# Advancing Data Exchange Innovations in FAIR Data Spaces

Christoph Lange, Fraunhofer FIT, and project team

27.05.2024

**FAIR Data Spaces**

# FAIR Data Spaces at a glance

- Vision: Development of a common cloud-based data space for industry and science

- Mission: Create and expand synergies between existing technologies and communities

- Financing: Funded by the Federal Ministry of Education and Research

- Runtime: May 2021 – December 2024

- Participation: 16 participating organizations

**FAIR Data Spaces**

# Project goals

FAIR* Data Spaces aims to build a shared cloud-based data space for economy and science by linking Gaia-X and NFDI

- Identify and leverage synergies of cooperation between the two initiatives

- Interweaving the content of the initiatives by clarifying legal and ethical issues and providing technical foundations

- Promoting a sovereign exchange of data between industry and science both nationally and in the EU in concrete applications and fields of work

*Guidelines findable, accessible, interoperable, reusable

**FAIR Data Spaces**

# Target communities

## Gaia-X

Economy

- EU-initiated project, brings together stakeholders from industry, science and administration
- The goal is to create an open, transparent and secure federated data infrastructure
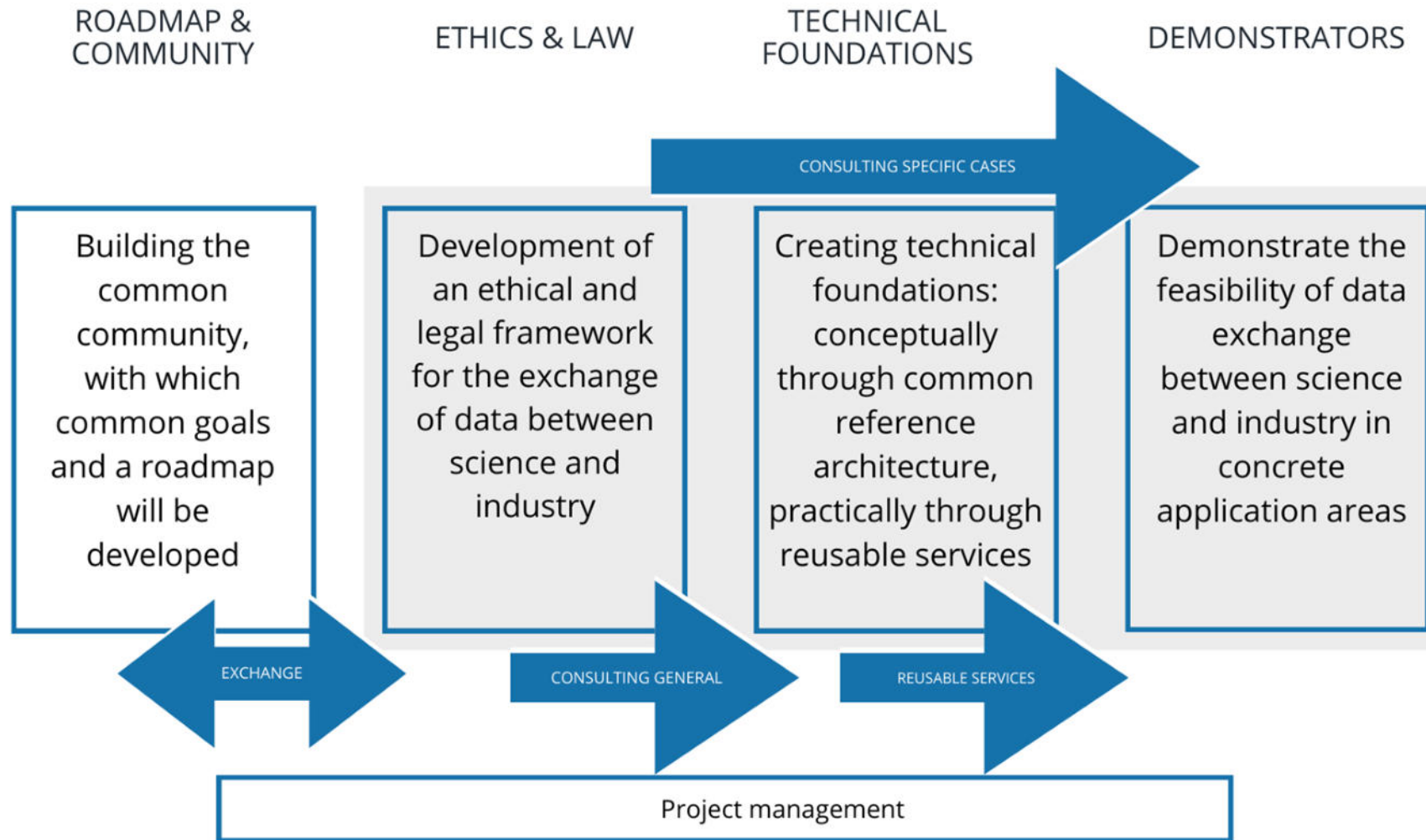- Gaia-X architecture consists of a variety of individual platforms that follow a common standard

