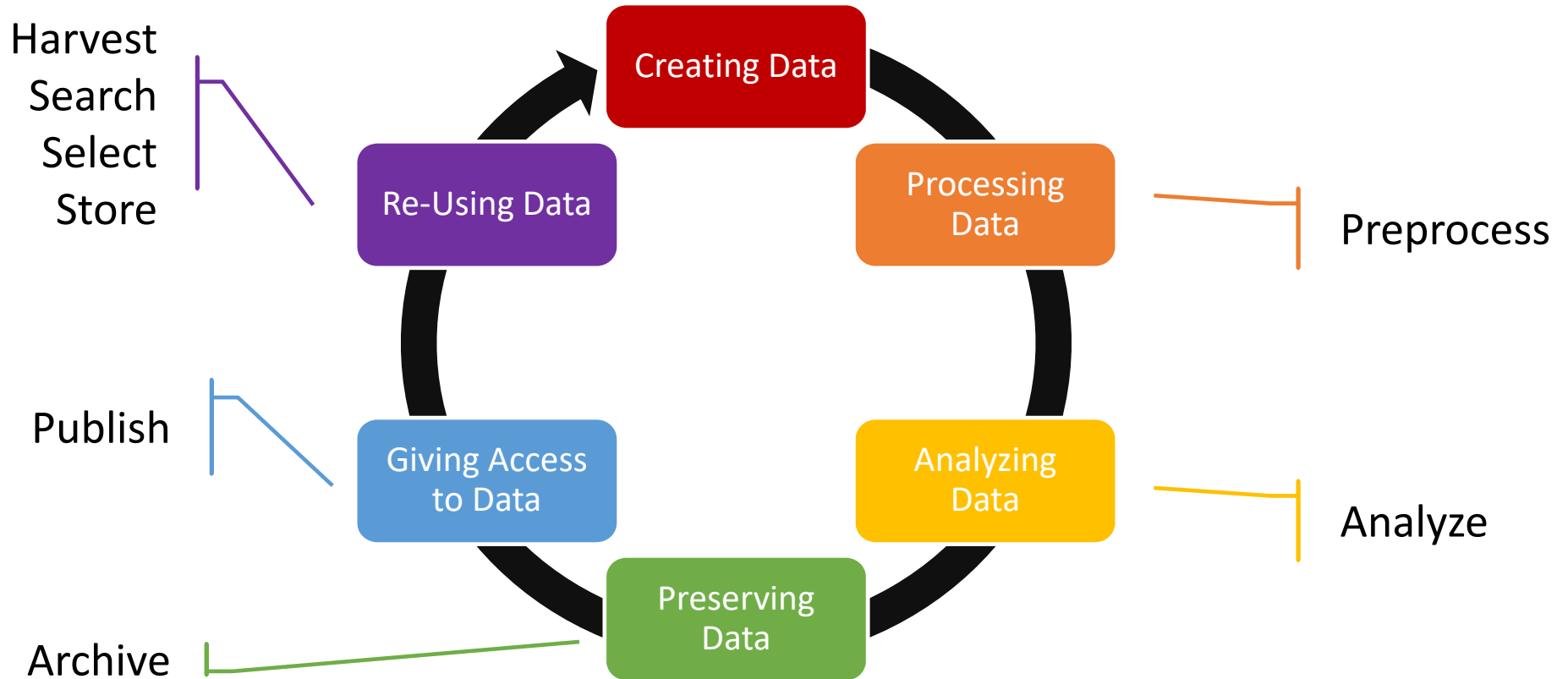
GeRDI

**Fishery Management**

Environment | Food Resources | Ecology | Economy | Medical | Hydrology | Tumor | ….

