

# Session 2.

## A Project perspective on Big Data architectural pipelines and benchmarks



# Panelists



**Arne Berre**

Chief Scientist, SINTEF



**Leonidas Kallipolitis**

*Aegis - I-BiDaaS Projects*



**Brian Elvesæter**

Research scientist, SINTEF  
*TheyBuyForYou Project*



**Athanasios Koumparos**

Senior Software Engineer,  
Vodafone Innovus  
*Track&Know Project*



**Jon Ander Gómez Adrián**

Universitat Politecnica de Valencia  
*DeepHealth Project*



**Caj SÖDERGÅRD**

Professor, VTT Data-driven  
Solutions  
*DataBio Project*

# Introduction to Architectural pipelines

- I-BiDaaS
- TBFY
- Track&Know
- DataBio
- DeepHealth

Processing types

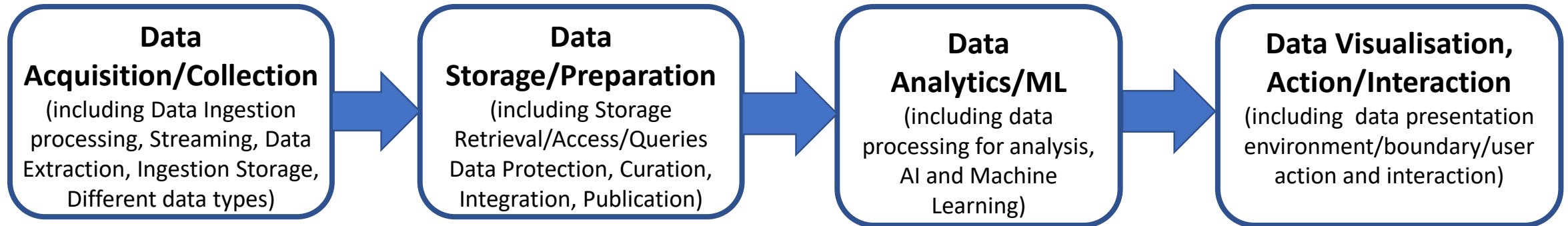
Batch,  
Real time,  
Interactive

Data types

Pipeline steps

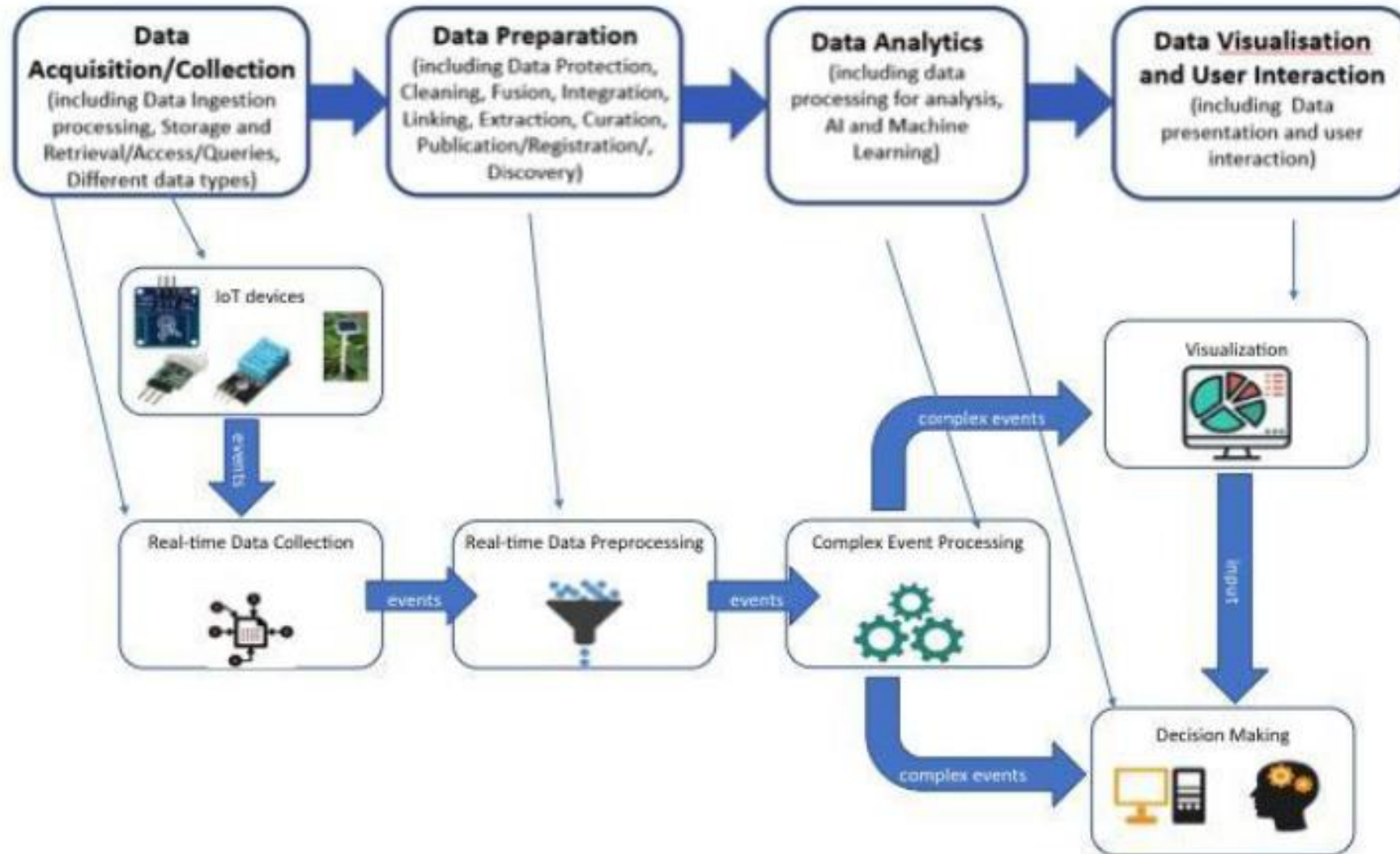
- Benchmarks
- Tools/Components
- Standards

Trans-continuum  
(Edge to Cloud)

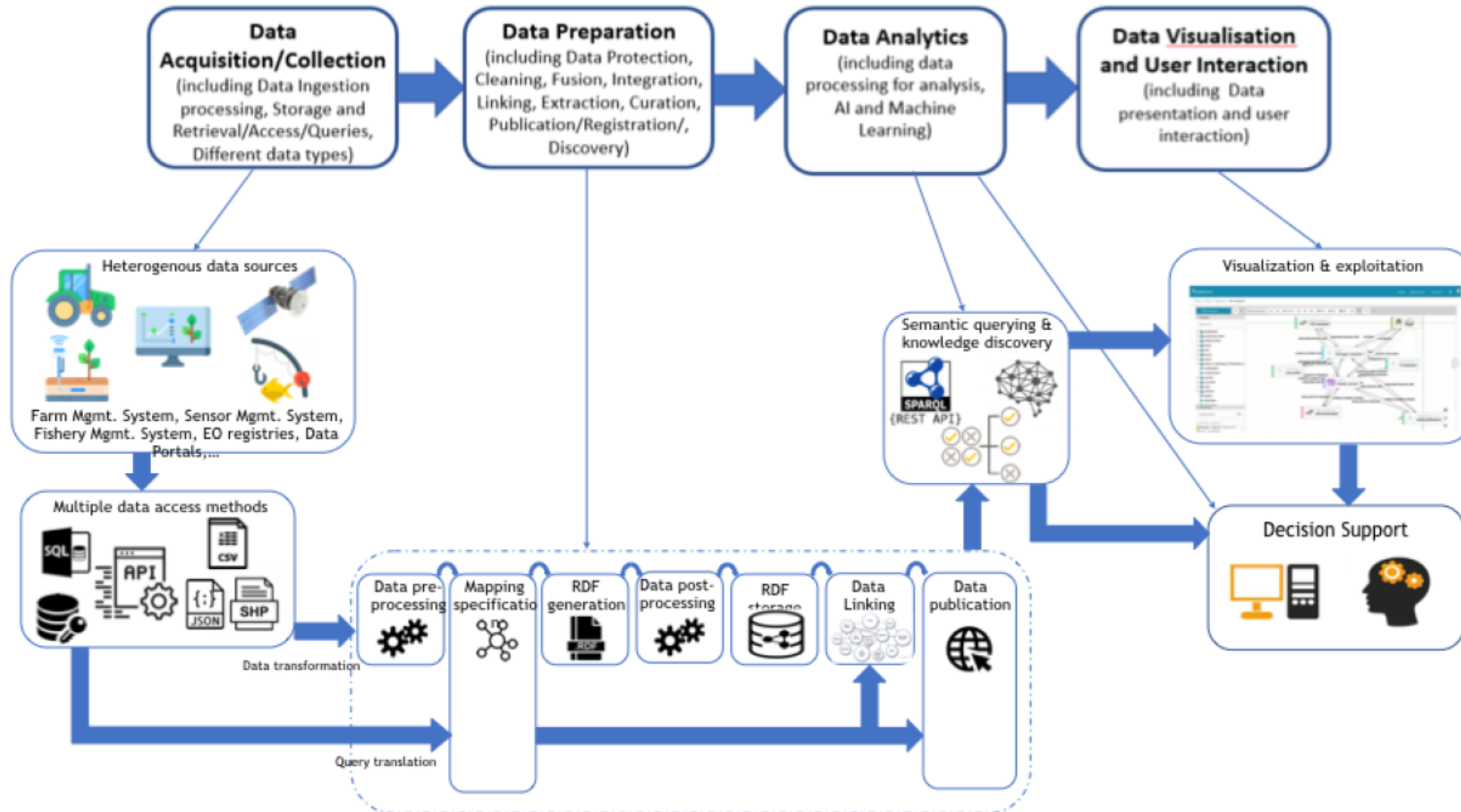




# IoT – Real time Data Pipeline



# Knowledge Graph/Ontology/Linked Data pipeline



**Data Acquisition/Collection**  
(including Data Ingestion processing, Streaming, Data Extraction, Ingestion Storage, Different data types)

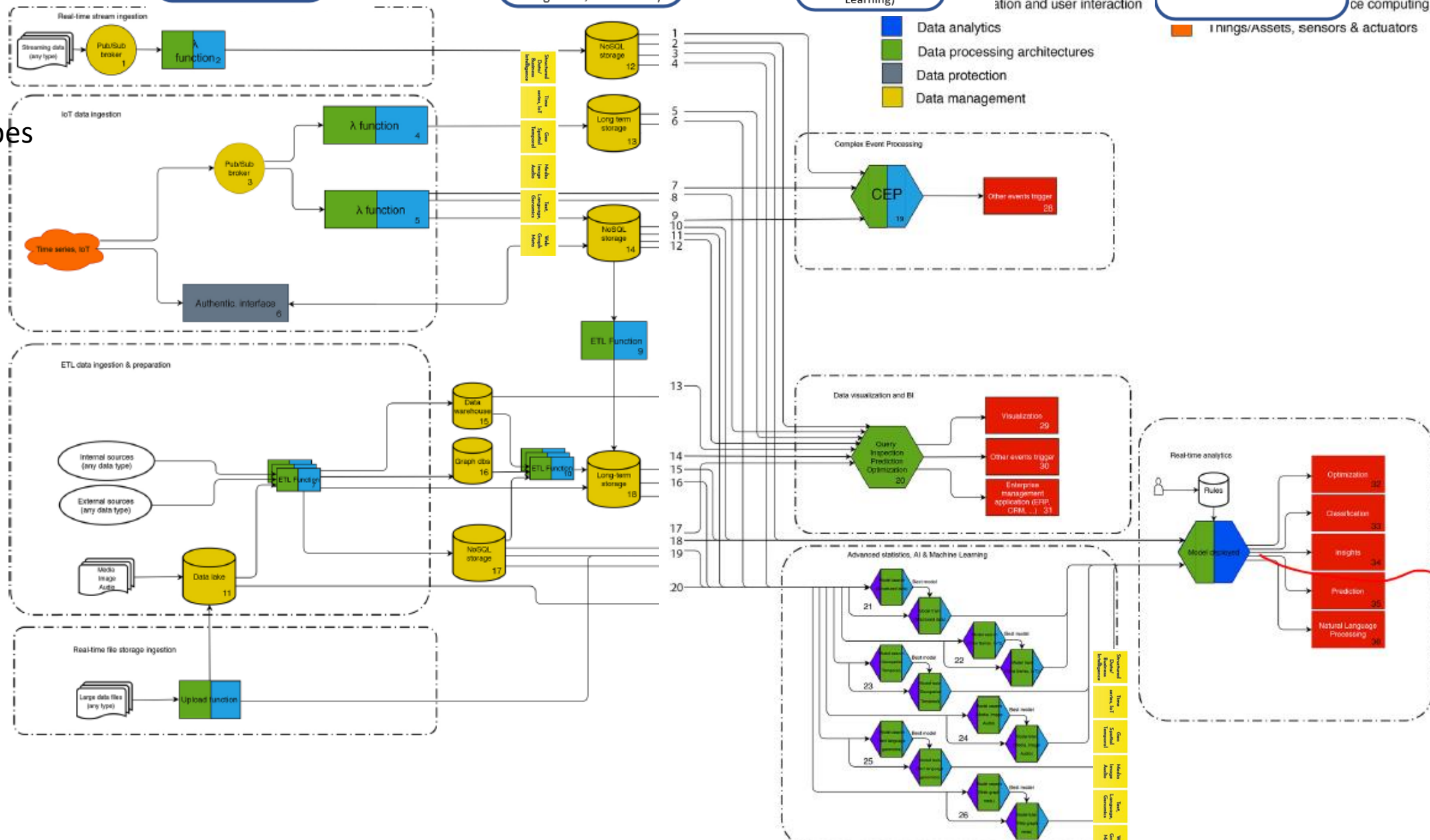
**Data Storage/Preparation**  
(including Storage Retrieval/Access/Queries Data Protection, Curation, Integration, Publication)

**Data Analytics/ML**  
(including data processing for analysis, AI and Machine Learning)

**Data Visualisation, Action/Interaction**  
(including data presentation environment/boundary/user action and interaction)