## NFDI

Science

- purpose of the association is to promote science and research through a National Research Data Infrastructure, which establishes and further develops an overarching research data management in Germany and increases the efficiency of the entire German science and research system
- Making research data FAIR

In addition, linkage with EU data spaces and connection to EOSC

FAIR Data Spaces

# Overview task areas



ROADMAP & COMMUNITY — Building the common community, with which common goals and a roadmap will be developed

ETHICS & LAW — Development of an ethical and legal framework for the exchange of data between science and industry

TECHNICAL FOUNDATIONS — Creating technical foundations: conceptually through common reference architecture, practically through reusable services

DEMONSTRATORS — Demonstrate the feasibility of data exchange between science and industry in concrete application areas

CONSULTING SPECIFIC CASES

EXCHANGE

CONSULTING GENERAL

REUSABLE SERVICES

Project management

**FAIR Data Spaces**

5

# Objectives of the Open Call

- Have further demonstrators developed, beyond our three own ones (see below)
- Either extensions/additions to our demonstrators, or independent innovation
- Based on the same principles, standards and technology as FAIR Data Spaces
- Round 1 (call closed; June to November 2024): up to 3 contracts @ 60 k€
- Round 2: starting as soon as possible, ending in December 2024
  - same overall budget (180 k€)
  - additional requirement: bridging operational initiatives

| existing industry initiative<br>e.g. ### Data Space | | existing research data<br>infrastructure<br>e.g. NFDI 4 ### |
|---|---|---|

**FAIR Data Spaces**

# Open Call Round 1

Proposal evaluation criteria:

| | |
|---|---|
| Interoperability with FAIR Data Spaces | 15% |
| System architecture | 15% |
| Development process | 5% |
| Documentation | 3% |
| Testing | 15% |
| Security | 7% |
| Innovation degree (evaluated by expert review board):<br>• Practical demonstration<br>• Industry↔research cross-benefit<br>• Going beyond state of the art (of FAIR Data Spaces and general)<br>• Awareness of ELSA | 40% |

FAIR Data Spaces

# Open Call Round 1

Requirements and Execution:

- Technology: common programming languages; container deployment; documented interfaces (e.g., OpenAPI); W3C semantic metadata (RDF etc.)
- Legal paperwork (not in the focus of this webinar), reference projects, team CVs
- 8-page concept paper:
  a. summary, architecture, frontend/backend implementation (addressing all evaluation criteria), innovation / added value.
  b. informative: time schedule, applicant profile, cost calculation
- Deliverable (at end of contract): open source code, documentation, website, video
- Schedule (~6 months):
  a. continuous implementation using repository and issue tracker
  b. regular status calls
  c. joint hackathon
  d. public presentation

**FAIR** Data Spaces

# Open Call Round 2

Discussion points for today – not formally implying any conditions that will apply to Round 2

- Who is here, representing what initiatives?
  a. ideally we'd have "consortia" formed from one company + and research organization (but not necessarily a member of an NFDI consortium or industry association)
     - However, we might also contract a single organization that provides strong evidence of support from "both sides" (e.g., letters of support)
     - It will not be sufficient to merely apply "industry-ready" *technology*, such as the EDC
  b. technical foundations on these slides are *representative*, so do build on what you have established in your existing initiatives, but there *must* be a "FAIR" interoperability layer.
     - as an alternative to Gaia-X & IDS specifically, you may refer to Data Spaces Blueprint
  c. Are other Fraunhofer institutes eligible – to be clarified! (Fraunhofer contracting "itself")
- Have you already reached out to "the other side"?
- What do you consider feasible within ~4 months (end of project minus formal/legal process)
- Any further questions you may have?
  a. First round was coordinated via e-Vergabe, second round most likely as well

    https://www.nfdi.de/fair-data-spaces-newsletter/

**FAIR Data Spaces**

# Technical Foundations / Architecture (1)

Architecture with a heavy focus on modern cloud-native technologies

- Virtual machines

  openstack. [1]