# 4c. Architecture Design

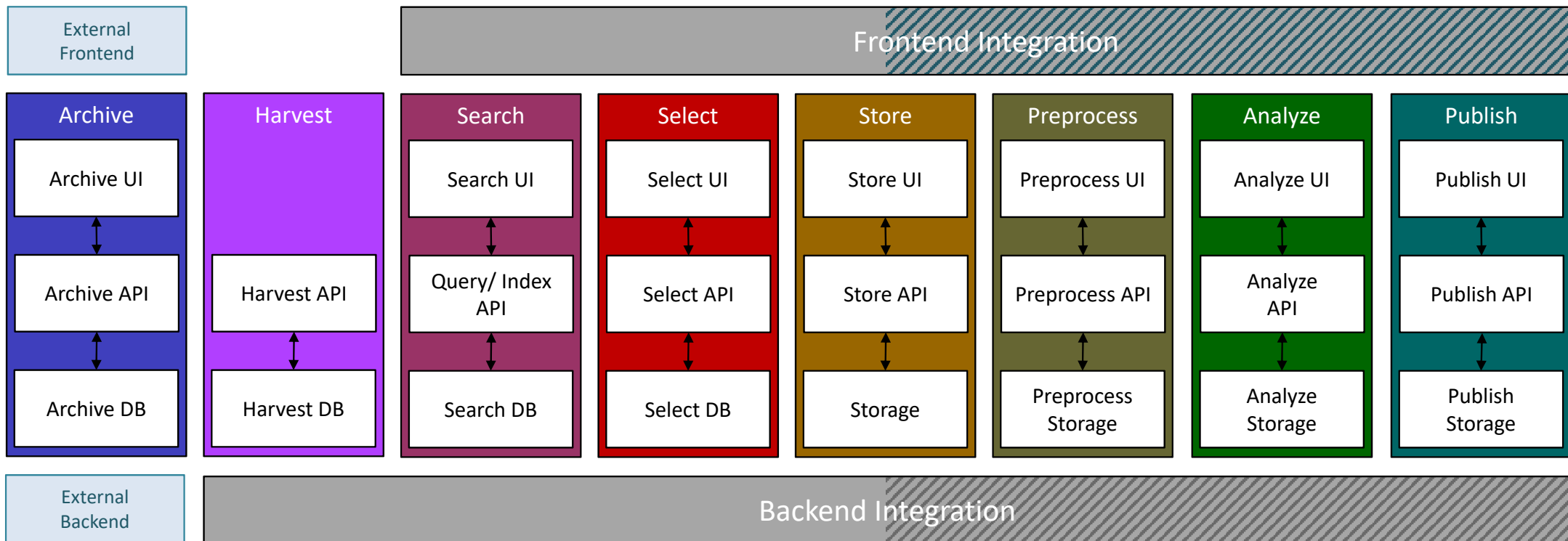Which services will GeRDI provide, and how are they interconnected?
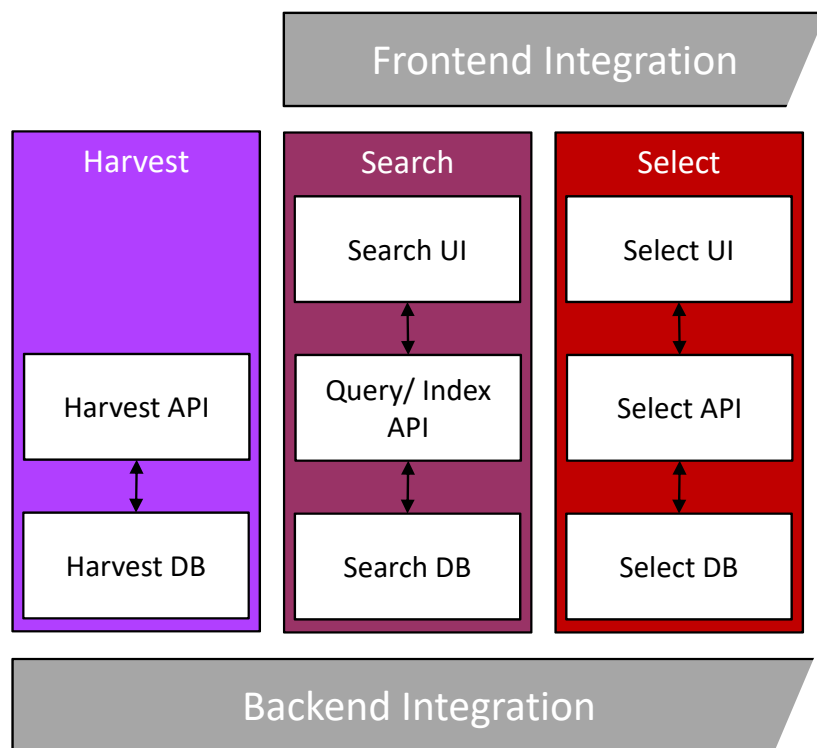
# GeRDI Services: Follow the Research Data Life Cycle