Realtime

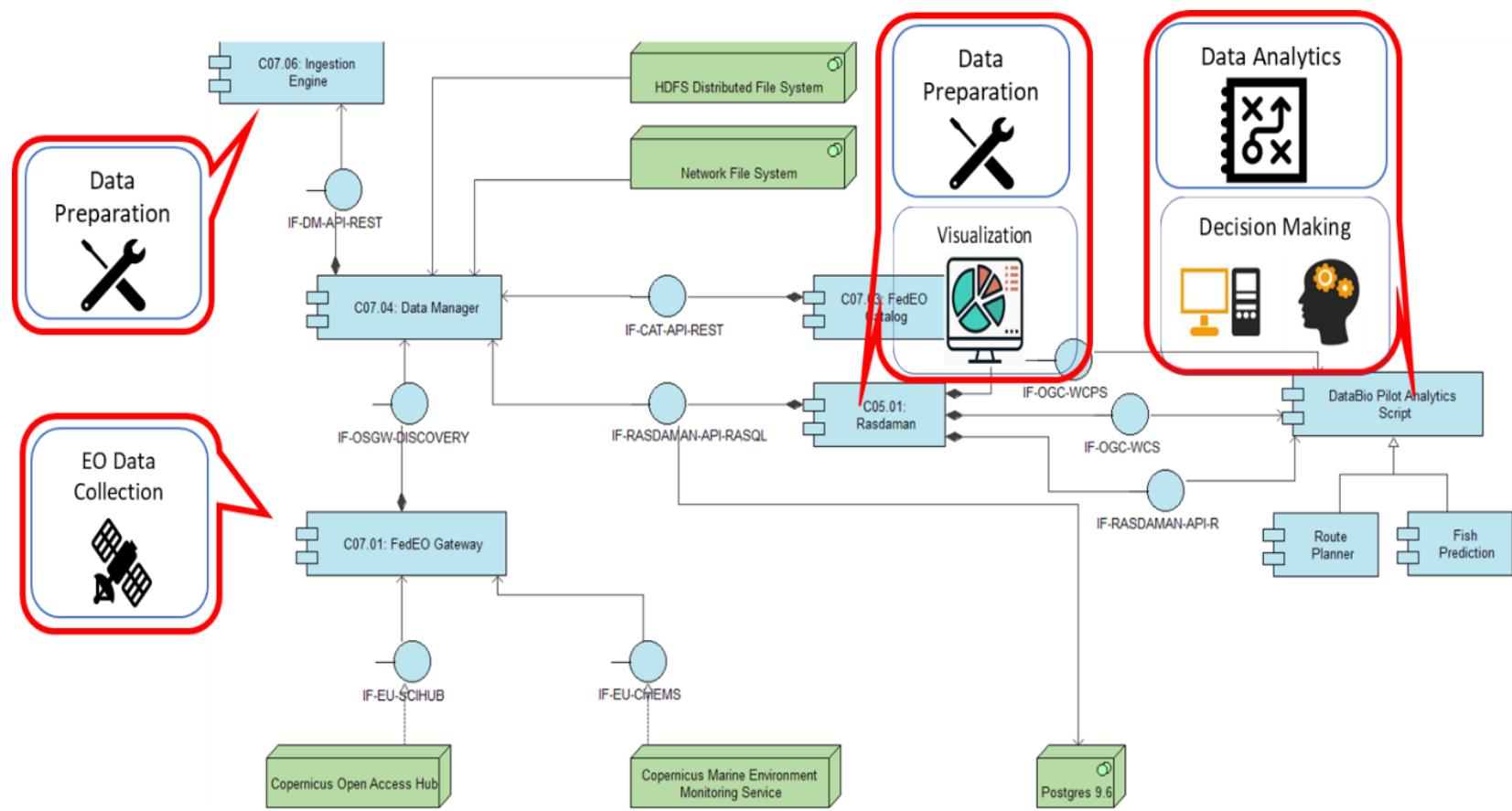
6 data types



Batch

Interactive

Realtime

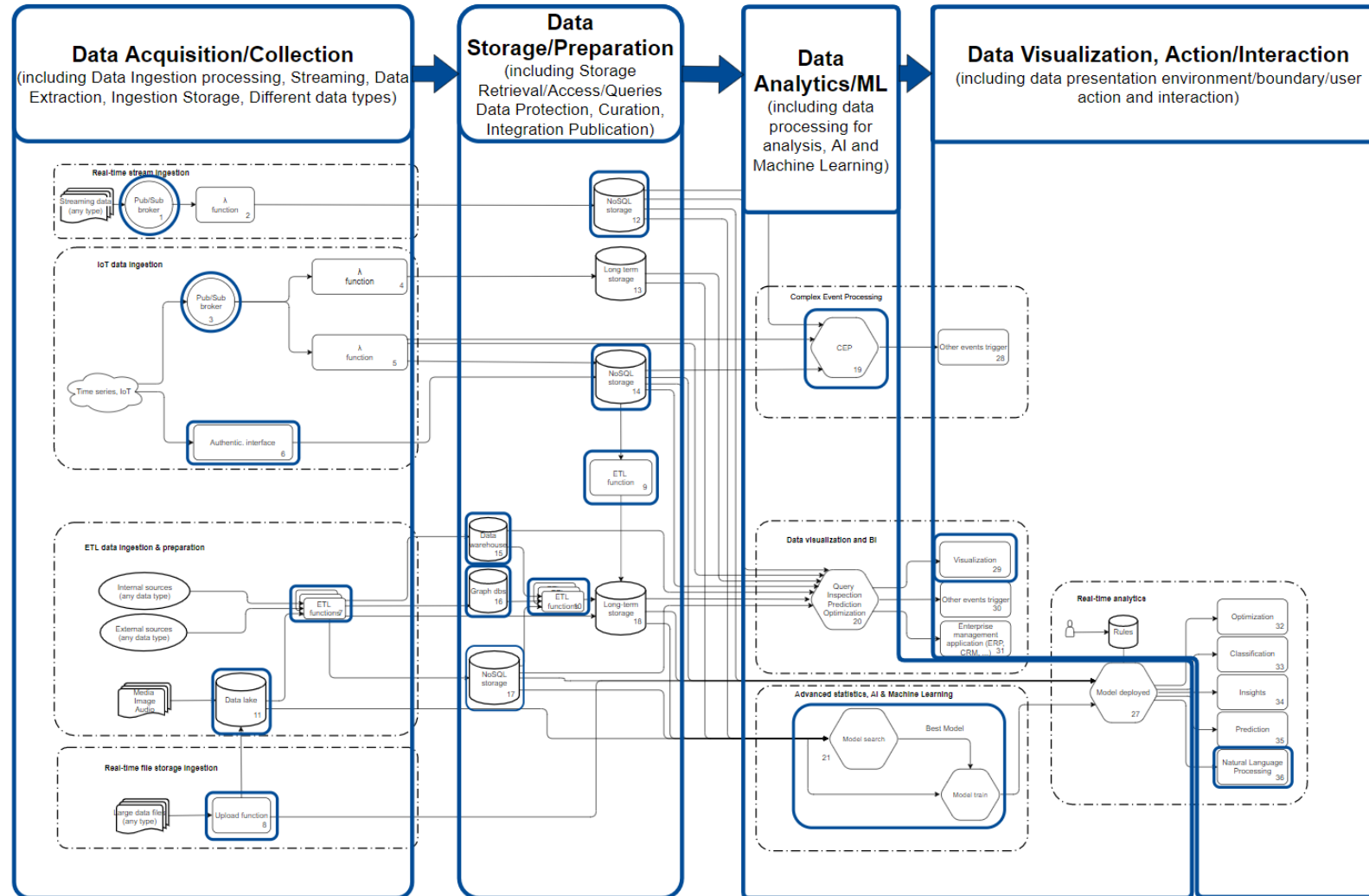


# Search Search by Blueprint/Pipeline

Search by step of the Data Value Chain or for the specific components of the generic big data pipeline generated by DataBench

DataBench has devised a generic architectural blueprint mapped to a top level pipeline along the Data Value Chain covering the steps of **Data Acquisition/Collection** (including data ingestion, processing, streaming, extraction and ingestion storage), **Data Preparation/Storage** (including storage retrieval/access/queries, data protection, curation, integration and publication), **Data Analytics** (including data processing for analysis, AI and Machine Learning) and **Data Visualization/Interaction** (including data presentation, environment/boundary/user action and interaction).

By clicking on the image below on one of the four steps of the pipeline, or in one of the specific elements from the generic blueprint, this search interface will help you discovery benchmarks and associated knowledge (nuggets)



# Conclusion on Pipelines and related benchmarks,

Arne J. Berre, SINTEF

- I-BiDaaS
- TBFY
- Track&Know
- DataBio
- DeepHealth



EUROPEAN  
**BIG DATA  
VALUE** FORUM  
3 - 5 NOVEMBER . 2020 - BERLIN + VIRTUAL

# Industrial-Driven Big Data as a Self- Service Solution

Leonidas Kallipolitis, AEGIS



**Project Consortium**  
13 partners

**PROJECT NAME**  
Industrial-Driven Big Data as a Self-Service Solution

**PROJECT TYPE**  
RIA

**DURATION**  
36 months

**START DATE**  
1 January 2018

**TOTAL BUDGET / TOTAL EC FUNDING**  
€ 4 997 035



<http://www.ibidaas.eu/>



[@ibidaas](https://twitter.com/ibidaas)



<https://www.linkedin.com/in/i-bidaas/>

# Identity Card



# Consortium

1. FOUNDATION FOR RESEARCH AND TECHNOLOGY HELLAS (**FORTH**)
2. BARCELONA SUPERCOMPUTING CENTER - CENTRO NACIONAL DE SUPERCOMPUTACION (**BSC**)
3. IBM ISRAEL - SCIENCE AND TECHNOLOGY LTD (**IBM**)
4. CENTRO RICERCHE FIAT SCPA (**CRF**)
5. SOFTWARE AG (**SAG**)
6. CAIXABANK, S.A (**CAIXA**)
7. THE UNIVERSITY OF MANCHESTER (**UNIMAN**)
8. ECOLE NATIONALE DES PONTS ET CHAUSSEES (**ENPC**)
9. ATOS SPAIN SA (**ATOS**)
10. AEGIS IT RESEARCH LTD (**AEGIS**)
11. INFORMATION TECHNOLOGY FOR MARKET LEADERSHIP (**ITML**)
12. UNIVERSITY OF NOVI SAD FACULTY OF SCIENCES SERBIA (**UNSPMF**)
13. TELEFONICA INVESTIGACION Y DESARROLLO SA (**TID**)



# Key messages



A **complete** and **safe environment** for methodological **big data experimentation**



Tool and services to **increase the quality** of data analytics



A Big Data as a **Self-Service solution** that helps in **breaking silos** and boosts EU's data-driven economy



Tools and services for **fast ingestion and consolidation** of both realistic and fabricated data