- Containers

  docker [2]    kubernetes [3]

- Storage

  ceph [4]    Object storage (S3) [5]

- CI/CD

  flux [6]

  GitLab Runners

  GitHub Actions

Powered by the cloud [7]    de.NBI    community
GERMAN NETWORK FOR BIOINFORMATICS INFRASTRUCTURE

FAIR Data Spaces

# Technical Foundations / Architecture (2)

## Gaia-X Architecture and Trust Framework (we are a "Gaia-X Qualified Project")



Source: Gaia-X European Association for Data and Cloud AISBL

# Technical Foundations / Architecture (3)

Sovereign Data Exchange via the Dataspace Protocol



Source: International Data Spaces Association

# FAIR Data Spaces Demonstrators

https://github.com/FAIR-DS4NFDI/FAIR-DSWiki/wiki



Vector Data

Raster Data

Visualize & Analyze

**FAIR-DS Demonstrator NFDI4Biodiversity and Gaia-X**

Data

Schemas

Git Repository

Data Quality Reports

Human Readable
HTML
Git

Machine Readable
JSON

**FAIR Research Data Quality Assurance and Workflows**

Station 1
Station 2
...
Station N

Push
Pull
Push
Pull
Push
Pull

Central Service

Scientist

Access Results
Submit Analysis

Station Registry

**Cross-Platform FAIR Data Analysis**

FAIR Data Spaces

# Demonstrator 4.1: NFDI4Biodiversity

Nikolaus Glombiewski, Bernhard Seeger (Philipps-Universität Marburg)

FAIR Data Spaces

# NFDI4Biodiversity Demonstrator Overview

- Application:
  - Spatio-Temporal Data Analysis
  - Heterogeneous Data Sources
  - Rust, Python, Angular
  - Docker, OpenIdConnect (Keycloak)

- Connection to Research Data Infrastructure:
  - Part of NFDI4Biodiversity Research Data Commons
  - The Visualization, Analysis and Transformation (VAT) System powered by Geo Engine

- Connection to Gaia-X:
  - Service Offering in a Federated Catalogue

**FAIR Data Spaces**

# Spatio-Temporal Data Analysis

| Species | Coordinates | Date |
|---------|-------------|------|
| Bird | 48.856614, 2.352221 | 10.05.1977 |
| Cat | 41.8933203, 12.4829321 | 18.04.1980 |
| Elephant | 52.517037, 13.38886 | 21.10.2015 |



Example 2:
Raster Data
(e.g. images from satellites)



Example 1:
Vector Data
(e.g. manual recordings)

**FAIR Data Spaces**

# Geo Engine: Adding Data Sources

# Geo Engine: Analysis

# NFDI4Biodiversity: Research Data Commons

# Connecting Layers with Geo Engine



- Application for Spatio-Temporal Data
- External Data Providers:
  - Standardized Protocols for Spatio-Temporal Data
  - Custom Data Exchange when necessary
- Also in "Mediation Layer" for offering FAIR Datasets

**FAIR Data Spaces**

# Geo Engine and Gaia-X Federated Catalogue

- Technological Basis:
  Gaia-X compliant catalogue developed by Eclipse XFSC (Cross Federation Services Components) using Spring, OpenIdConnect, PostgreSQL, Neo4j

- GeoEngine Self-Description:
  - A Service Offering in JSON-LD format
  - Verifiable Credentials: Set of claims or attributes, digitally signed by a trusted entity
    Who? What? Where? Which standard?

- Adding Data from supported dcat:DataService
  - OGC Environmental Data Retrieval (EDR)
  - In principal: Aruna Object Storage

# Demonstrator 4.2:
## Data Quality Assurance and Workflows

*Jonathan Hartman, RWTH Aachen University*

22

# Data Quality Assurance and Workflows

## Goals

- Build off of existing Infrastructure
  - git.rwth-aachen.de
  - Open Telekom Cloud (OTC)
- Provide an example of Automated Analysis / Data QA
  - the demonstrator
- Provide a framework for hosting / sharing Workflows

FAIR Data Spaces

# Data Quality Assurance and Workflows

## Infrastructure

## GitLab & Runners
- Lightweight agents controlled by CI/CD Scripts from Repositories
- Scalable, based on the workload
- Isolated, each runner context is run in its own container
- Customizable, capable of loading a huge variety of containers
- Multiple runners can be assigned to a project / group.