# GeRDI Services: Microservice Architecture Overview

| External Frontend | Frontend Integration | | | | | | |
|---|---|---|---|---|---|---|---|
| **Archive** | **Harvest** | **Search** | **Select** | **Store** | **Preprocess** | **Analyze** | **Publish** |
| Archive UI | | Search UI | Select UI | Store UI | Preprocess UI | Analyze UI | Publish UI |
| Archive API | Harvest API | Query/ Index API | Select API | Store API | Preprocess API | Analyze API | Publish API |
| Archive DB | Harvest DB | Search DB | Select DB | Storage | Preprocess Storage | Analyze Storage | Publish Storage |
| External Backend | Backend Integration | | | | | | |

# GeRDI Core Services



**Frontend Integration**

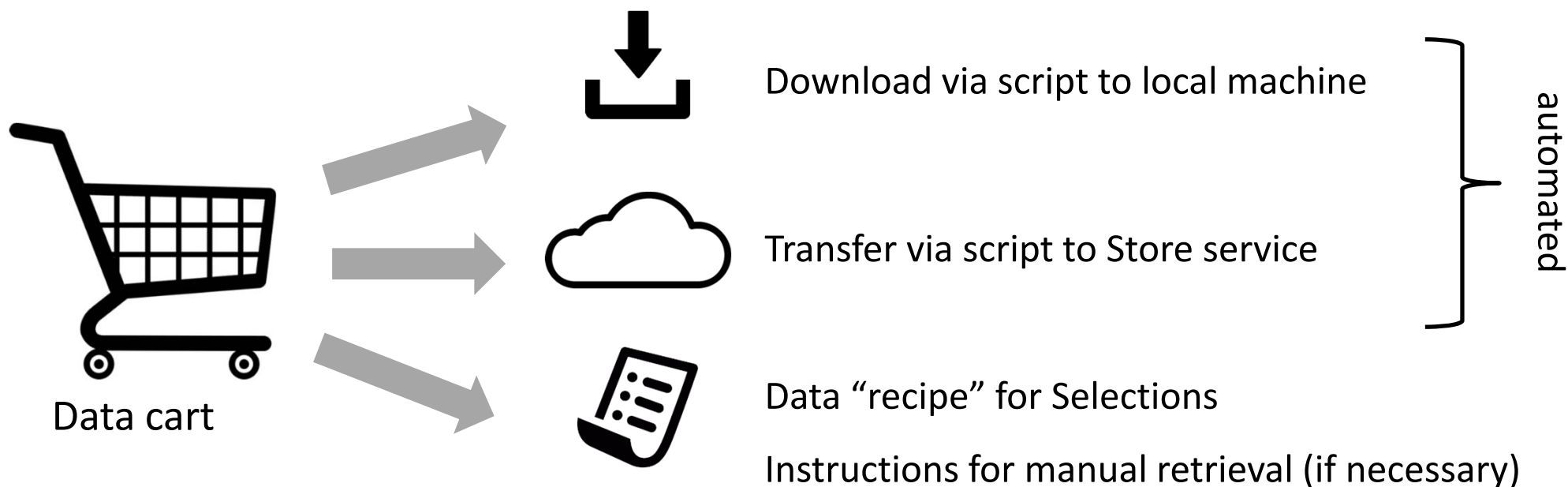| Harvest | Search | Select |
|---------|--------|--------|
| | Search UI | Select UI |
| Harvest API | Query/ Index API | Select API |
| Harvest DB | Search DB | Select DB |

**Backend Integration**

Components that are <u>developed and operated</u> by the GeRDI project team:

- Integration of data archives and repositories, based on standards such as OAI-PMH