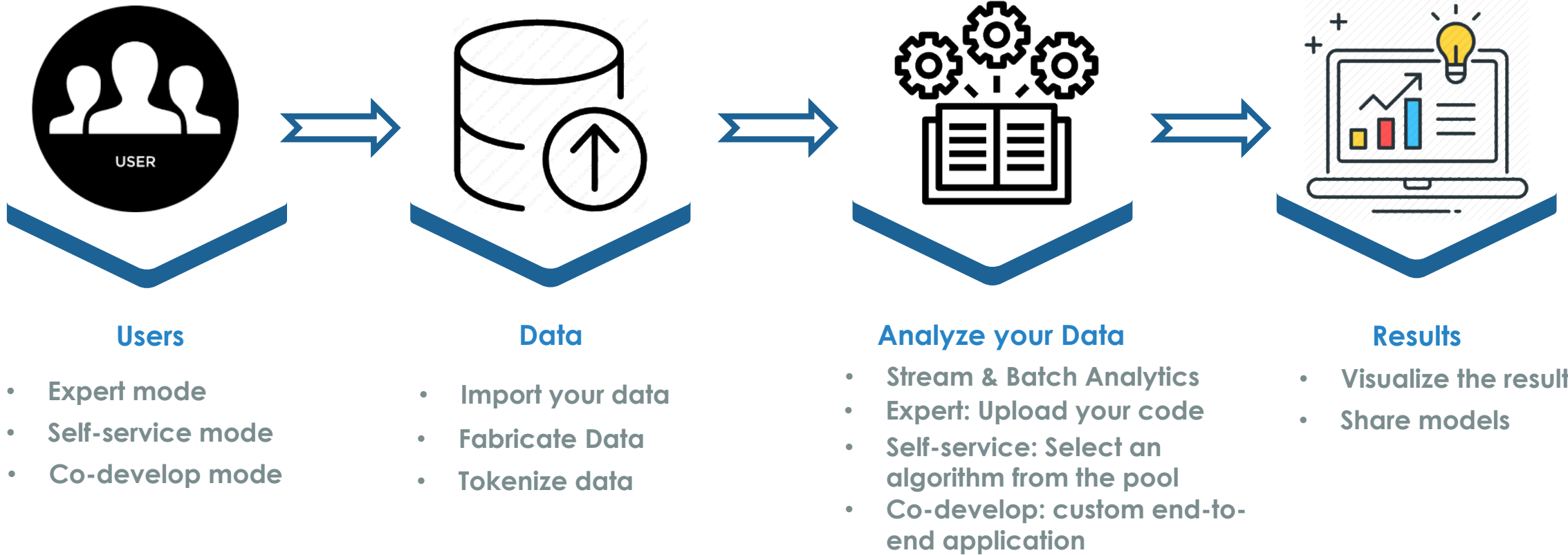


**Increases impact** in research community and contributes to industrial innovation capacity



Tools and services for the management of **heterogeneous infrastructures**

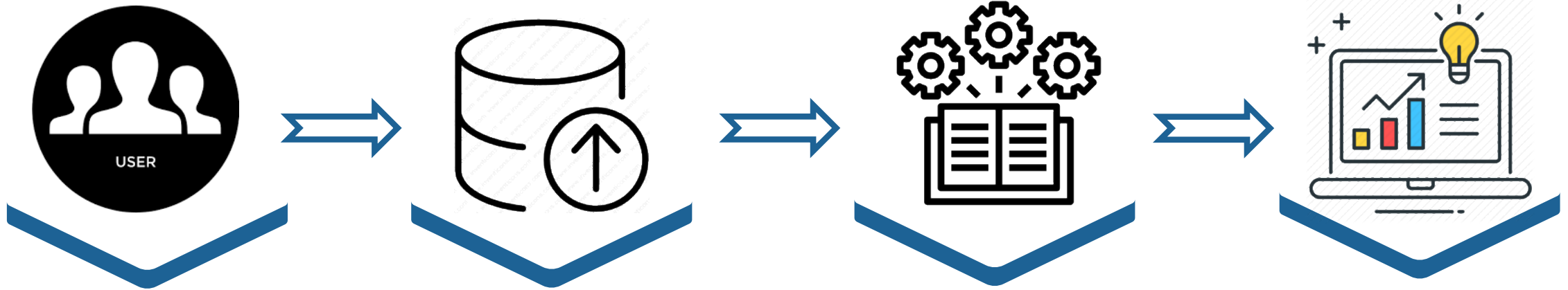
# I-BiDaaS pipeline



## Benefits of using I-BiDaaS



# I-BiDaaS pipeline



## Users

- Expert mode
- Self-service mode
- Co-develop mode

## Data

- Import your data

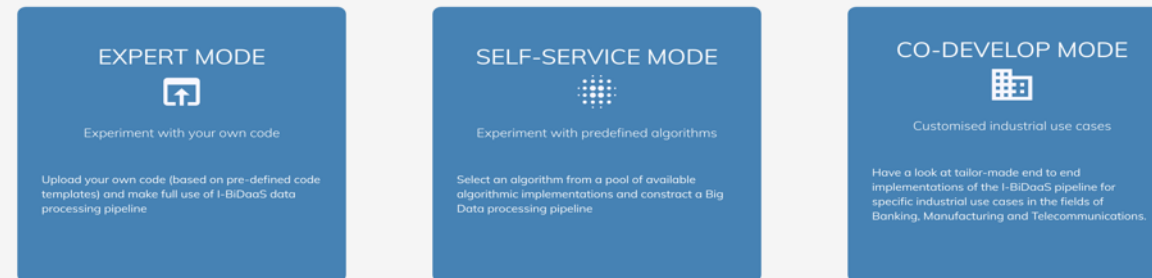
## Analyze your Data

- Stream & Batch Analytics
- Expert: Upload your code

## Results

- Visualize the results
- Share models

## Flexible solution



## Benefits of using I-BiDaaS



Do it yourself  
In a flexible  
manner



Break data silos

Safe environment

Interact with Big Data  
technologies

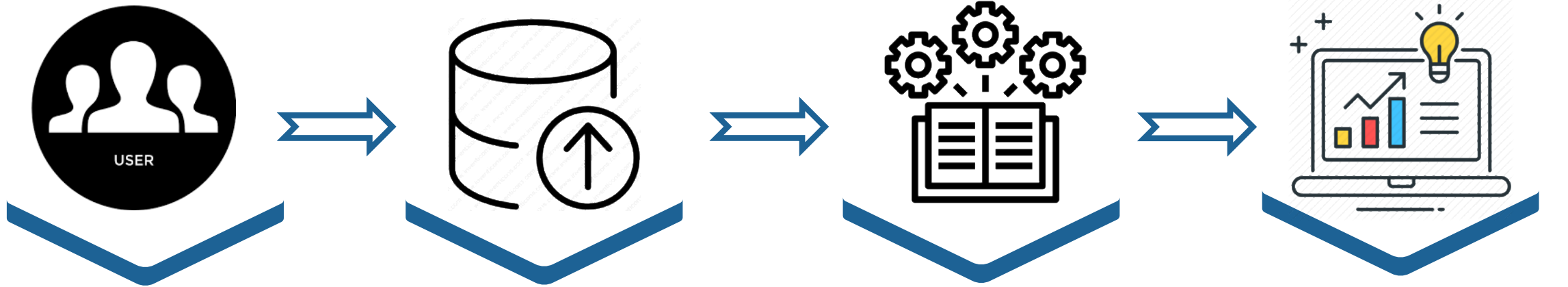
Increase speed of  
data analysis

Intra- and inter-  
domain data-flow



Cope with the rate of data  
asset growth

# I-BiDaaS pipeline



## Users

- Expert mode
- Self-service mode
- Co-develop mode

## Data

- Import your data
- Fabricate Data
- Tokenize data

## Analyze your Data

- Stream & Batch Analytics
- Expert: Upload your code
- Self-service: Select an

## Results

- Visualize the results
- Share models

**Data sharing  
& breaking silos**

## Benefits of using I-BiDaaS



Do it yourself  
In a flexible  
manner



Break data silos



Safe environment



Interact with Big Data  
technologies



Increase speed of  
data analysis

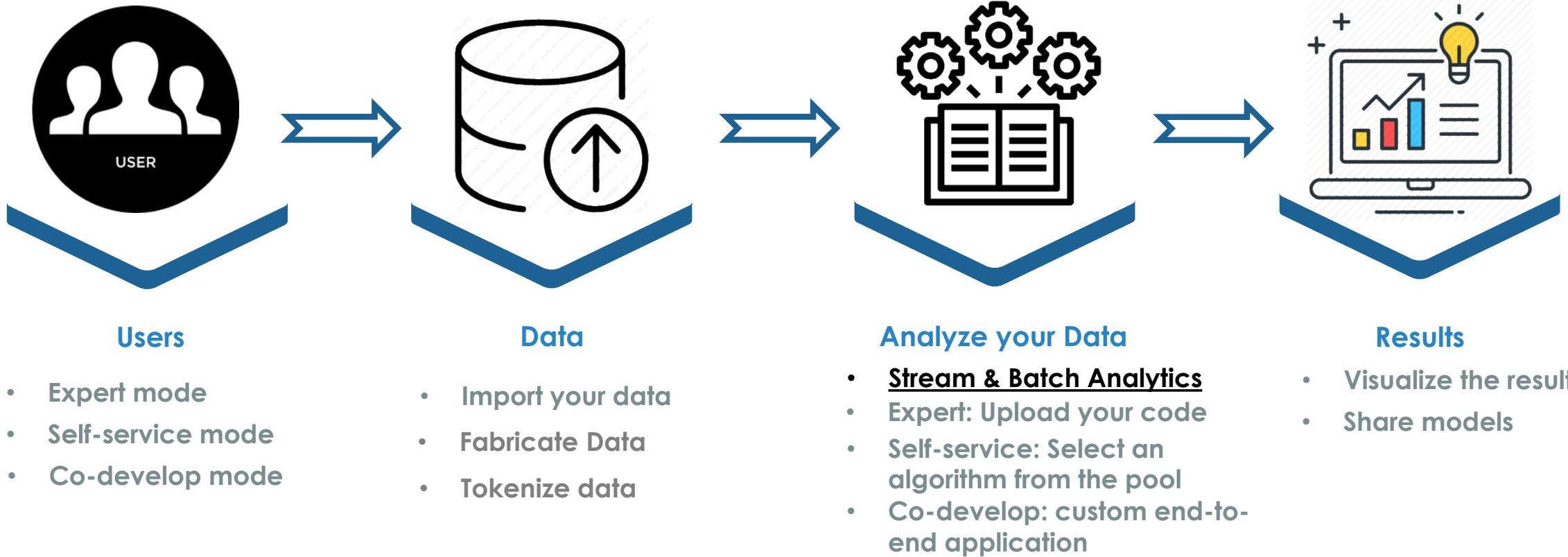


Intra- and inter-  
domain data-flow



Cope with the rate of data  
asset growth

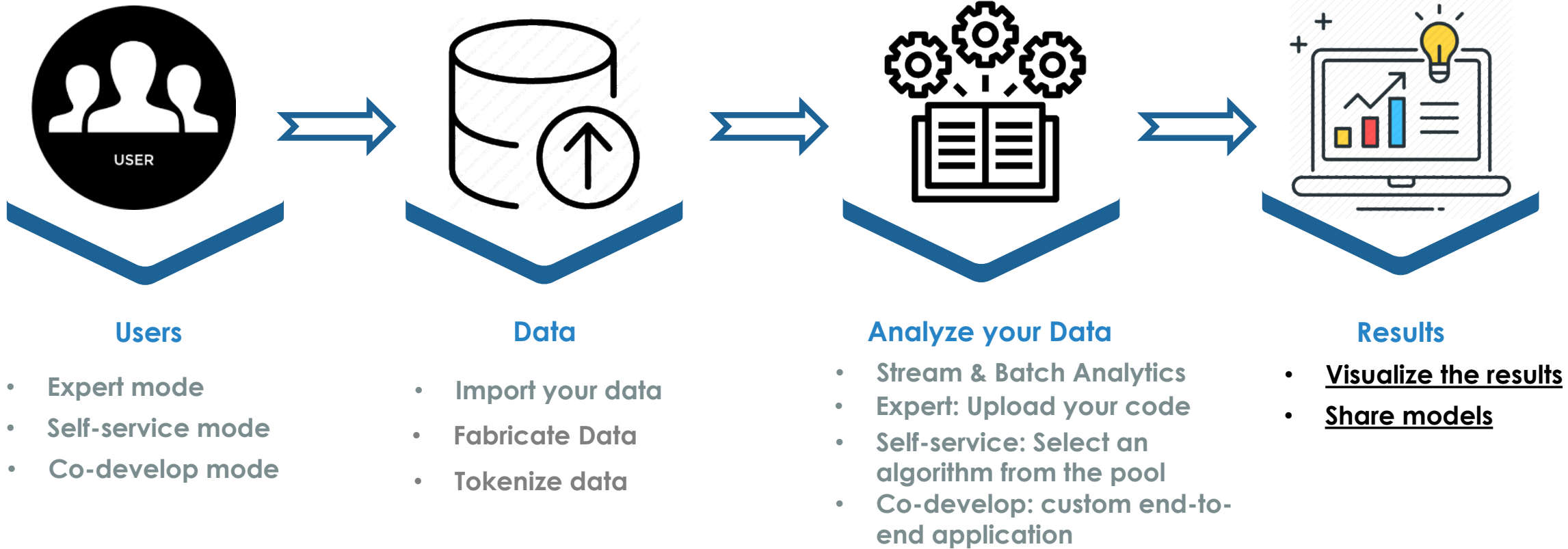
# I-BiDaaS pipeline



## Benefits of using I-BiDaaS



# I-BiDaaS pipeline



## Benefits of using I-BiDaaS



Do it yourself  
In a flexible  
manner



Break data silos



Safe environment



Interact with Big Data  
technologies



Increase speed of  
data analysis



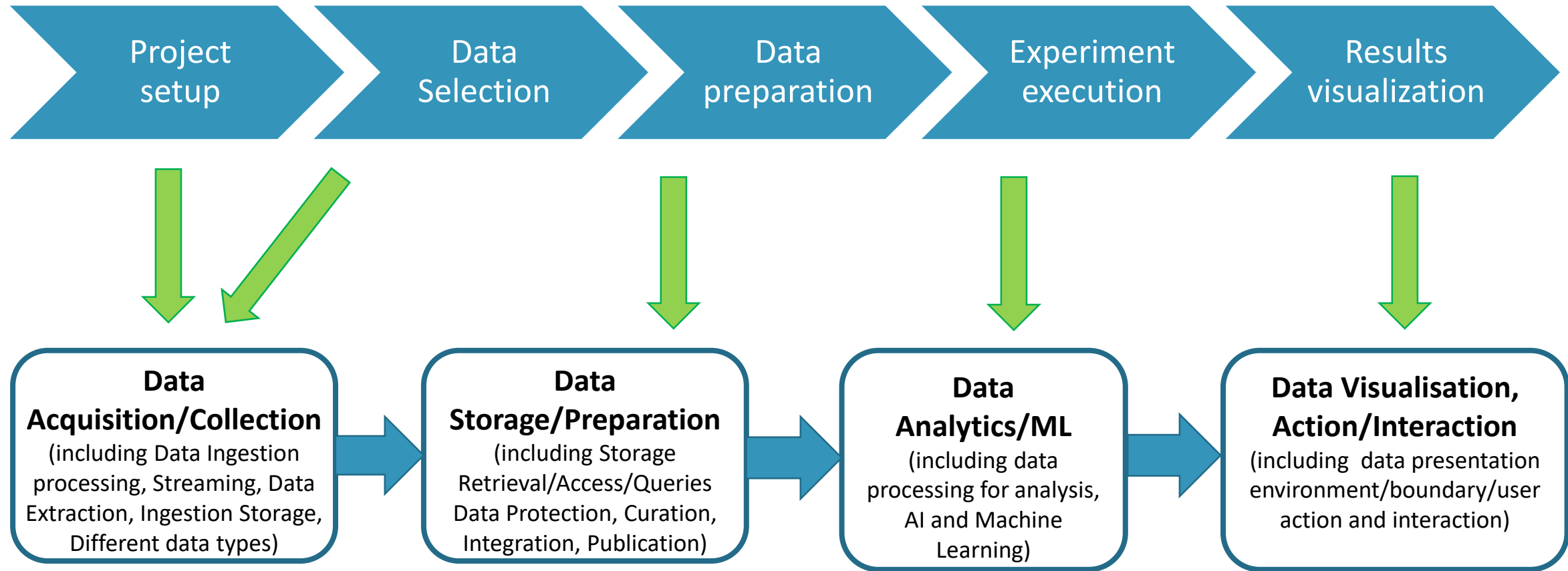
Intra- and inter-  
domain data-flow



Cope with the rate of data  
asset growth



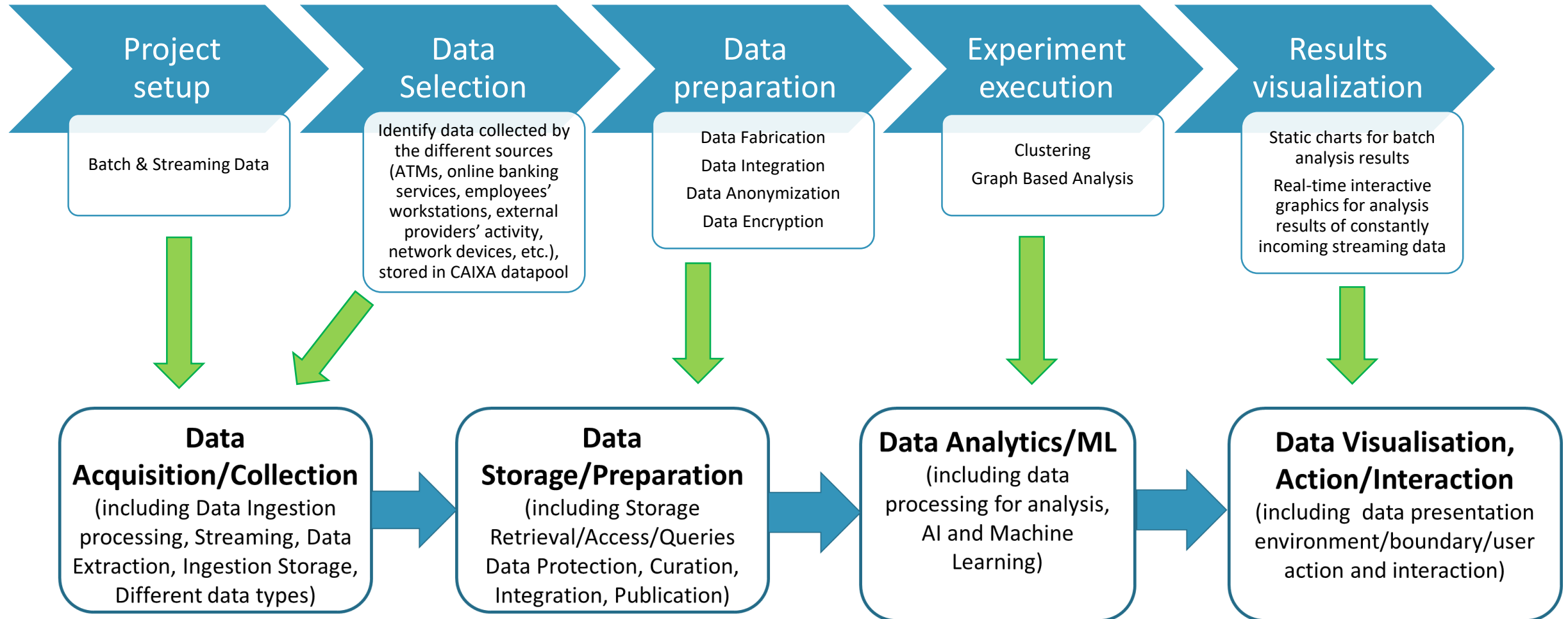
# I-BiDaaS Experimental Workflow



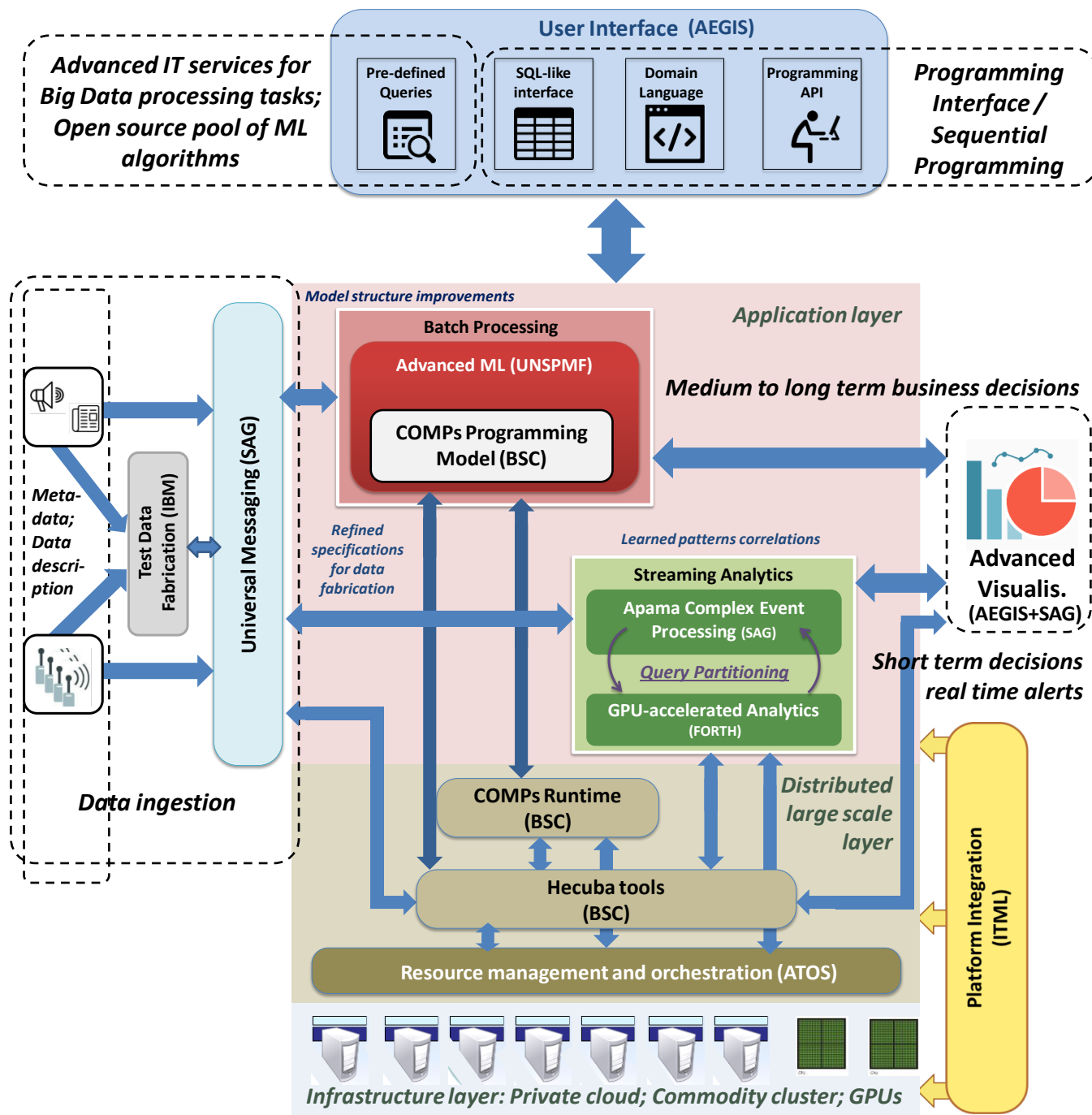
## DataBench Pipeline



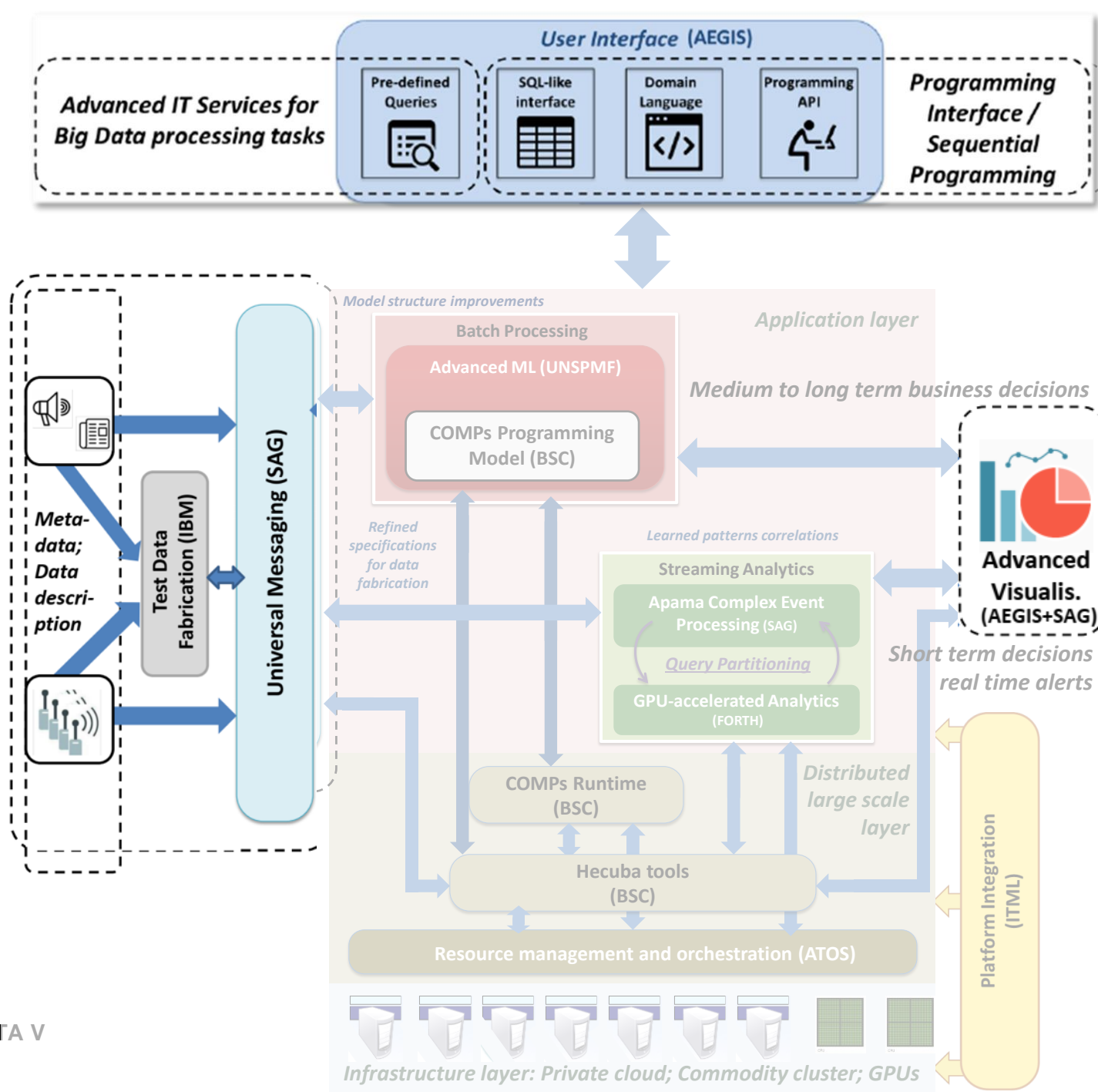
# Example: Banking Experiments Workflow



## DataBench Pipeline



**The I-BiDaaS solution: Architecture & technologie**

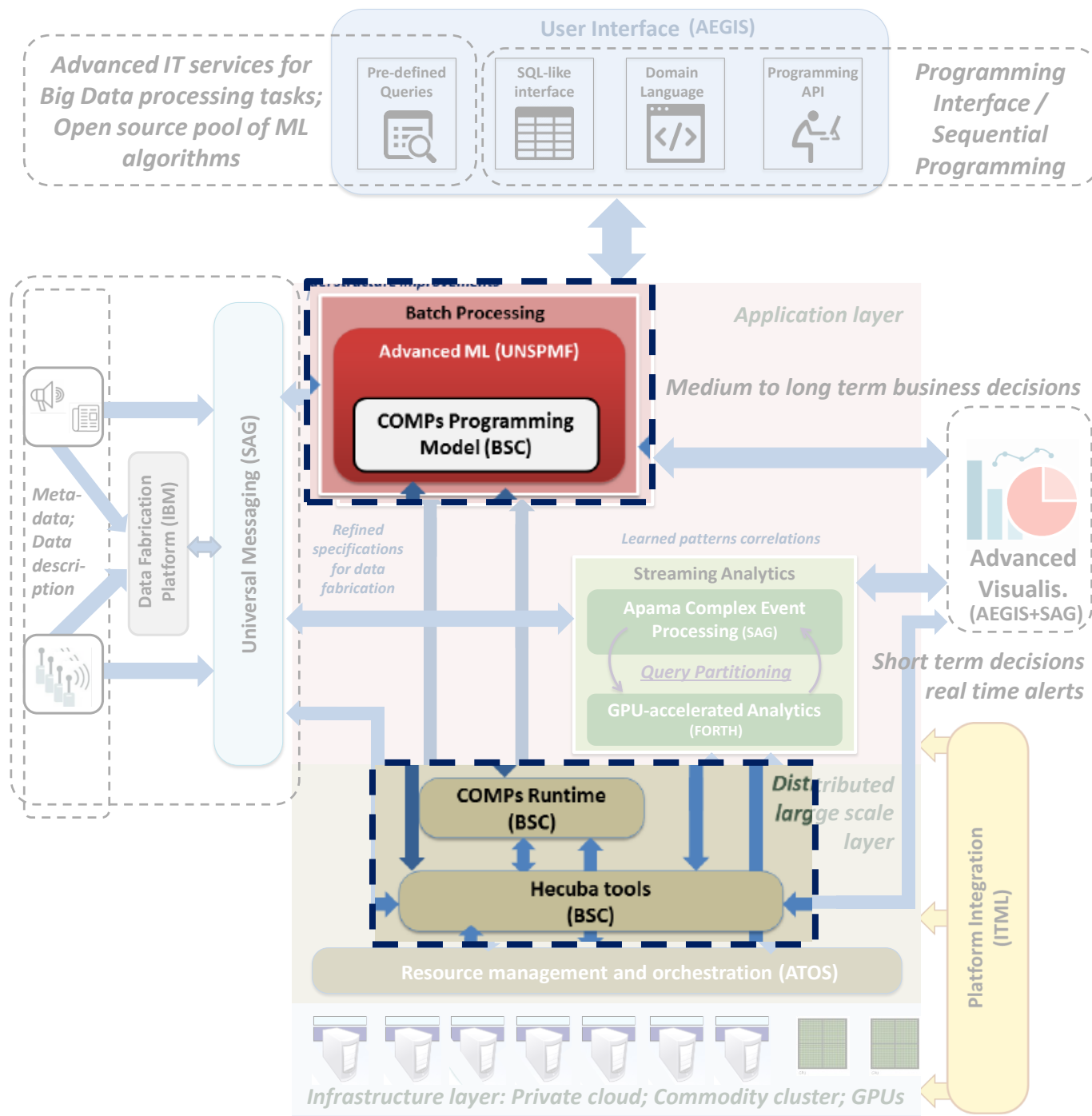


# WP2

## Data, user interface, visualization

## Technologies:

- **IBM TDF**
- **SAG UM**
- **AEGIS AVT**



# WP3 Batch Analytics

## Technologies:

- **BSC COMPSs**
- **BSC Hecuba**
- **BSC Qbeast**
- **Advanced ML (UNSPMF)**

<http://ibidaas.eu/tools>

# Benchmarking: Technology level



I-BiDaaS partner	Technology name	Platform role	DataBench pipeline	Current benchmarks
FORTH	GPU accelerator technology	Data pre-processing, Streaming Analytics	Step 3	Custom benchmark (throughput, latency)
BSC	COMPSs	Sequential programming model for distributed architectures	Step 3	Applications (Own use cases)
BSC	Hecuba	Data management framework with easy interface	Step 2	Applications (Own use cases)
BSC	Qbeast	Multidimensional indexing and storage	Step 2	TPC-H
IBM	Test Data Fabrication	Synthetic test data fabrication	Step 1	Several open source + commercial products (e.g., Grid tools of CA) / No known benchmarks yet
SAG	Apama Streaming Analytics Platform	Streaming Analytics	Step 3	Custom benchmark (throughput)
SAG	Universal Messaging	Message Broker	Step 1	Custom benchmark (throughput)
AEGIS	Advanced visualization and monitoring	Visualization and interface	Step 4	N/A
UNSPMF	Pool of ML algorithms in COMPSs/Python	Batch analytics	Step 3	Respective MPI implementation; Sklearn; HiBench
ATOS	Resource management and orchestration module	Resource management	Step 2	N/A

# Benchmarking: Business, data & analytics

	Business Objectives	Data Sets	Data Size	Processing Type	Type of Analysis
Telecoms	<ul style="list-style-type: none"> <li>- improve and optimize current operations</li> </ul>	<ul style="list-style-type: none"> <li>- Anonymized mobility data (structured)</li> <li>- Anonymized call center data (unstructured)</li> </ul>	TB	<ul style="list-style-type: none"> <li>- batch &amp; streaming</li> </ul>	<ul style="list-style-type: none"> <li>- predictive</li> <li>- descriptive / diagnostic</li> </ul>
Finance	<ul style="list-style-type: none"> <li>- improve decision making</li> <li>- improve efficiency of Big Data solutions</li> </ul>	<ul style="list-style-type: none"> <li>- Tokenized online banking control data (structured)</li> <li>- Tokenized bank transfer data (structured)</li> <li>- Tokenized IP address data (structured)</li> </ul>	PB	<ul style="list-style-type: none"> <li>- batch</li> <li>- batch &amp; streaming</li> </ul>	<ul style="list-style-type: none"> <li>- descriptive / diagnostic</li> </ul>