Computational Cloud

Runner 1
*(waiting...)*

Runner 2
*(waiting...)*

Runner 3
*(waiting...)*

# Data Quality Assurance and Workflows

1. A Repository is triggered by some event *(Commits, Merge requests, Scheduled, Hooks)*
2. An assigned runner picks up the job
3. The appropriate Container is loaded
4. Any scripting steps can be executed in the created environment

On Event

Computational Cloud

Runner 1
*(waiting…)*

Runner 2
RUNNING

Runner 3
*(waiting…)*

User Repository

- Use Container A
- Run my_script.py
- Save the Output

my_script.py

Runner 2
RUNNING

Container A

my_script.py

# Data Quality Assurance and Workflows

**Demonstrator Docker Container**

**Python Environment**
*(Python Executable, all required Libraries)*

**Demonstrator**
*(Code & Report Templates)*

## Demonstrator
- written in Python
- based on the Frictionless standard & library
- Available as a "pippable" library

## Provided as a Docker Container
- Hosted on git.rwth-aachen.de

FAIR Data Spaces

# Data Quality Assurance and Workflows

## Maintained by the User:

- Data to be Analyzed
  - Tabular

- A GitLab repository
  - CI/CD Script
  - Config file *(optional)*
  - Access Credentials to the data *(optional)*
  - Data Schemas *(optional, can also be stored with the data)*



**FAIR Data Spaces**

# Data Quality Assurance and Workflows

# Data Quality Assurance and Workflows

# Data Quality Assurance and Workflows

# Demonstrator 4.3:

## Cross-Platform FAIR Data Analysis
## PADME PHT

Yeliz Ucer Yediel, Muhammad Hamza Akhdar (Fraunhofer FIT),
Macedo Maia, Toralf Kirsten (University of Leipzig),
Mehrshad Jaberansary, Oya Beyan (University of Cologne)

**FAIR Data Spaces**

# Cross-Platform FAIR Data Analysis PADME PHT

- Idea: "Bring the algorithms to the data" by using Distributed Analytics (DA)
- Benefits:
  - The data remains in the control of the data providers
  - Research can leverage otherwise inaccessible data
  - The results are made more robust by incorporating a variety of datasets.
- Provides ecosystem from the first idea to the analysis results
  - Central Components: Playground, Train Creator, Train Store House, Train Requester,
  - Client Software: PHT Station

# PADME in a Nutshell

https://padme-analytics.de/ , https://docs.padme-analytics.de/

- Implementation of the PHT/FL concepts by using FAIR standards
- Result of a collaboration between four research institutes



- Based on containerization technologies ([www.docker.com](www.docker.com)), deployed on Kubernetes env.



- Benefits:
  - Operating system agnostic
  - Data source and data structure agnostic
  - Programming-language agnostic