- Meta data collection, integration, normalization, and enrichment, based on standards such as DataCite

- Search over multiple sources

- Selection of (meta) data into a data cart

# The Select Service

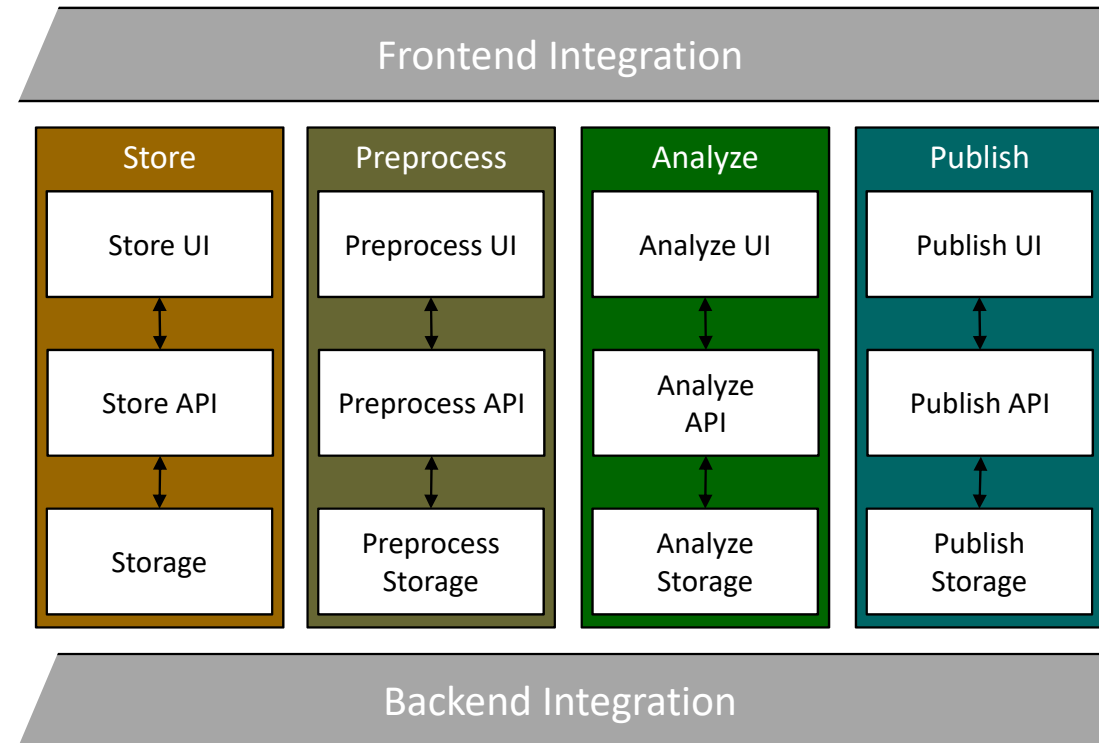General idea: select and save a named compound collection of data sets

Download via script to local machine

Transfer via script to Store service

automated

Data cart

Data "recipe" for Selections

Instructions for manual retrieval (if necessary)

DFN DEUTSCHES FORSCHUNGSNETZ

lrz Leibniz-Rechenzentrum der Bayerischen Akademie der Wissenschaften

ZBW Leibniz-Informationszentrum Wirtschaft Leibniz Information Centre for Economics

ZIH

Funded by DFG

# Extended Services

Components for which the GeRDI team
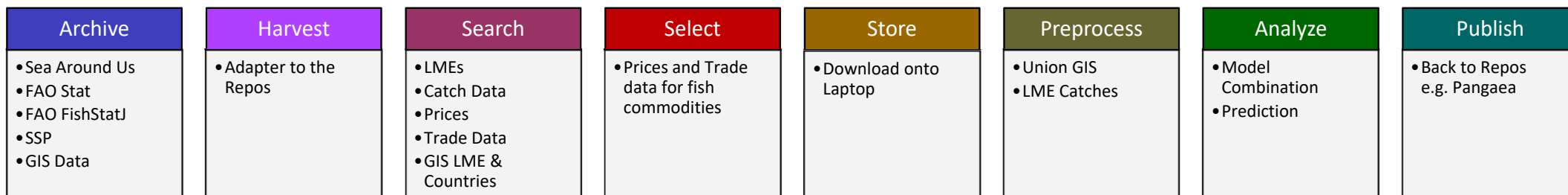
provides <u>reference implementations</u>:

- Storage on different systems
  (e.g. DFN-Cloud or HPC Staging)
  based on data recipes

- Merging, truncating, normalizing data

- Analyzing data
  (Jupyter, R-Studio, Simulations, etc.)