Manufacturing	<ul style="list-style-type: none"> <li>- improve and optimise current operations</li> <li>- improve the quality of the process and product</li> </ul>	<ul style="list-style-type: none"> <li>- Anonymized SCADA/MES data (structured)</li> <li>- Anonymized Aluminum Die-casting (structured)</li> </ul>	GB	<ul style="list-style-type: none"> <li>- batch</li> <li>- batch &amp; streaming</li> </ul>	<ul style="list-style-type: none"> <li>- predictive</li> <li>- diagnostic</li> </ul>

# Benchmarking: Business level

I-BiDaaS Partner	Use Case	Most relevant business KPIs
TID	Accurate location prediction with high traffic and visibility	<ul style="list-style-type: none"> <li>- Acquisition of insights on the dynamics of cellular sectors</li> <li>- <b>Processing costs (cost reduction)</b></li> <li>- <b>Customer satisfaction</b></li> </ul>
TID	Optimization of placement of telecommunication equipment	
TID	Quality of service in Call Centers	
CAIXA	Enhanced control on online banking	<ul style="list-style-type: none"> <li>- <b>Cost reduction</b></li> <li>- Data accessibility</li> <li>- <b>Time efficiency</b></li> <li>- End-to-end execution time (from data request to data provision)</li> </ul>
CAIXA	Advanced analysis of bank transfer payment in financial terminal	
CAIXA	Analysis of relationships through IP addresses	
CRF	Production process of aluminium die-casting	<ul style="list-style-type: none"> <li>- 3 <b>Product quality</b> levels (High, Medium, Low)</li> <li>- Overall Equipment Effectiveness (OEE),</li> <li>- Maintenance cost</li> <li>- <b>Cost reduction</b></li> </ul>
CRF	Maintenance and monitoring of production assets	





I-BiDaaS aims to empower IT and non-IT big data experts to easily utilize and interact with big data technologies.



- [www.ibidaas.eu](http://www.ibidaas.eu)
- [twitter.com/ibidaas](https://twitter.com/ibidaas)
- [linkedin.com/in/i-bidaas](https://www.linkedin.com/company/i-bidaas)
- [zenodo.org/communities/i-bidaas](https://zenodo.org/communities/i-bidaas)



This project has received funding from the European Union's Horizon 2020 Research and Innovation program under grant agreement **No 780787**.



We made applications for **3 SECTORS**

- Financial
- Telecommunications
- Manufacturing



**PARTNERS**



**USE CASES**



**SOFTWARE COMPONENTS**

5 Open Source  
6 Proprietary

Thank you!  
Questions?



# Enabling procurement data value chains for economic development, demand management, competitive markets and vendor intelligence

---

**Brian Elvesæter, SINTEF**

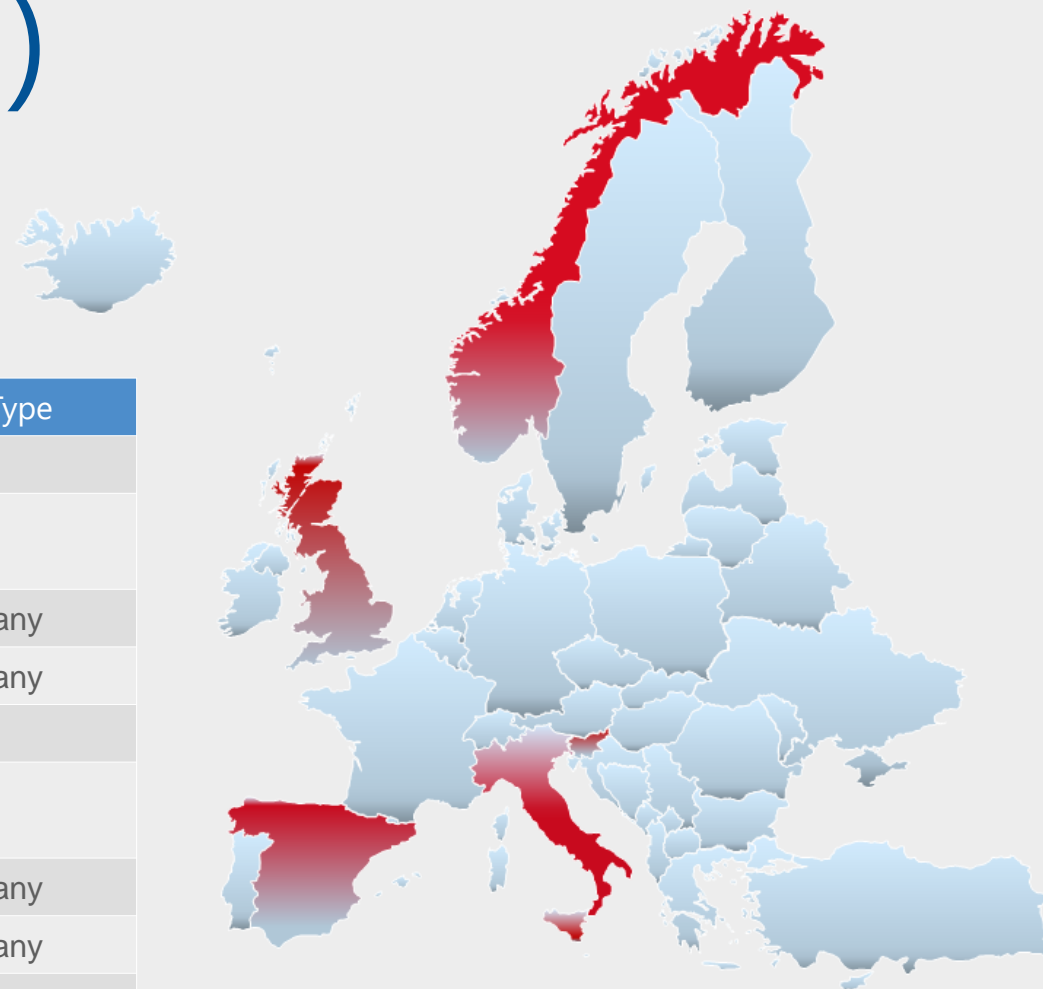
brian.elvesater@sintef.no

<https://theybuyforyou.eu>

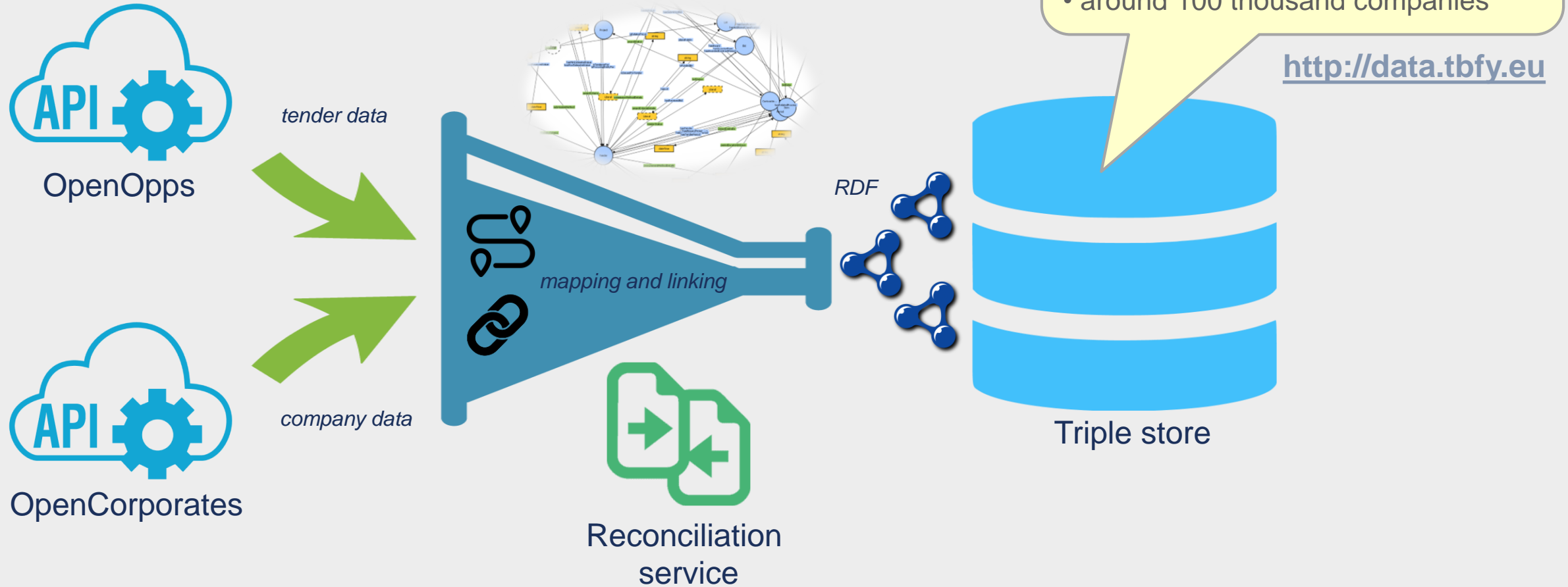
# TheyBuyForYou (TBFY)

- Horizon 2020 Innovation Action
  - Grant agreement No 780247
- Duration: 36 months
  - Jan 2018 – Dec 2020
- Overall budget
  - € 3 274 440
- Developing Big Data tools for Public Procurement
  - data access
  - data analytics
  - data interaction
  - data visualization
- 10 Partners from 5 countries

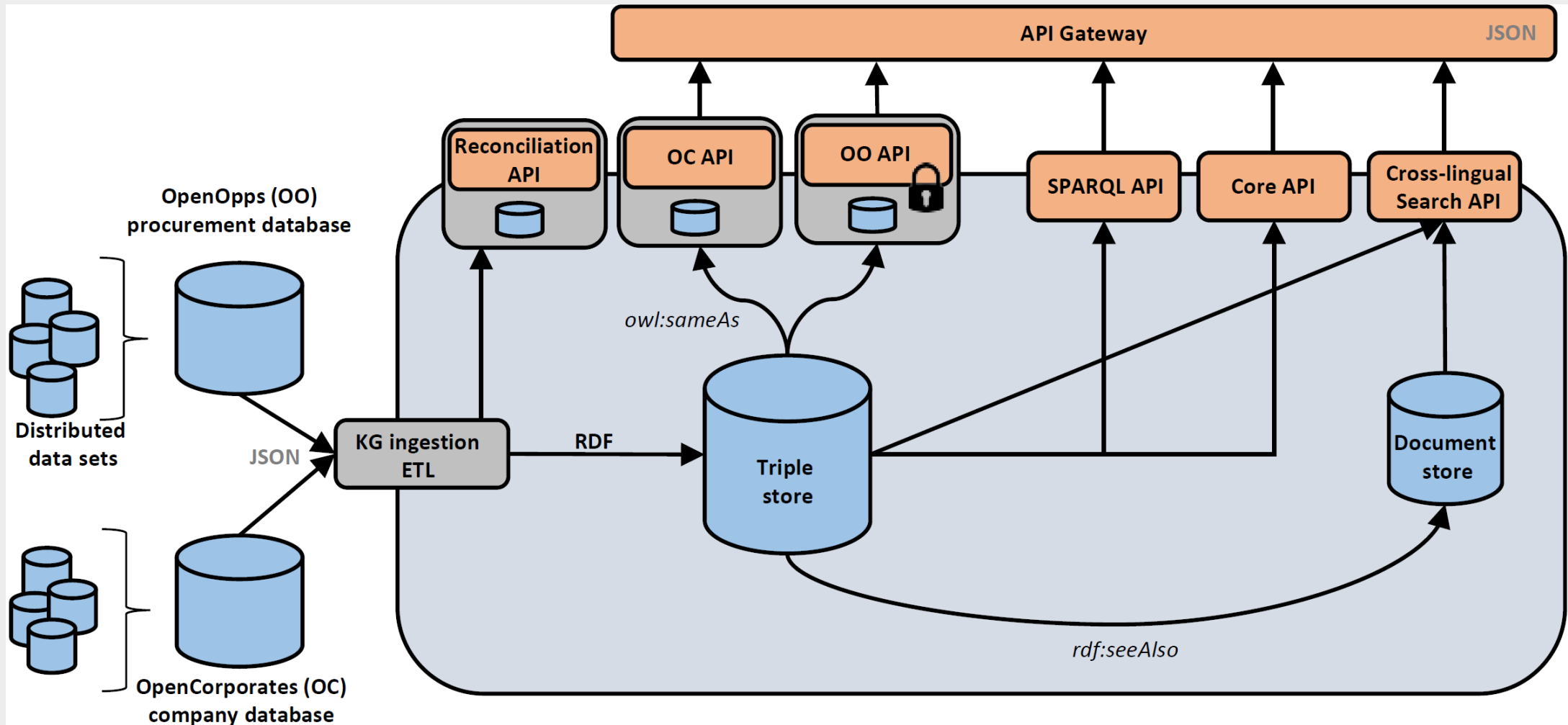
Partner	Country	Organisation Type
SINTEF AS	NO	RTO
Municipality of Zaragoza	ES	Public sector
CERVED	IT	Private Company
OpenCorporates	UK	Private Company
Josef Stefan Institute	SI	RTO
Ministry of Public Sector Innovation	SI	Public Sector
OESÍA Networks	ES	Private Company
OpenOpps	UK	Private Company
Universidad Politécnica de Madrid	ES	University
King's College London	UK	University



# TBFY Knowledge Graph



# TBFY Platform



# Benchmark case #2: Query/visualization



- SPARQL Query Performance
  - KG API (RESTful API wrapping pre-defined SPARQL queries)
  - Visualization
- Benchmark questions?
  - Which triple store (RDF) database to choose?
  - Which cloud computing (resource) plan to choose?
    - <https://aws.amazon.com/pricing/>
  - Database stability issues
    - JVM heap size
- DataBench ToolBox
  - LDDB Semantic Publishing Benchmark (SPB) to benchmark RDF triple stores
    - GraphDB
    - Apache Jena Fuseki & TDB
    - ...

**THEY BUY Knowledge Graph & Statistics & Tender**

**DataBench** | Benchmarks | Knowledge Nuggets | Other Tools | Search

**Semantic Publishing Benchmark (SPB)**

Navigation: <- Back

Tags: Gigabytes (40) | Data Management (38) | Execution time (68) | Synthetic data (55) | Graphs or linked data (24) | Batch (35) | Interactive/near/Real-time (19) | Execution performance (50) | Graph Databases (16) | Data analytics (52) | Data processing architectures (48)

Sub-tags: Data acquisition/Collection (64) | Data preparation (71) | Data analytics (49) | NoSQL (23) | Graph (17)

**Description**

Semantic Publishing Benchmark (SPB) v2.0 is a LDDB benchmark for RDF database engines inspired by the Media/Publishing industry, particularly by the BBC's Dynamic Semantic Publishing approach.