# PHT PADME and EDC Integration
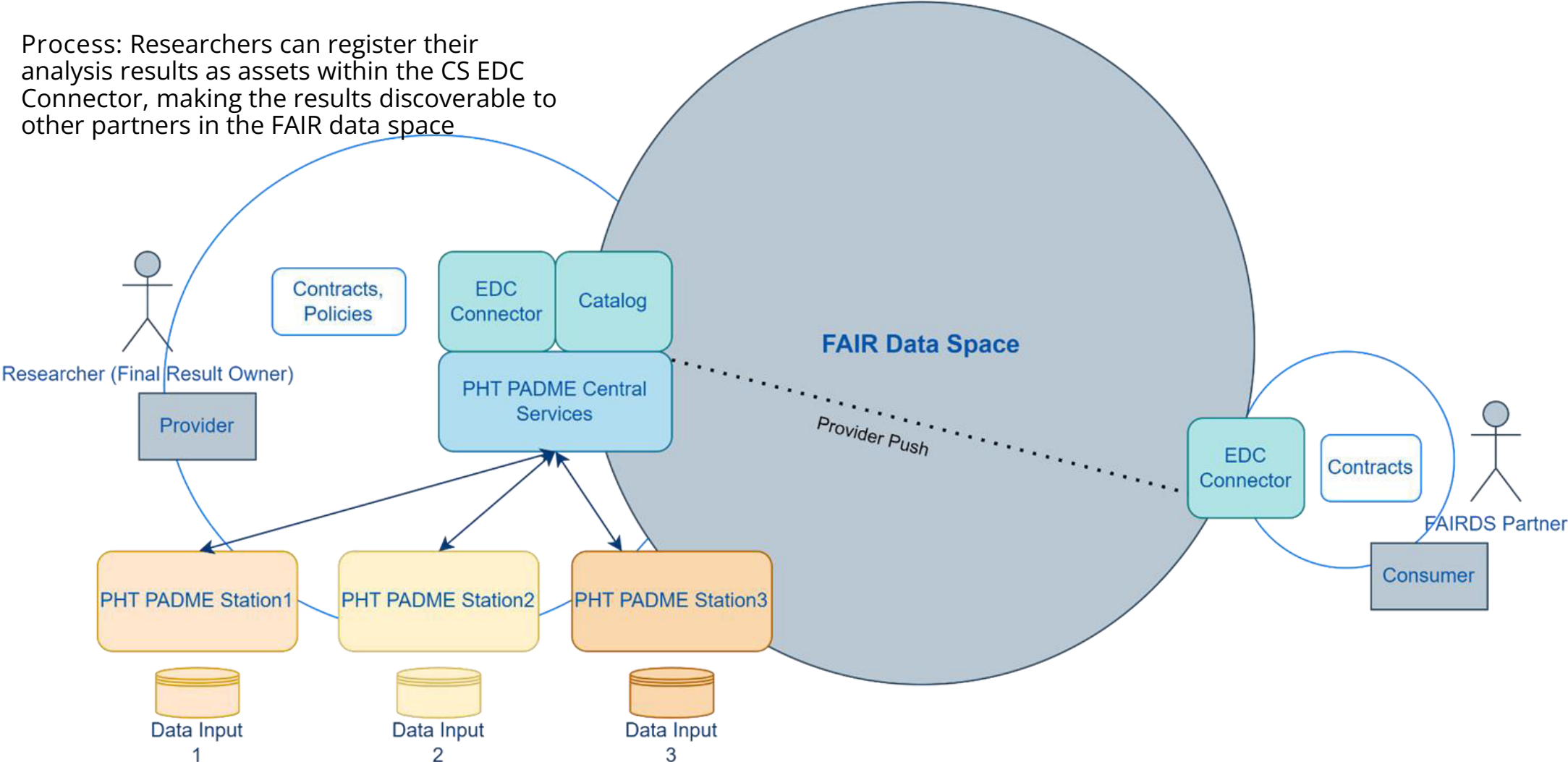
Use Case 1: PHT as a Data Provider

- Objective: Enable sharing of analysis results using EDC connector

- EDC Data Transfer Mode : Provider Push

- Process: Researchers can register their analysis results as assets within the CS EDC Connector, making the results discoverable to other partners in the FAIR data space

Use case 2: PHT as a Data Consumer

- Objective: Enable the PHT Station to consume data from other providers within the FAIR Data Space

- EDC Data Transfer Mode : Consumer Pull

- Process: The PHT Station can discover data catalogs of the other participants, negotiate contracts, and initiate data transfer requests. Upon a successful negotiation, the provided credentials are used to provide access to the Train, which can then execute the analysis