- Publish the derived data product:
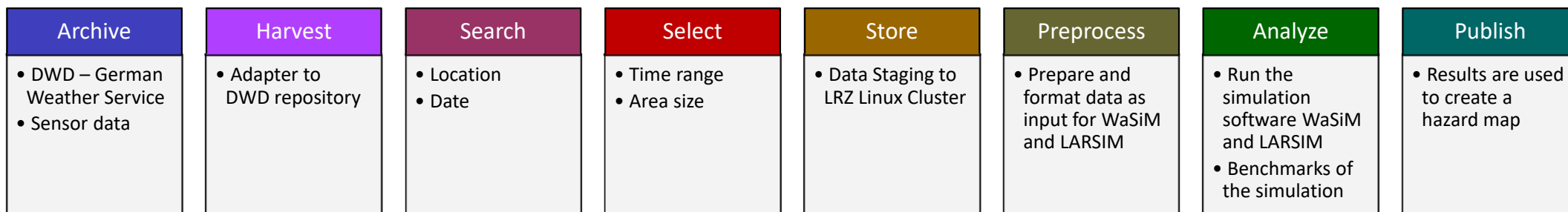  Reusability of research data

# Research Case Mappings

## Fishery Management (Environmental Resource and Ecological Economics)

| Archive | Harvest | Search | Select | Store | Preprocess | Analyze | Publish |
|---------|---------|--------|--------|-------|------------|---------|---------|
| • Sea Around Us<br>• FAO Stat<br>• FAO FishStatJ<br>• SSP<br>• GIS Data | • Adapter to the Repos | • LMEs<br>• Catch Data<br>• Prices<br>• Trade Data<br>• GIS LME & Countries | • Prices and Trade data for fish commodities | • Download onto Laptop | • Union GIS<br>• LME Catches | • Model Combination<br>• Prediction | • Back to Repos e.g. Pangaea |

## Hydrology – Flash flood modeling

| Archive | Harvest | Search | Select | Store | Preprocess | Analyze | Publish |
|---------|---------|--------|--------|-------|------------|---------|---------|
| • DWD – German Weather Service<br>• Sensor data | • Adapter to DWD repository | • Location<br>• Date | • Time range<br>• Area size | • Data Staging to LRZ Linux Cluster | • Prepare and format data as input for WaSiM and LARSIM | • Run the simulation software WaSiM and LARSIM<br>• Benchmarks of the simulation | • Results are used to create a hazard map |