The application scenario considers a media or a publishing organization that deals with large volume of *streaming content*, namely news, articles or "media assets". This content is enriched with *meta* links to *reference knowledge* – taxonomies and databases that include relevant concepts, entities and factual information. This metadata allows publishers to efficiently retrieve relevant content, across business models. For instance, some, like the BBC, can use it to maintain rich and interactive web-presence for their content, while others, e.g. news agencies, would be able to provide better defined content.

From a technology standpoint, the benchmark assumes that an RDF database is used to store both the reference knowledge (mostly static) and the metadata (that grows constantly, to stay in sync with content). The main interactions with the repository are (i) *updates*, that add new metadata or alter it, and (ii) *queries*, that retrieve content according to various criteria.

**Web references**

<http://ldbouncil.org/developer/spb>

**Date of last description update**

2018

**Originating group**

LDDB

**Time – first version, last version**

2013-2018

**Type/Domain**

Graph benchmark

**Workload**

Basic: Consisting of an interactive query-mix for evaluation RDF systems in most common use-cases


# Questions?

---

**Brian Elvesæter, SINTEF**

brian.elvesater@sintef.no

<https://theybuyforyou.eu>



# Vodafone Innovus Fleet Management solutions

Athanasios Koumparos  
Vodafone Innovus

European Big Data Value Forum, 2020.11.04





# Getting to know Vodafone Innovus



## Who we are

a 100% Vodafone company located in Athens, Greece, est. 2004

Operating in 7 Vodafone OpCos



## What we do

We design and implement innovative IoT solutions based on **VF Group strategy & customer needs**



## What we deliver

We have built **Vodafone Group products** as well as the **Global Device Management solution**



## Our advantages

- ✓ World class Platform
- ✓ Global & local expertise
- ✓ Operations experience
- ✓ Agile way of working

World class capabilities build a world class eco-system





# Enterprise Fleet | Overview

Logistics  
Vertical

Soho / SME /  
Corporate  
segments

Developed by  
VF Innovus

Refreshed  
platform,  
GDSP  
connected

A well-established solution, applicable to customers with  
*advanced needs for location tracking, sensor monitoring & in-depth reports / analytics !*

## Markets deployed

VF Greece & VF Albania

>13k active assets  
>6k assets in cold chain  
logistics

## 11 years of experience

providing an excellent  
solution since 2009

## A fresh, award winning approach

new platform with re-designed  
UI, advanced Driver Safety &  
enhanced capabilities



## Innovation

We participate in EU programs  
(Track & Know, Truconomy) and  
collaborate with the Hellenic  
Organization of Intelligent Transport  
Systems to enhance the solution  
capabilities

API, GUI

Analytics – Offline Processing

Online Data Processing

Low Level Device – Platform Communication

GSM Network



# Devices on vehicles



## Device in vehicles

- Devices are installed under the hood (or bonnet)
- Vehicles vary in type, age and electronics

## Data

- More than 3M data packets received daily by the platform
- GPS data (coordinates, angle, speed, date)





# Challenges

## Reception problems

- Low signal strength due to installation (under the bonnet)
- Urban areas (bridges, tall buildings etc.)
- GPS signal reflections

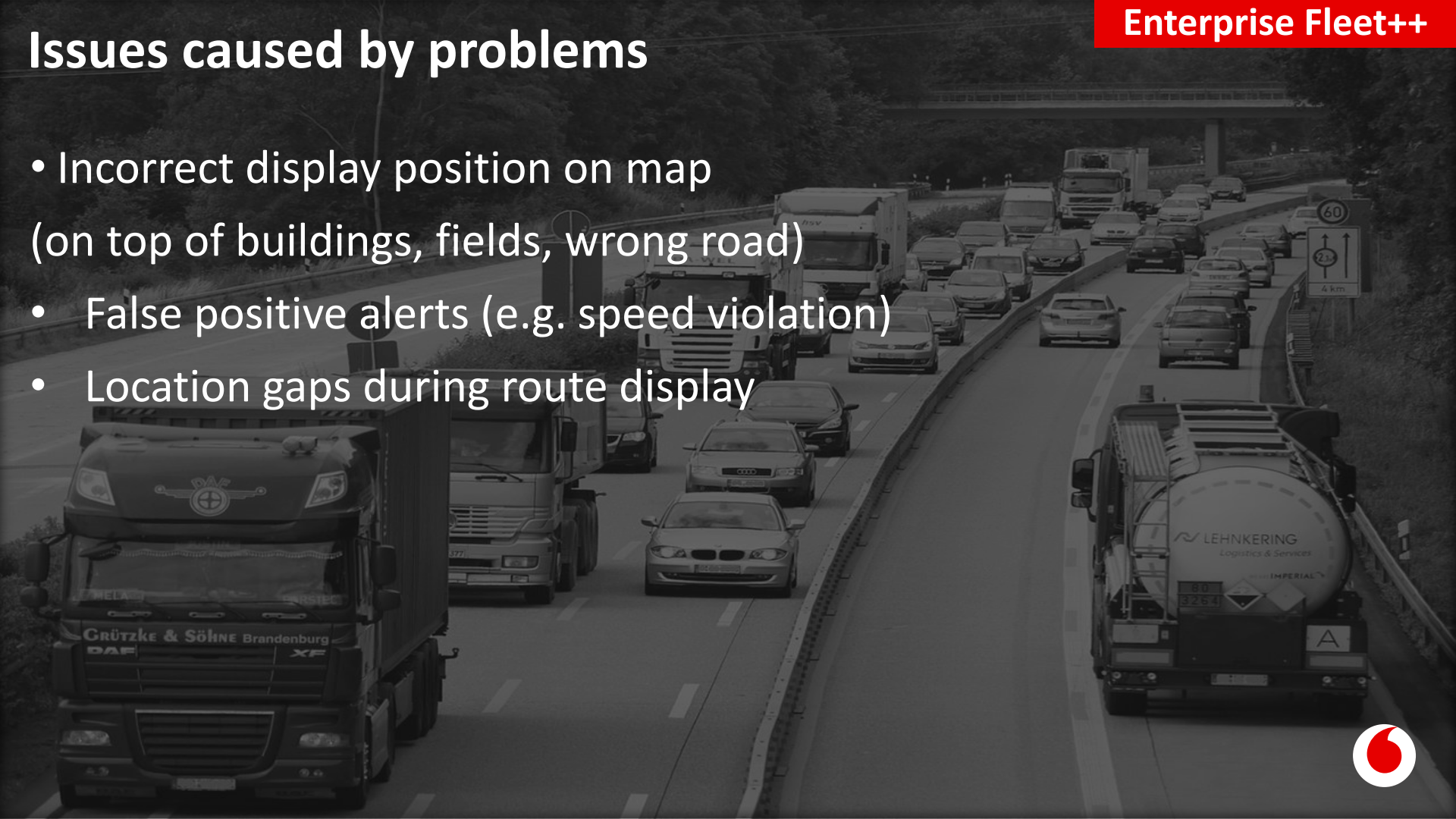
## Vehicle

- Poor quality of power supply
- Vehicle aging – noisy electricals



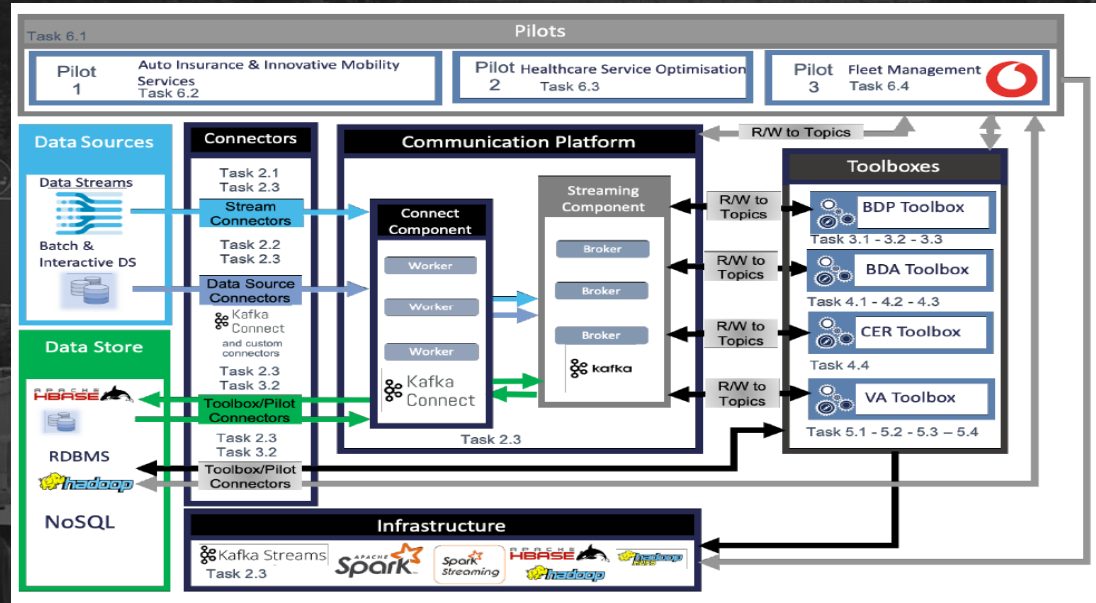
# Issues caused by problems

- Incorrect display position on map  
(on top of buildings, fields, wrong road)
- False positive alerts (e.g. speed violation)
- Location gaps during route display



# Track & Know

- Big Data Analytics
- ML based components
- Scalable
- Online Processing
- Offline Processing
- Versatile
- Easy to integrate (Kafka Topics)



### Data Collection / Collection

Live data streaming from devices (GPS/Sensors)

### Data Storage / Preparation

Cleansing, enrichment, storage in NoSQL stores

### Data Analytics / ML

Pattern recognition  
Location forecasting  
Clustering  
Mobility networks

### Data Visualization Action / Interaction

Visual analytics  
End user GUI





# Integration VFI IoT Platform to Track & Know

Enterprise Fleet++

## Online Processing Flows:

- Data cleansing
- Data enrichment
- Data pattern recognition
- Future location prediction
- Driver behavior analysis

## Offline Processing Flows

- Individual mobility networks, predict next service period
- Hot spot analysis
- Trajectory matching
- Visualization

API, GUI

Analytics – Offline  
Processing



Online  
Track & Know

Online Processing



Offline  
Track & Know

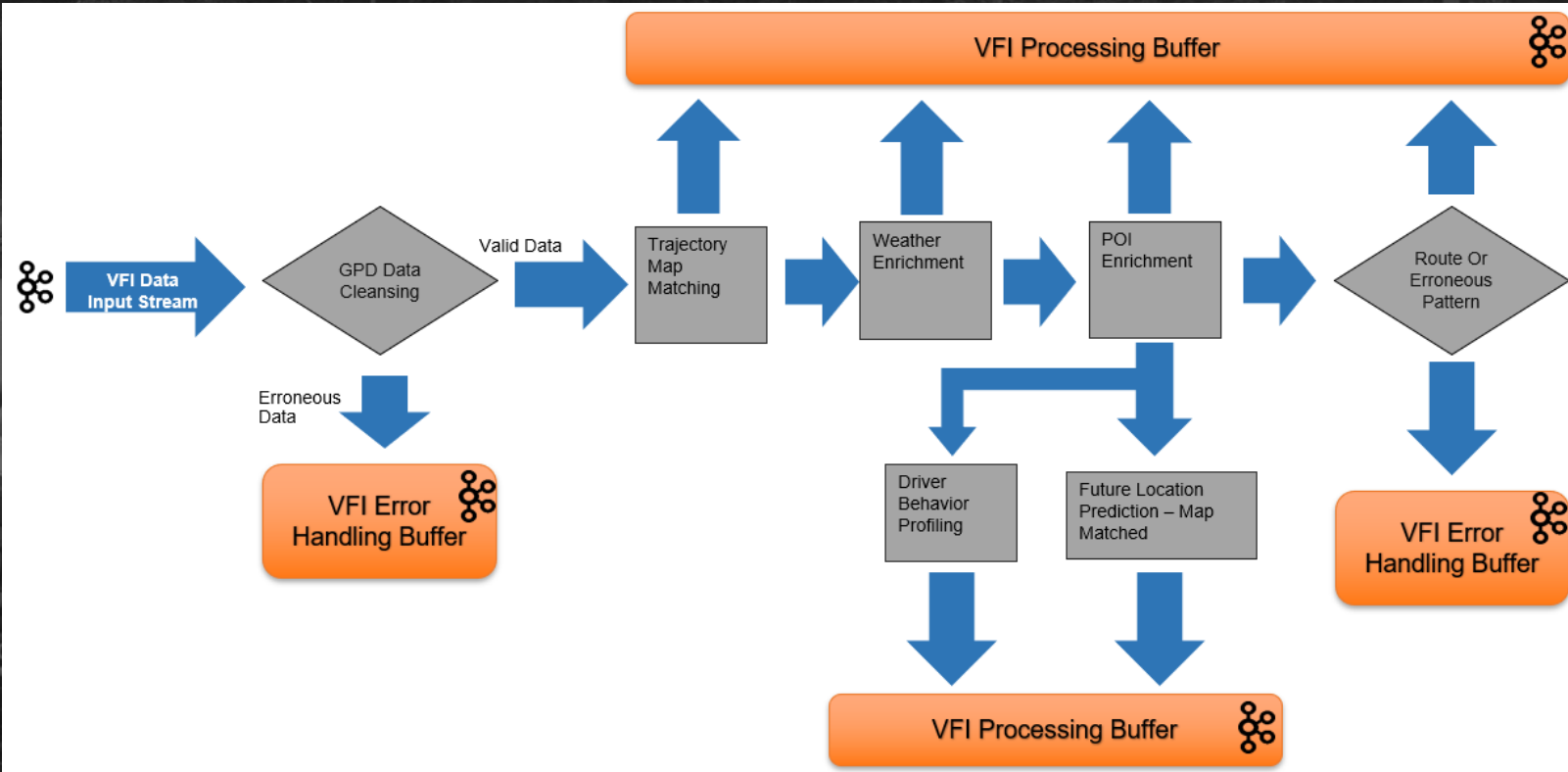
Low Level Device – Platform Communication

GSM Network

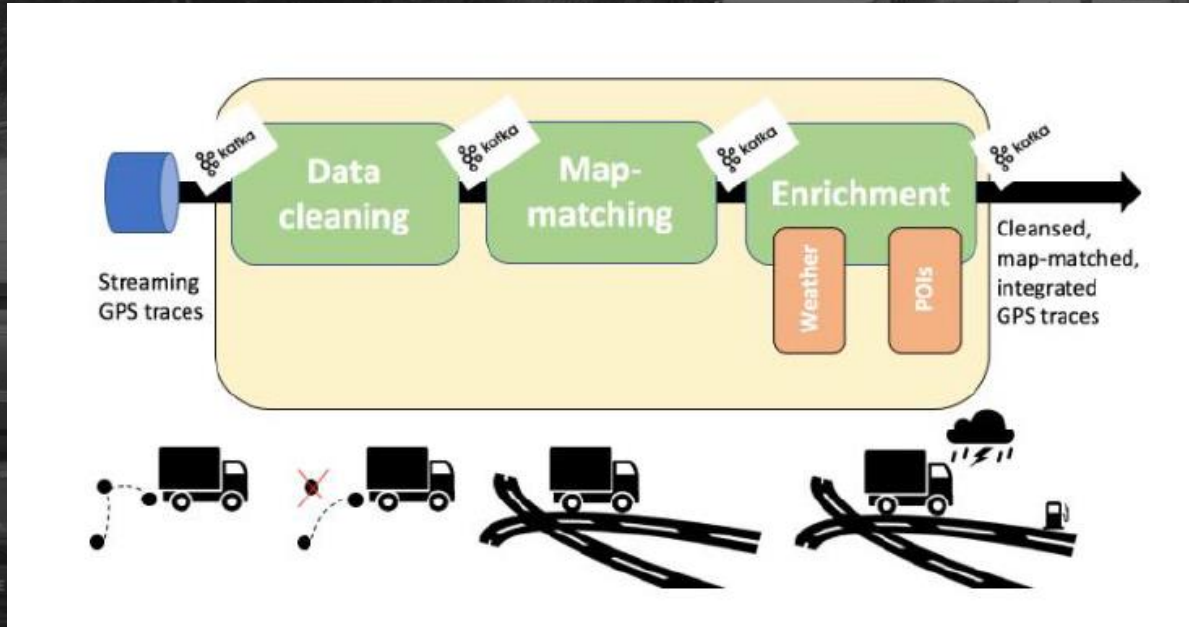




# Data Flows T&K Platform

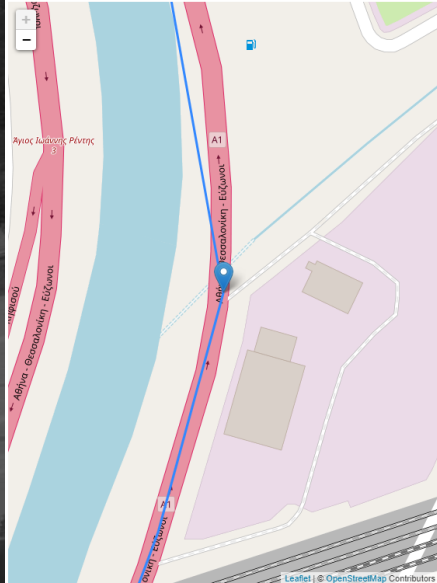


# Data cleansing tool - Map matching tool

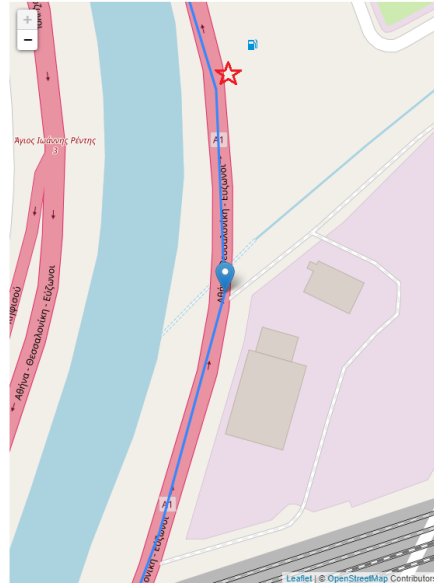


# GPS post enhancements

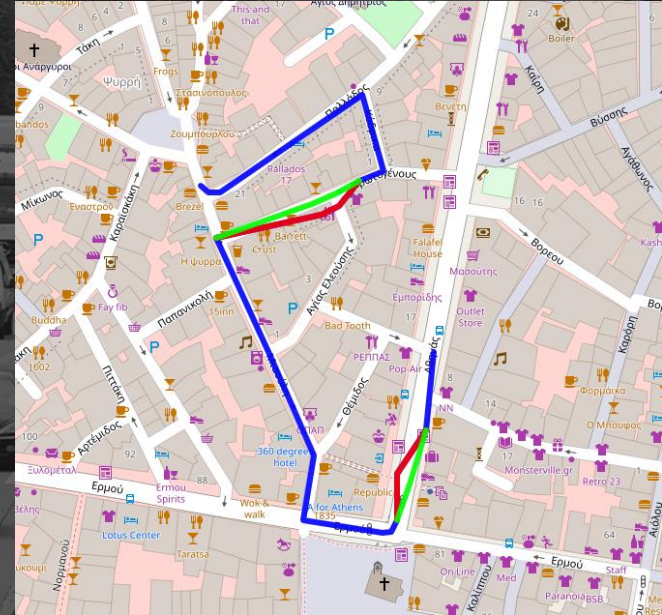
Raw data



Including predicted data



- Vehicle moving on a highway
- Sampling rate at 30 seconds



- Vehicle moving in dense urban area
- GPS reception not always good







Thank you





This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 732064



BIG DATA VALUE  
PUBLIC-PRIVATE PARTNERSHIP

This project is part  
of BDV PPP

## DATA BIO PIPELINES

**Prof. Dr. Caj Södergård,**  
Technical Manager of DataBio

*European Big Data Forum 2020  
DataBench session  
4.11.2020*

**VTT**

VTT Technical Research Centre of Finland

# Data-driven Bioeconomy - DataBio

DataBio aim: Develop big data tools for enhancing production of raw materials for food, energy and biomaterials industries