**FAIR Data Spaces**

# EDC Integration into CS - Provider Push Scenario

Process: Researchers can register their
analysis results as assets within the CS EDC
Connector, making the results discoverable to
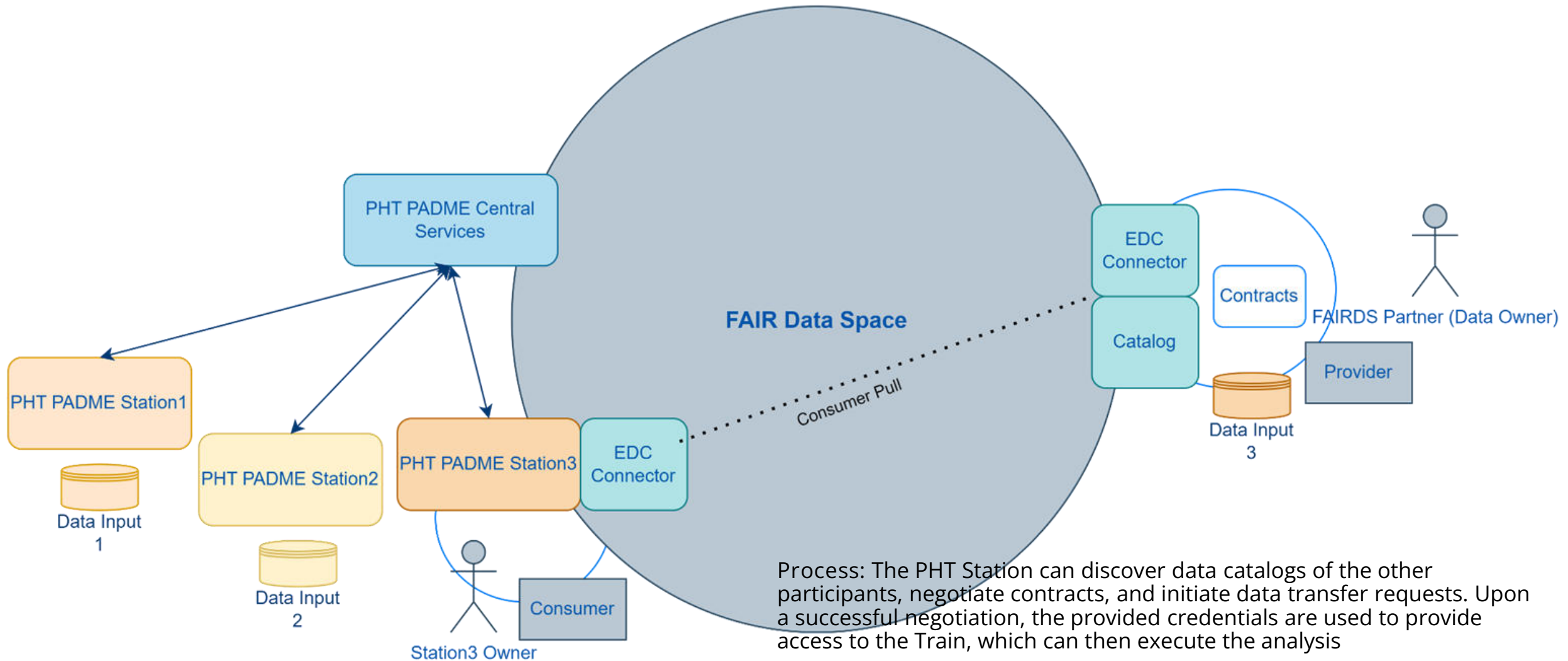other partners in the FAIR data space

# EDC Integration into Station Consumer Pull Scenario



Process: The PHT Station can discover data catalogs of the other participants, negotiate contracts, and initiate data transfer requests. Upon a successful negotiation, the provided credentials are used to provide access to the Train, which can then execute the analysis

# Federated Learning over Multiple PADME PHT Stations

- Federated learning involves training a central model using data distributed across multiple Stations (Client/Provider) and Central Service (Server/Consumer)

- Local models are trained on each PHT Station

- Each local model are sending to Central Service

- The application of a aggregation function over each local model weights determines a federated learning

- The federated model are sending back to each Station and retrained in the next round