# Schema for Harvested Meta Data
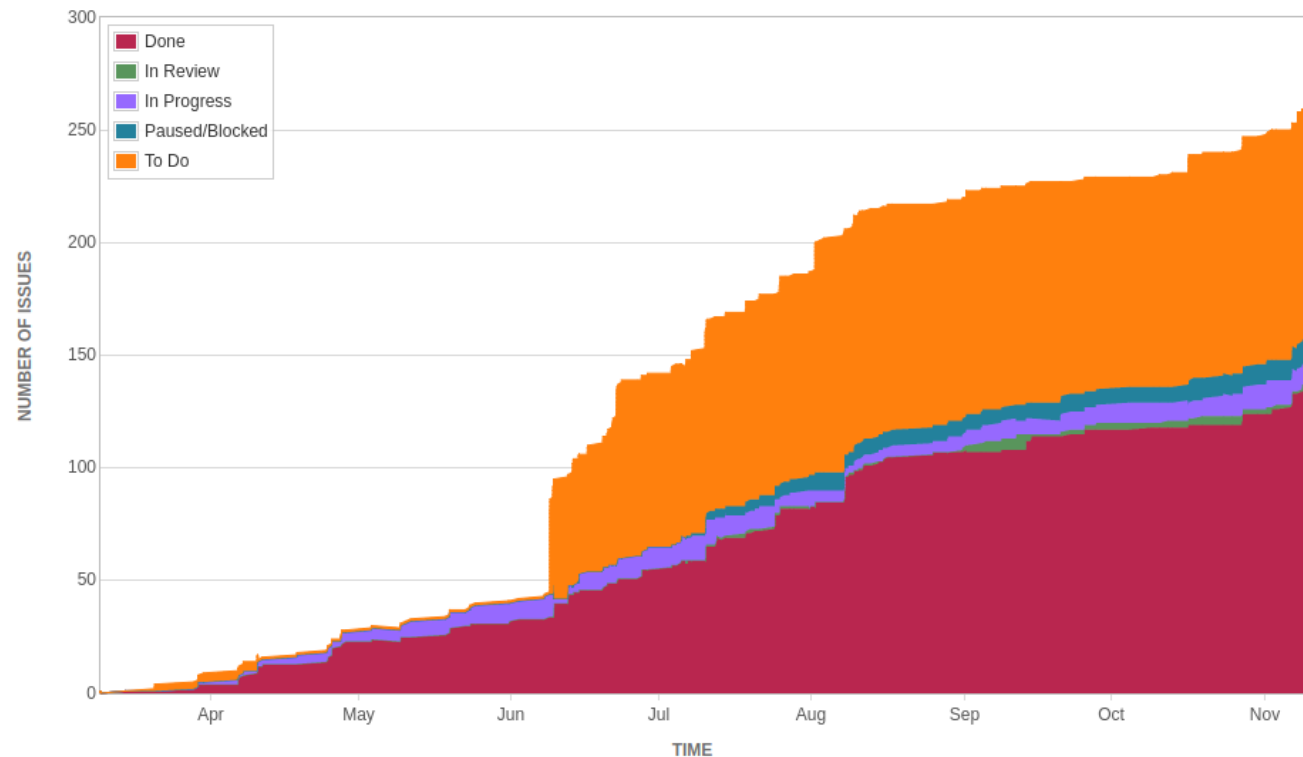
# Possible Implementation View

# Software Management

Agile Workflow Sprints every other week

- Estimation of task
- Assignment of tasks
- Review of the sprint via Atlassian tools