- A EU-funded project with 48 partners
- 27 bioeconomy pilots in 17 countries
- Duration 2017-2020



**Responsible and sustainable production of food, energy and biomaterials**



**Better raw material utilisation from agriculture, forestry and fishery sources**



**New business opportunities through market-ready big data technologies**

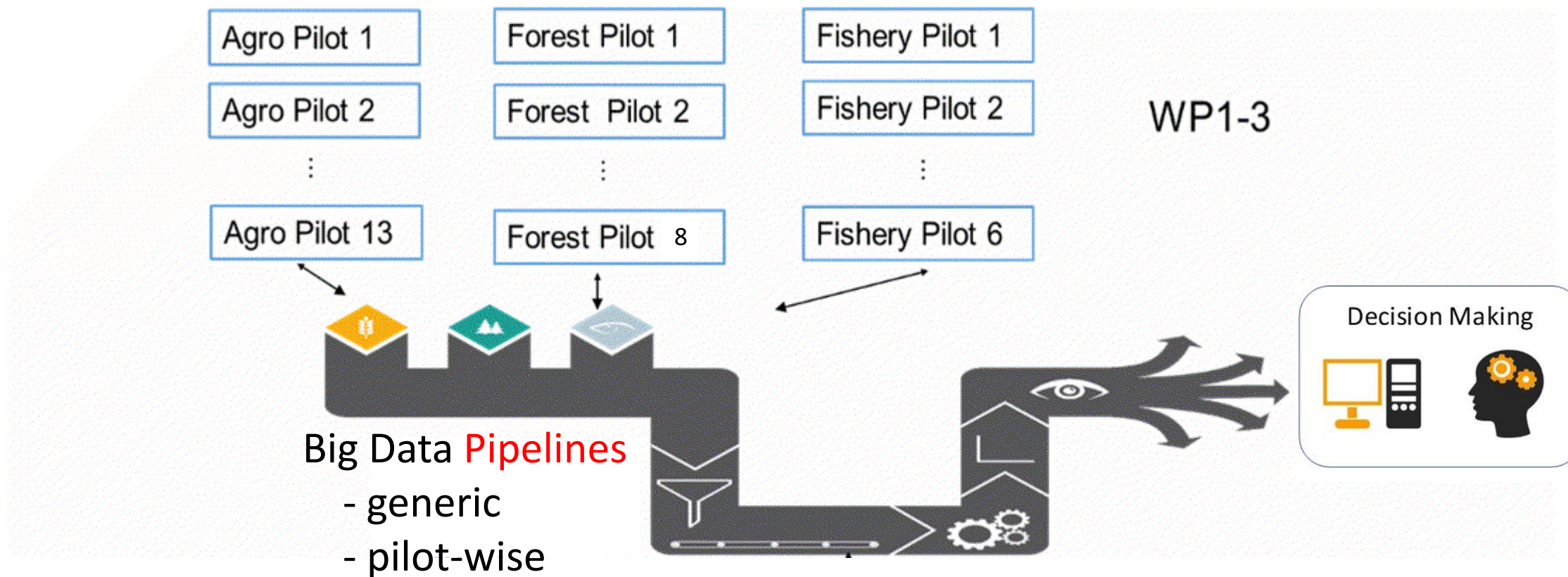


## Motivation

Population growth and urbanisation are increasing the demand for natural resources, which is putting a strain on the Earth's carrying capacity.

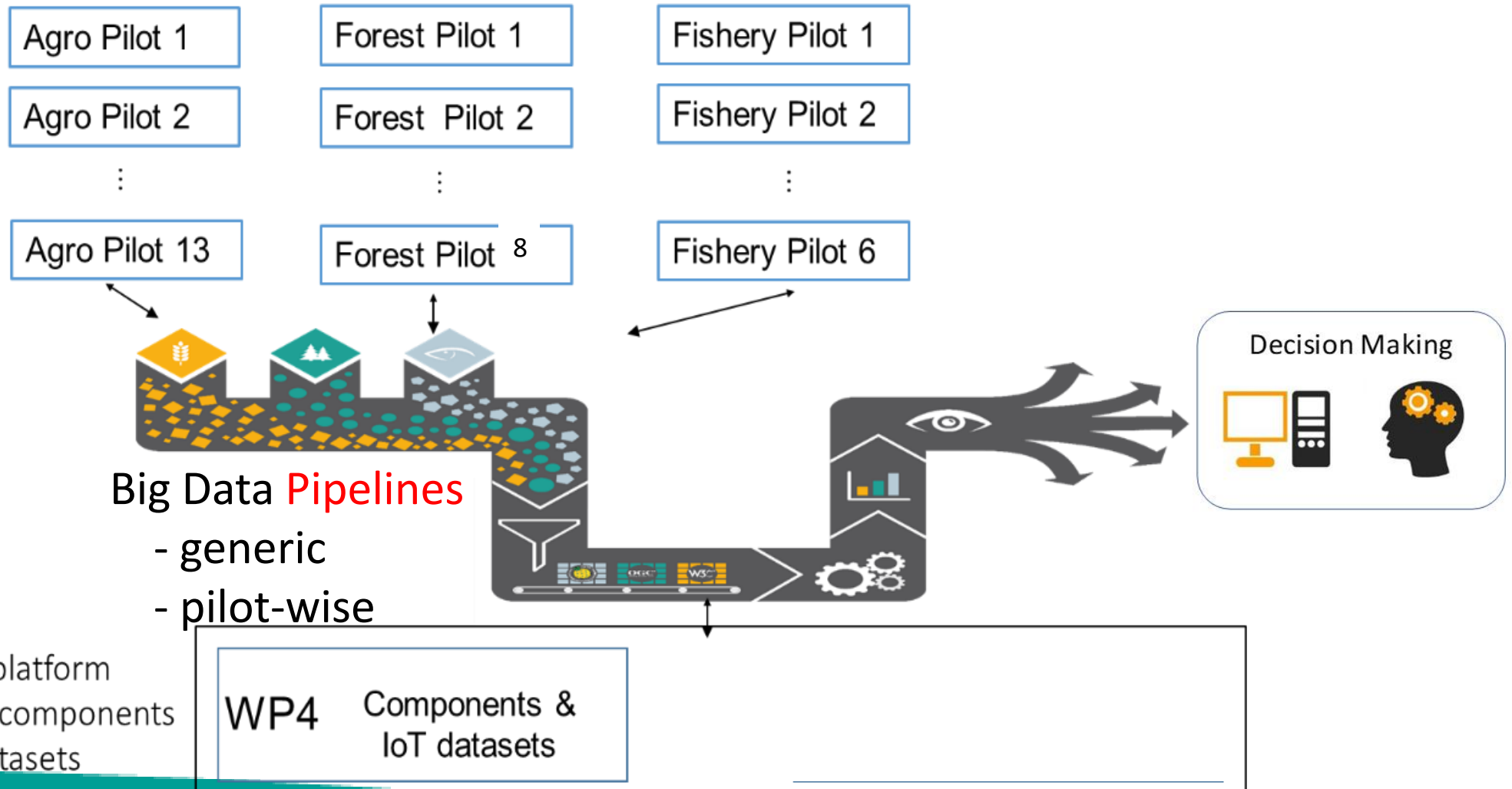
We aim to develop new sustainable ways to use forest, farm and fishery resources and to communicate real-time information to decision-makers and producers.

# DataBio platform serves the 27 pilots

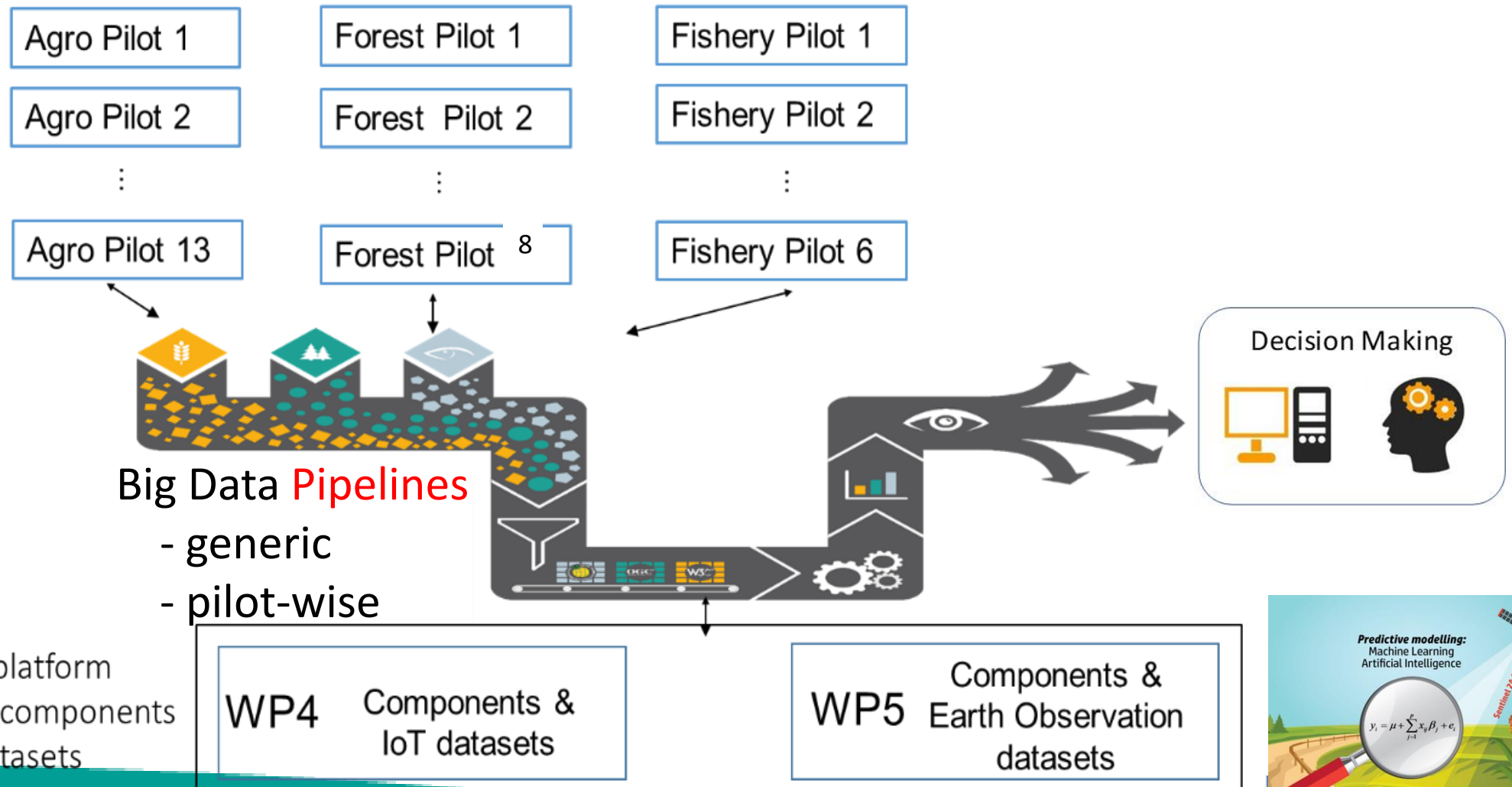




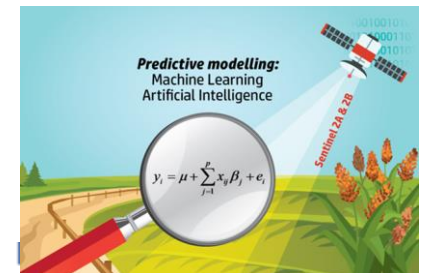
# DataBio platform serves the 27 pilots



# DataBio platform serves the 27 pilots



DataBio platform  
with big data components  
and datasets



# Pipelines vs. Services

## Pipeline

- A chain of real-time processing components:
  - Acquisition/Collection,
  - Preparation
  - Analytics
  - Visualisation, User interaction
- Clear interfaces between components and to outside
- A "*white box*" showing internal wiring for developers

## Service

- Provides usability to end users
- No display of internal wirings of components
- Accessed through API:s (web services, remote calls)
- Activated remotely through database queries ("end points") and executed in the cloud.
- Represents a "*black box*".

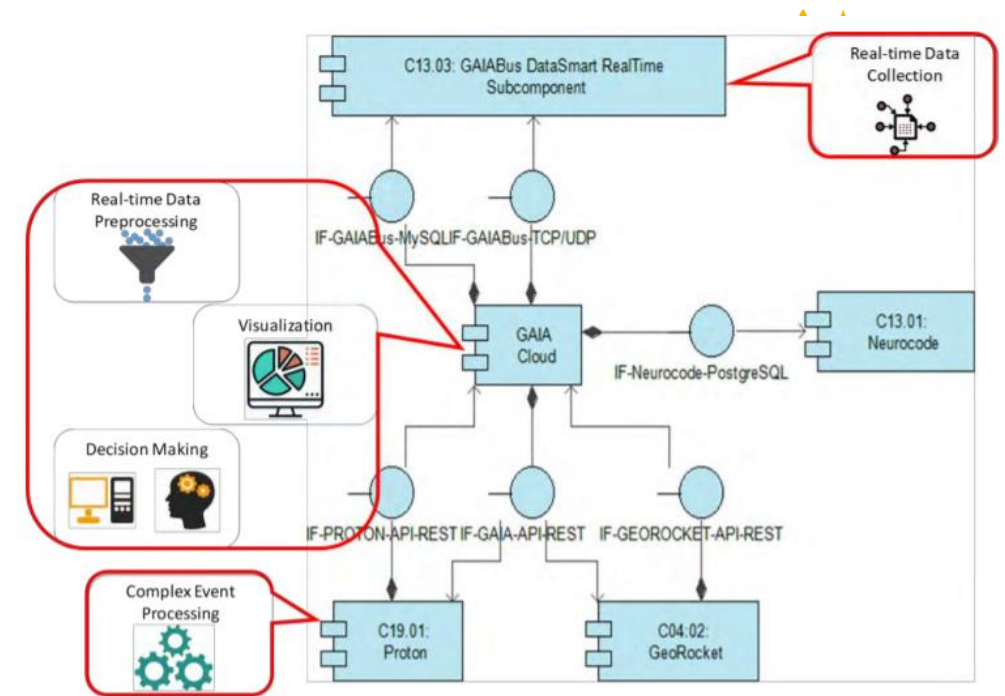
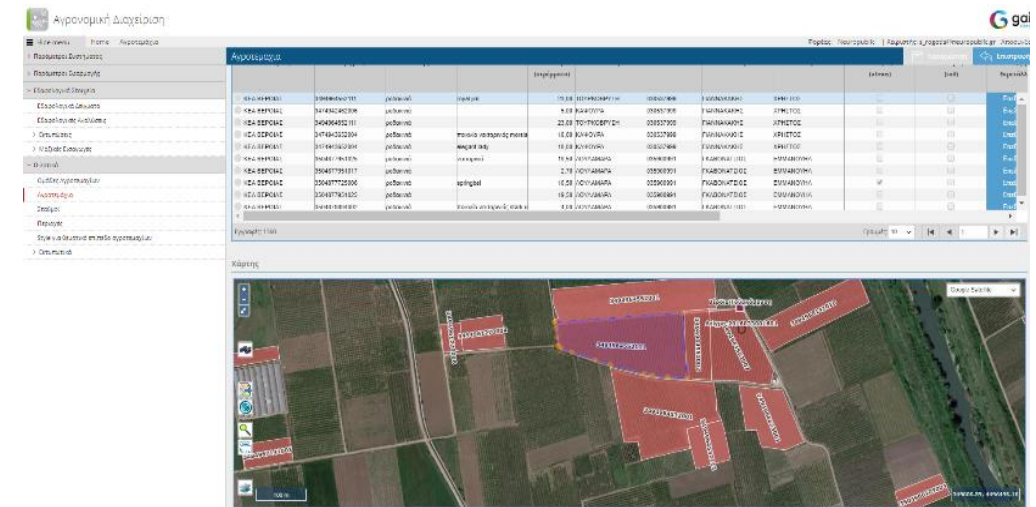
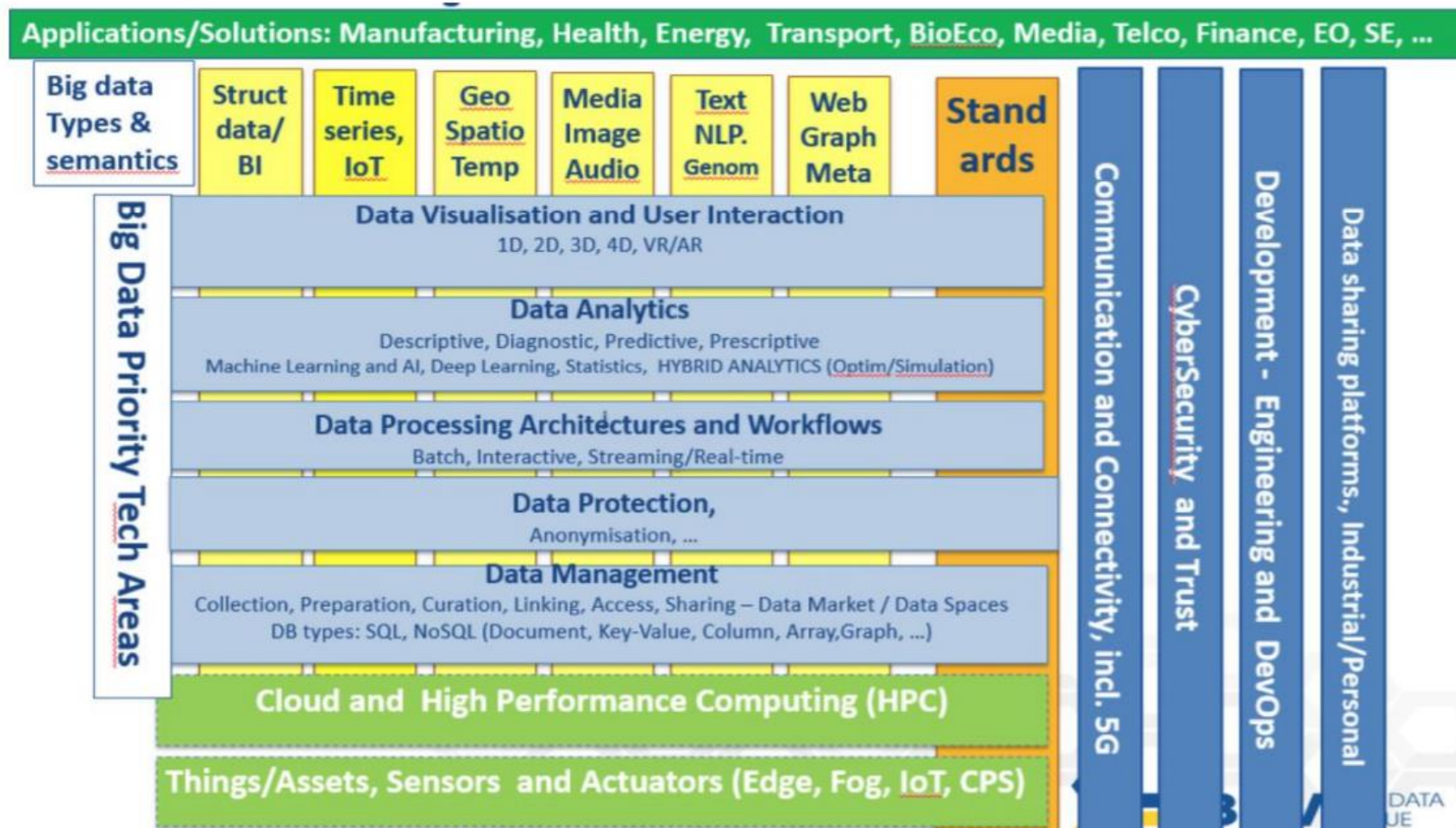