**DIC UKL Leipzig**

**Bruegel Station Aachen**

**University Hospital Cologne**

Model Train Code — Local Segmentation Model

Model Train Code — Local Segmentation Model

Model Train Code — Local Segmentation Model

**Central Service (De.NBI Cloud)**

Request a New Train — Federated Segmentation Model

**Scientist**

37

# Use case: Liver Tumours Segmentation
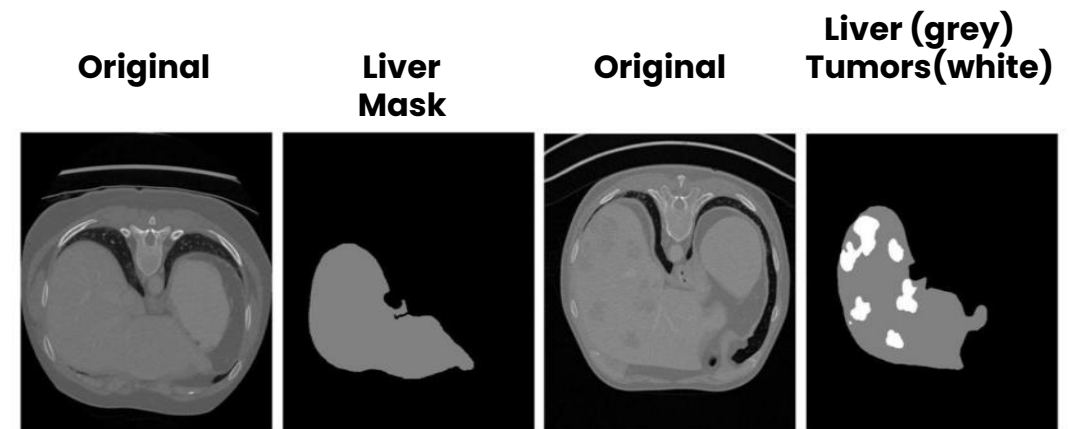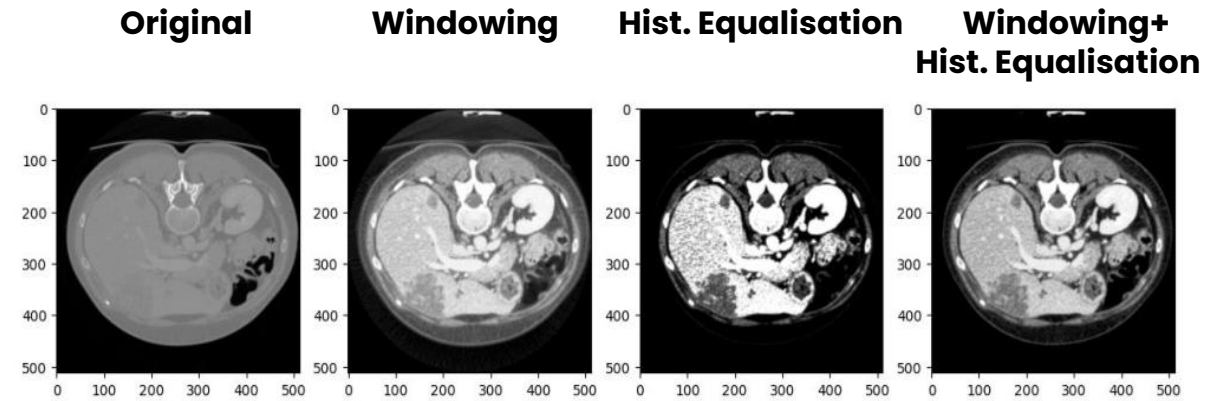
- Problem Statement
  - Based on The American Cancer Society's:
    - About 41,630 new liver cancer cases in US were diagnosed in 2023
    - About 29,840 people have died of these cancers
  - AI-based approaches helps to early detect tumours
  - However, the data can be distributed in different sources (e.g., hospitals)
  - Data access depends on distinct rules or regulations from each data provider
- Possible Solution
  - Explore Computed Tomography (CT) scans for image segmentation
  - Federated learning models over data from multiple data providers
- Liver CT Scan Data for Segmentation:
  - The CT Liver dataset consists of 3D NIFTI images or 2D DICOM scans
  - Segmentation masks are the labels
  - Segmentation models for medical scans:
    - UNET
    - nn-UNET
    - Dense-UNET

**Original**  **Windowing**  **Hist. Equalisation**  **Windowing+ Hist. Equalisation**



**Original**  **Liver Mask**  **Original**  **Liver (grey) Tumors(white)**



**No cancer**          **With cancer**

# All project participants

# Thank you for your interest!

Contact: Christoph Lange & Zeyd Boukhers, christoph.lange-bever@fit.fraunhofer.de, zeyd.boukhers@fit.fraunhofer.de

use subject "Open Call"

## Stay in touch:

🌐 www.nfdi.de/fair-data-spaces

Community Newsletter: https://www.nfdi.de/newsletter-abo/

Wiki: https://github.com/FAIR-DS4NFDI/FAIR-DSWiki/wiki

🐦 @FAIRDataSpaces

#FAIRDataSpaces

**FAIR Data Spaces**