Continuous Integration with nightly builds

- based on Bamboo and Kubernetes

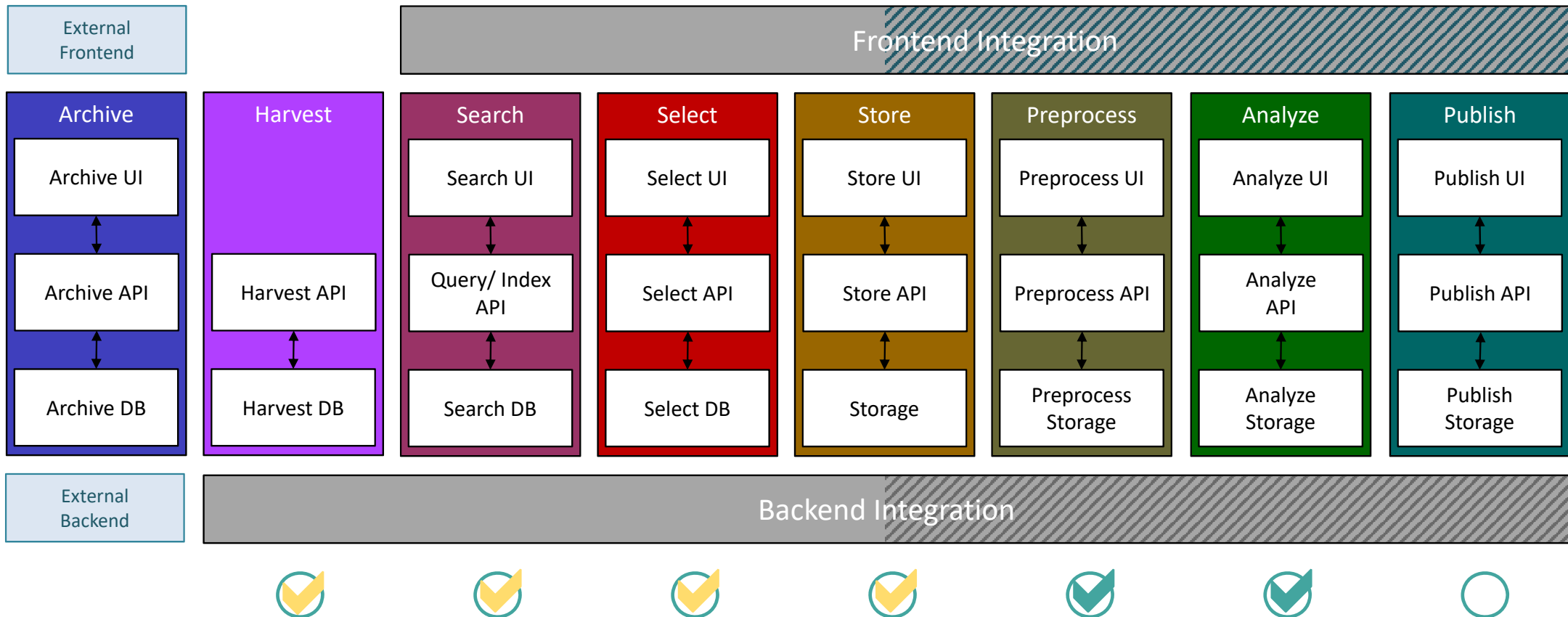GeRDI — Generic Research Data Infrastructure

# 5. Status

# Achievements of GeRDI

- Generic metadata schema based on Datacite

- Harvester technology, including metadata mapping capacity

- Microservice architecture

- Services for Search, Bookmark, and Store

- First version of AAI services

- Accepted as GO FAIR implementation network

# GeRDI Services: Microservice Architecture Overview



Generic Research Data Infrastructure · www.gerdi-project.eu · Hans-Joachim Bungartz

# Success criteria

| | Target End of 2019 | June 2019 |
|---|:---:|:---:|
| **Number of current or pending collaborative projects (incl. NFDI consortia) that would like to use GeRDI software** | 4 | 9 |
| **Number of interviews / feedback meetings with the community partners** | 90 | 79 |
| **Number of connected research data repositories (without Zenodo)** | 15 | 13 (16) |
| **Number of harvested metadata records (without Zenodo)** | 700.000 | 567.036 |
| **Number of presentations about GeRDI** | 45 | 39 |

Funded by

# Success criteria

Target end of 2019      June 2019

| Repository in test.gerdi.org (NOT demo.gerdi.org) | Number of datasets |
|---|---|
| *Zenodo* | *1.356.474* |
| 1 PANGAEA | 384.138 |
| 2 DWD | 102.519 |
| 3 European Nucleotide Archive (ENA) | 50.002 |
| 4 ArcGIS | 8.051 |
| 5 Sea Around Us | 3.692 |
| 6 IMR, | 1.252 |
| 7 SOEP, | 473 |
| 8 OCEANTEA | 164 |
| 9 FAOSTAT | 78 |
| 10 FishStatJ | 27 |
| 11 LMU-ifo Economics & Business Data Center (EBDC) | 119 |
| 12 AlpenDAC | 336 |
| 13 Eurostat (only partly harvested for not flooding GeRDI index) | 9.257 |
| 14 U.S. National Library of Medicine | 5.908 |
| 15 OGLP | 947 |
| 16 Open Data LMU | 73 |
| **SUMME** | **567.036** |

**Number of current or pending collaborative projects (incl. NFDI consortia) that would like to use GeRDI software**

**Number of interviews / feedback meetings with the community partners**

**Number of connected research data repositories (without Zenodo)**

**Number of harvested metadata records (without Zenodo)**

**Number of presentations about GeRDI**

**Thanks for your attention!**

**Cocktail & dinner time …**