Figure 17: Mapping of generic components into pilot A1.1 component view

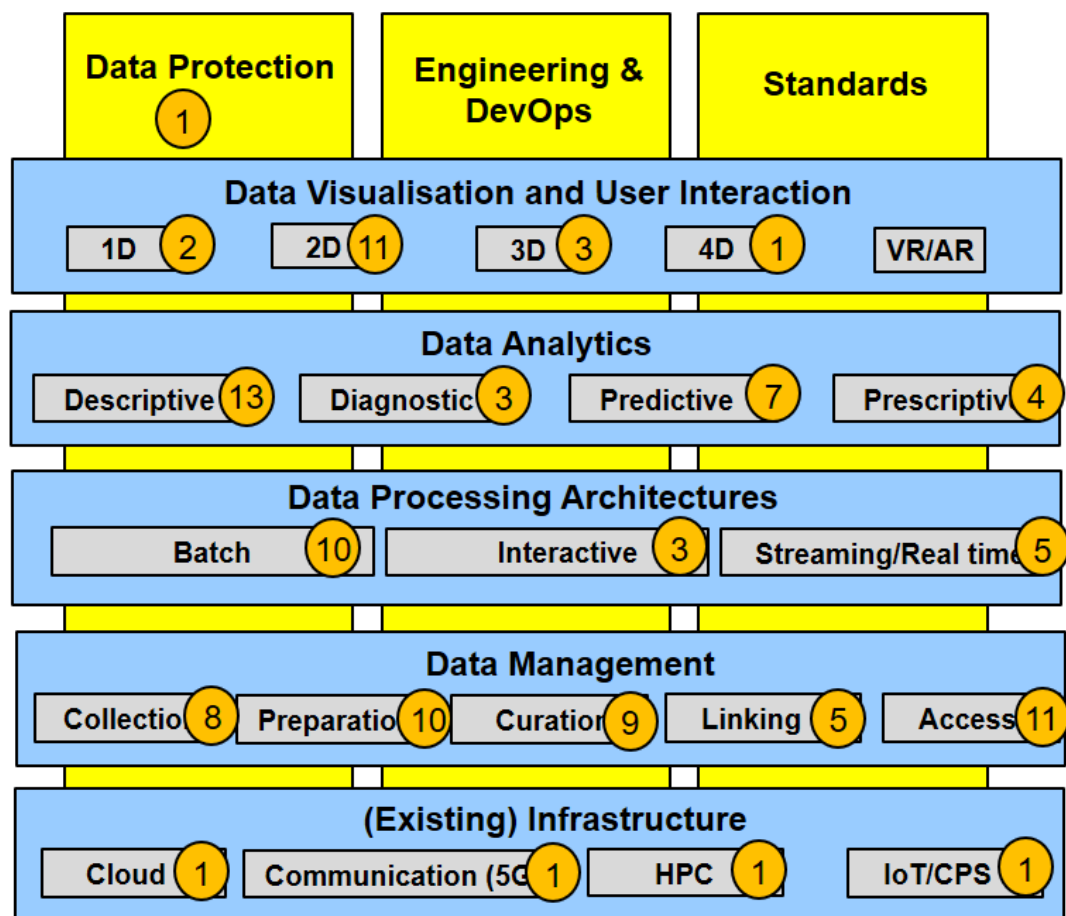




# We used the BDVA Reference Model



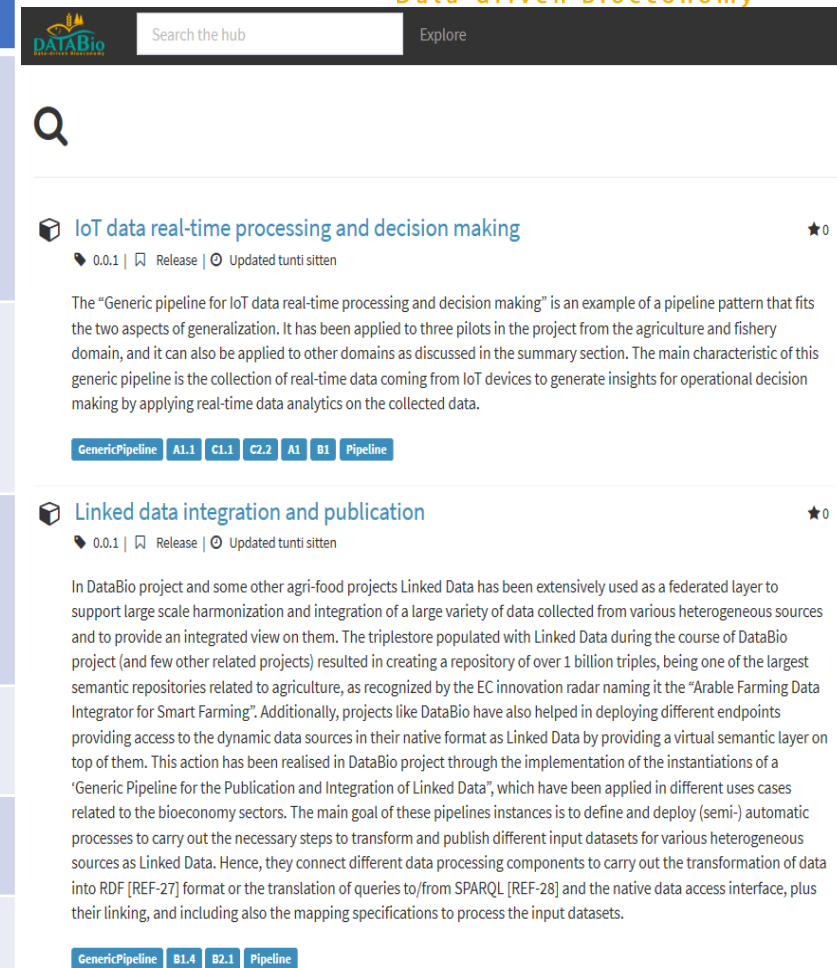
# Platform development in DataBio in numbers



- 62 components from 28 partners in two trial rounds
- 1-6 components per pilot (average 2)
- 14 new user interfaces
- 59 new APIs
- +2,7 in Technology Readiness Level (1-9)

# Generic pipelines in DataBio

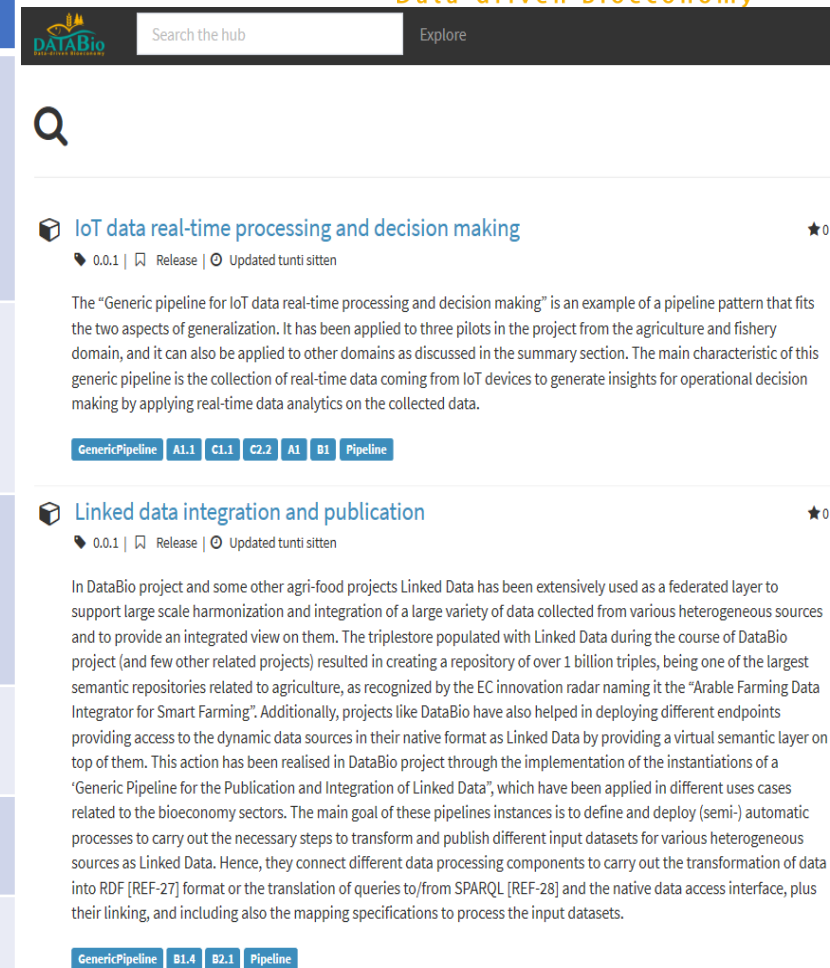
Generic pipeline	Applied in pilots
Earth Observation and Geospatial data processing	<b>A1.1</b> Precision agri <b>B1.2</b> Precision agri , <b>C1.1</b> Insurance, <b>C2.2</b> CAP Support Fishery: <b>A1</b> Tuna fishing operations , <b>B1</b> Tuna fishing planning,
IoT data real-time processing and decision making	Agri: <b>A1.1</b> Precision agri, <b>B1.1</b> Cereals & biomass crops Fishery: <b>A1</b>
Linked Data Integration and publication	Agri: <b>B1.4</b> Precision agri <b>B2.1</b> Machinery mgmt, <b>Other</b> open farm and geospatial datasets Fishery: <b>Virtual pilot</b>
Privacy-aware analytics	Fishery (Norway): <b>Virtual pilot</b>
Genomics	Agri: <b>A2.1</b> Greenhouse , <b>B1.3</b> Biomass sorghum
Forestry Data management	Forestry (Finland): <b>2.2.1</b> Data Sharing , <b>2.2.2</b> Monitoring, <b>2.4.2</b> Shared Multiuser environment
Fisheries decision support in catch planning	Fishery: <b>B2, C1, C2, Virtual pilot</b>



The screenshot shows the DataBio website interface. At the top, there is a search bar and a navigation menu. Below the search bar, there is a list of pipelines. The first pipeline listed is 'IoT data real-time processing and decision making'. It has a version number of 0.0.1, a 'Release' button, and a 'Updated tunti sitten' status. The description of this pipeline states: 'The "Generic pipeline for IoT data real-time processing and decision making" is an example of a pipeline pattern that fits the two aspects of generalization. It has been applied to three pilots in the project from the agriculture and fishery domain, and it can also be applied to other domains as discussed in the summary section. The main characteristic of this generic pipeline is the collection of real-time data coming from IoT devices to generate insights for operational decision making by applying real-time data analytics on the collected data.' Below the description, there is a tag 'GenericPipeline' and a list of categories: A1.1, C1.1, C2.2, A1, B1, Pipeline. The second pipeline listed is 'Linked data integration and publication'. It also has a version number of 0.0.1, a 'Release' button, and a 'Updated tunti sitten' status. The description of this pipeline states: 'In DataBio project and some other agri-food projects Linked Data has been extensively used as a federated layer to support large scale harmonization and integration of a large variety of data collected from various heterogeneous sources and to provide an integrated view on them. The triplestore populated with Linked Data during the course of DataBio project (and few other related projects) resulted in creating a repository of over 1 billion triples, being one of the largest semantic repositories related to agriculture, as recognized by the EC innovation radar naming it the "Arable Farming Data Integrator for Smart Farming". Additionally, projects like DataBio have also helped in deploying different endpoints providing access to the dynamic data sources in their native format as Linked Data by providing a virtual semantic layer on top of them. This action has been realised in DataBio project through the implementation of the instantiations of a "Generic Pipeline for the Publication and Integration of Linked Data", which have been applied in different uses cases related to the bioeconomy sectors. The main goal of these pipelines instances is to define and deploy (semi-) automatic processes to carry out the necessary steps to transform and publish different input datasets for various heterogeneous sources as Linked Data. Hence, they connect different data processing components to carry out the transformation of data into RDF [REF-27] format or the translation of queries to/from SPARQL [REF-28] and the native data access interface, plus their linking, and including also the mapping specifications to process the input datasets.' Below the description, there is a tag 'GenericPipeline' and a list of categories: B1.4, B2.1, Pipeline.

# Generic pipelines in DataBio

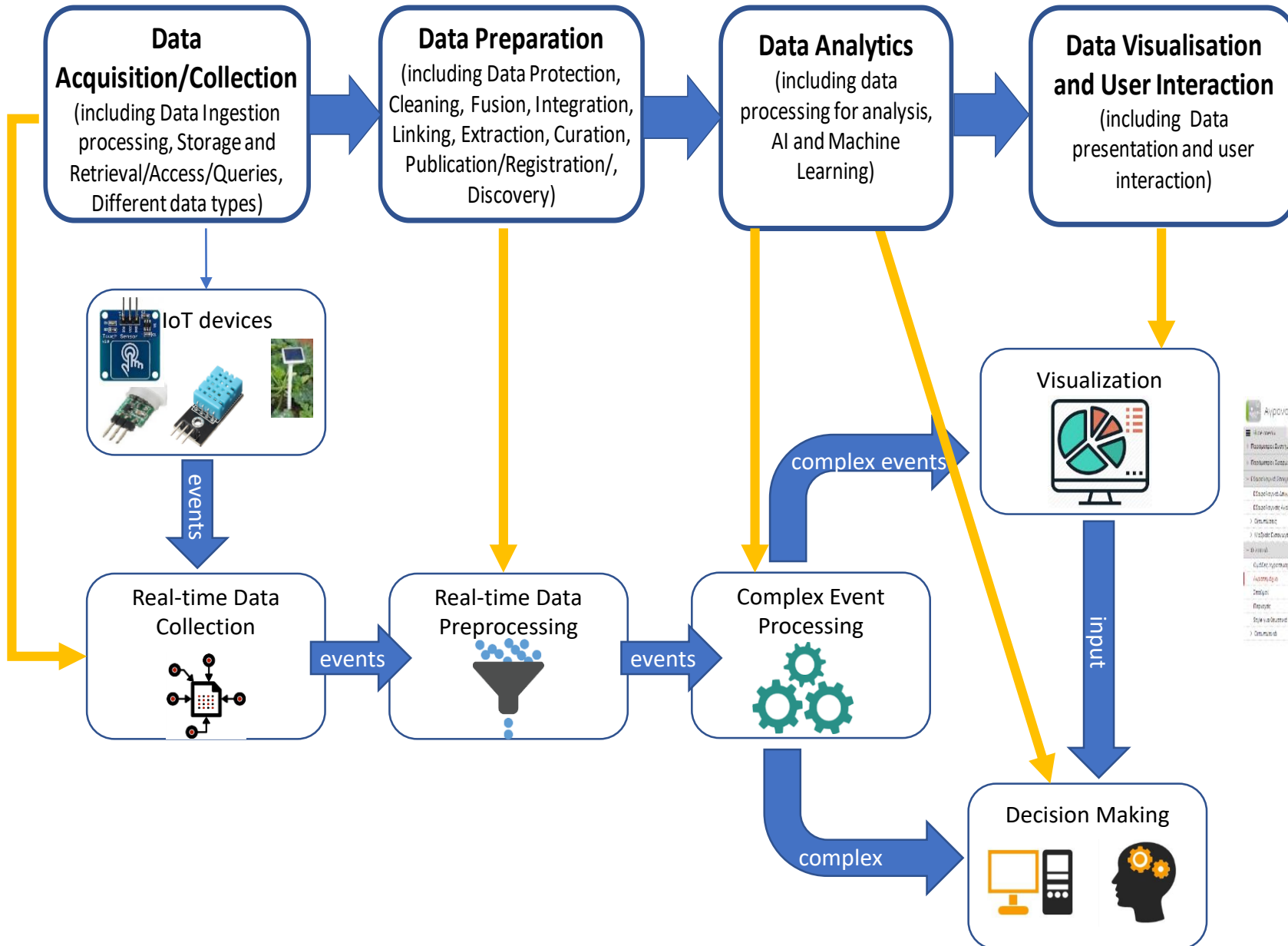
Generic pipeline	Applied in pilots
Earth Observation and Geospatial data processing	<b>A1.1</b> Precision agri <b>B1.2</b> Precision agri , <b>C1.1</b> Insurance, <b>C2.2</b> CAP Support Fishery: <b>A1</b> Tuna fishing operations , <b>B1</b> Tuna fishing planning,
IoT data real-time processing and decision making	Agri: <b>A1.1</b> Precision agri, <b>B1.1</b> Cereals & biomass crops Fishery: <b>A1</b>
Linked Data Integration and publication	Agri: <b>B1.4</b> Precision agri <b>B2.1</b> Machinery mgmt, <b>Other</b> open farm and geospatial datasets Fishery: <b>Virtual pilot</b>
Privacy-aware analytics	Fishery (Norway): <b>Virtual pilot</b>
Genomics	Agri: <b>A2.1</b> Greenhouse , <b>B1.3</b> Biomass sorghum
Forestry Data management	Forestry (Finland): <b>2.2.1</b> Data Sharing , <b>2.2.2</b> Monitoring, <b>2.4.2</b> Shared Multiuser environment
Fisheries decision support in catch planning	Fishery: <b>B2, C1, C2, Virtual pilot</b>



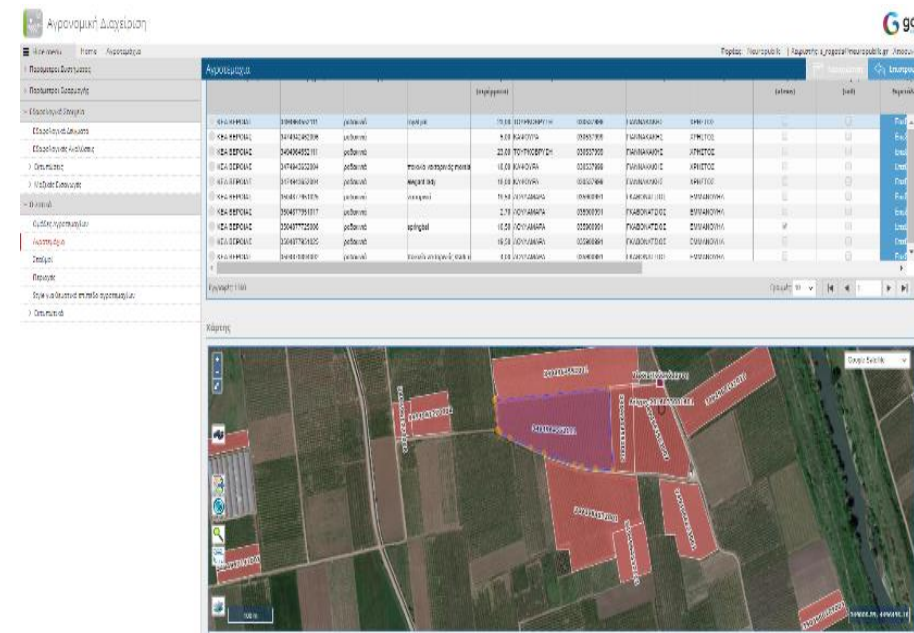
The screenshot shows the DataBio web portal interface. At the top, there is a search bar with the text 'Search the hub' and a button labeled 'Explore'. Below the search bar, there is a search icon (magnifying glass). The main content area displays two search results. The first result is titled 'IoT data real-time processing and decision making' and includes a version number '0.0.1', a 'Release' button, and a 'Updated tunti sitten' (Updated a few minutes ago) status. The second result is titled 'Linked data integration and publication' and also includes a version number '0.0.1', a 'Release' button, and a 'Updated tunti sitten' status. Both results have a star icon and a '0' next to it, indicating no favorites or ratings. The 'IoT data real-time processing and decision making' result has a detailed description: 'The "Generic pipeline for IoT data real-time processing and decision making" is an example of a pipeline pattern that fits the two aspects of generalization. It has been applied to three pilots in the project from the agriculture and fishery domain, and it can also be applied to other domains as discussed in the summary section. The main characteristic of this generic pipeline is the collection of real-time data coming from IoT devices to generate insights for operational decision making by applying real-time data analytics on the collected data.' Below the description, there are tags: 'GenericPipeline', 'A1.1', 'C1.1', 'C2.2', 'A1', 'B1', and 'Pipeline'. The 'Linked data integration and publication' result has a detailed description: 'In DataBio project and some other agri-food projects Linked Data has been extensively used as a federated layer to support large scale harmonization and integration of a large variety of data collected from various heterogeneous sources and to provide an integrated view on them. The triplestore populated with Linked Data during the course of DataBio project (and few other related projects) resulted in creating a repository of over 1 billion triples, being one of the largest semantic repositories related to agriculture, as recognized by the EC innovation radar naming it the "Arable Farming Data Integrator for Smart Farming". Additionally, projects like DataBio have also helped in deploying different endpoints providing access to the dynamic data sources in their native format as Linked Data by providing a virtual semantic layer on top of them. This action has been realised in DataBio project through the implementation of the instantiations of a "Generic Pipeline for the Publication and Integration of Linked Data", which have been applied in different uses cases related to the bioeconomy sectors. The main goal of these pipelines instances is to define and deploy (semi-) automatic processes to carry out the necessary steps to transform and publish different input datasets for various heterogeneous sources as Linked Data. Hence, they connect different data processing components to carry out the transformation of data into RDF [REF-27] format or the translation of queries to/from SPARQL [REF-28] and the native data access interface, plus their linking, and including also the mapping specifications to process the input datasets.' Below the description, there are tags: 'GenericPipeline', 'B1.4', 'B2.1', and 'Pipeline'.



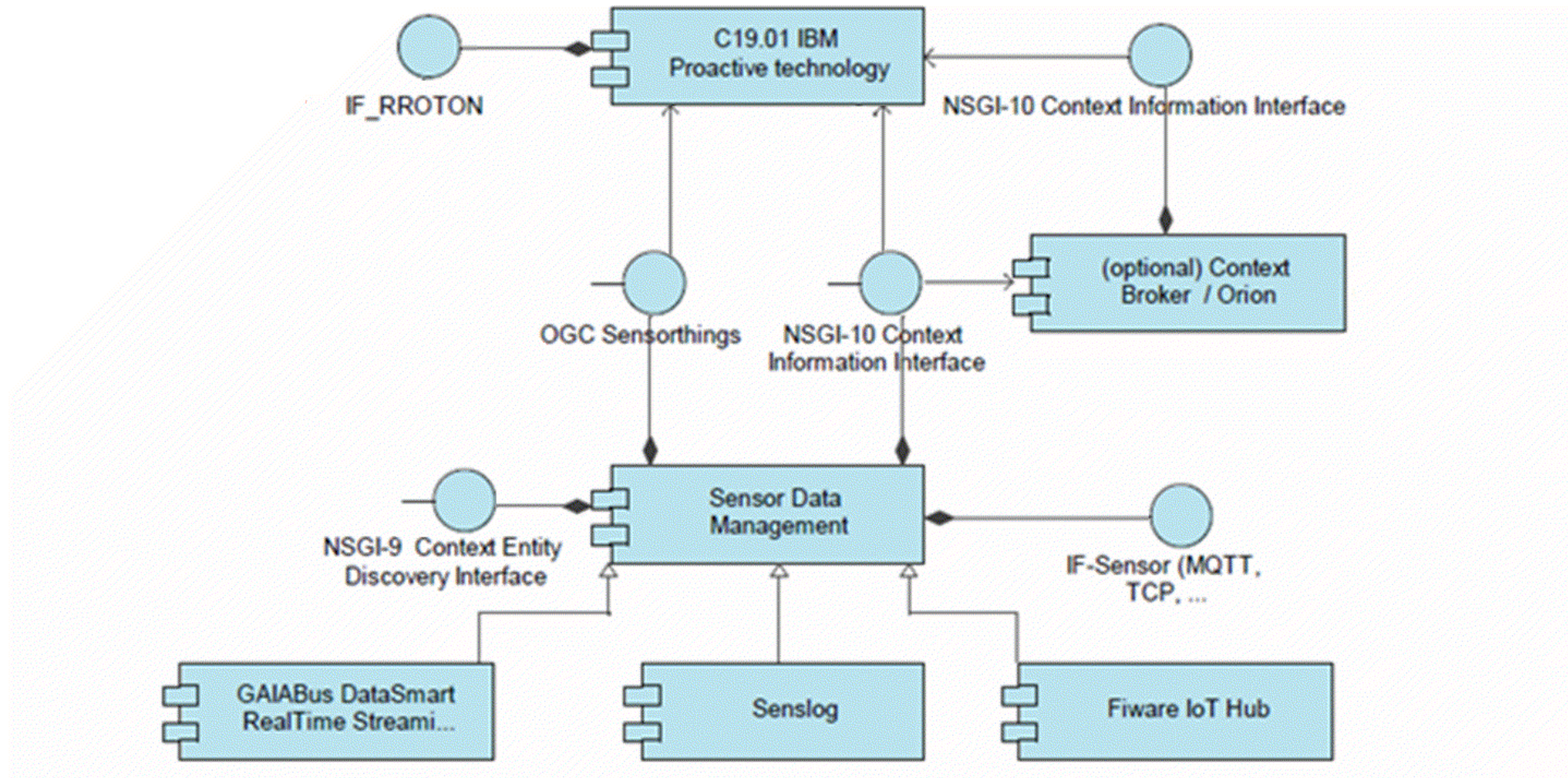
# Generic Pipeline – IoT data realtime processing



Instance: Prediction and real-time alerts of diseases and pests outbreaks in crops



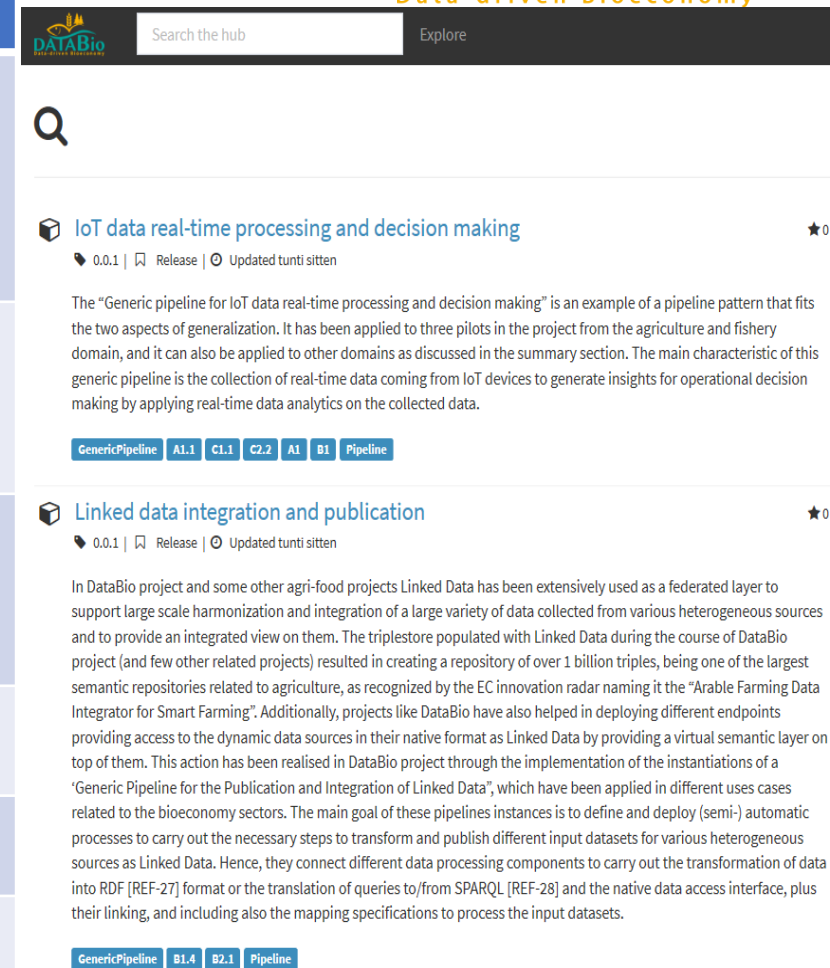
# Crop monitoring using "IoT Generic Pipeline"





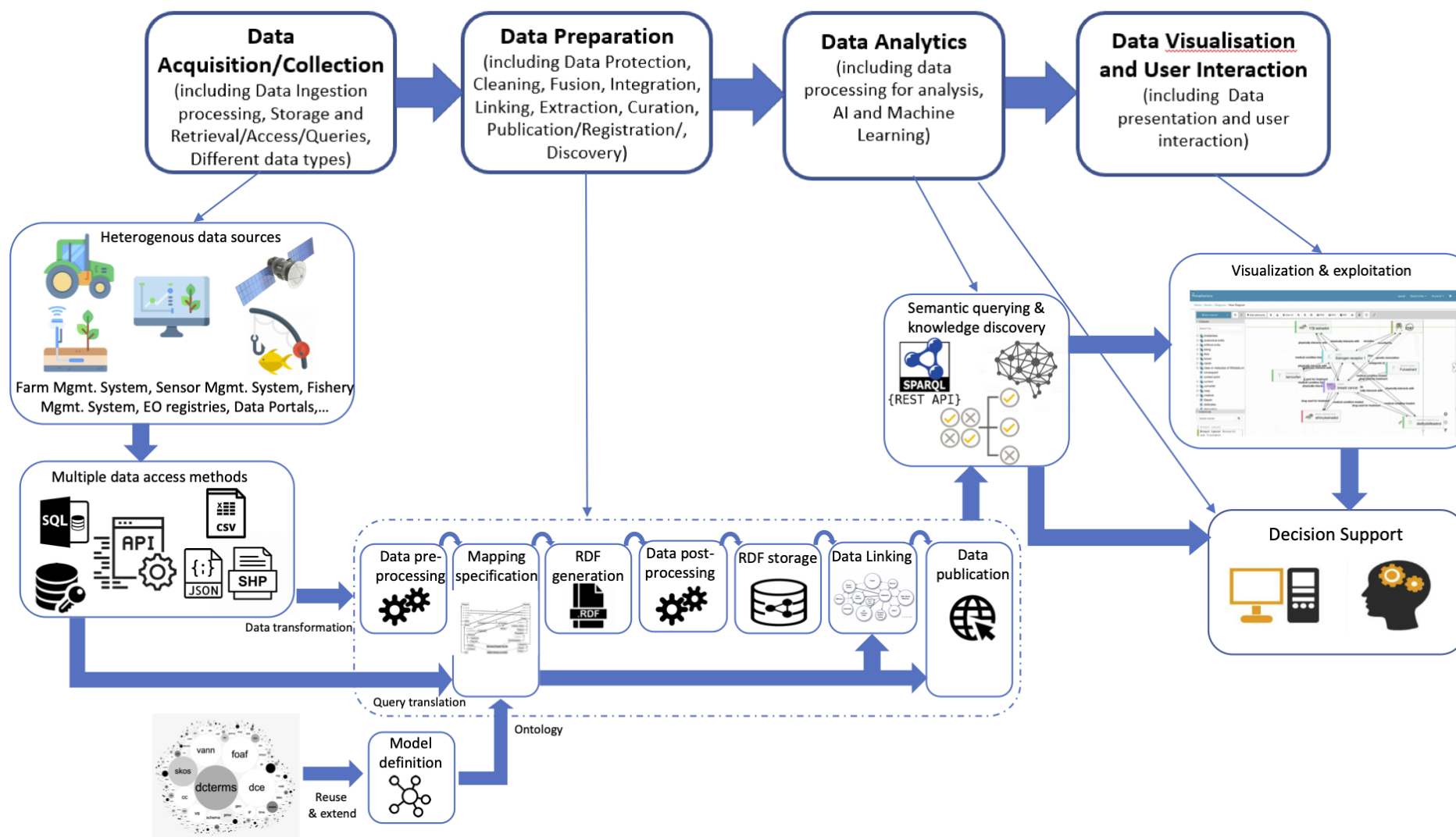
# Generic pipelines in DataBio

Generic pipeline	Applied in pilots
Earth Observation and Geospatial data processing	<b>A1.1</b> Precision agri <b>B1.2</b> Precision agri , <b>C1.1</b> Insurance, <b>C2.2</b> CAP Support Fishery: <b>A1</b> Tuna fishing operations , <b>B1</b> Tuna fishing planning,
IoT data real-time processing and decision making	Agri: <b>A1.1</b> Precision agri, <b>B1.1</b> Cereals & biomass crops Fishery: <b>A1</b>
Linked Data Integration and publication	Agri: <b>B1.4</b> Precision agri <b>B2.1</b> Machinery mgmt, <b>Other</b> open farm and geospatial datasets Fishery: <b>Virtual pilot</b>
Privacy-aware analytics	Fishery (Norway): <b>Virtual pilot</b>
Genomics	Agri: <b>A2.1</b> Greenhouse , <b>B1.3</b> Biomass sorghum
Forestry Data management	Forestry (Finland): <b>2.2.1</b> Data Sharing , <b>2.2.2</b> Monitoring, <b>2.4.2</b> Shared Multiuser environment
Fisheries decision support in catch planning	Fishery: <b>B2, C1, C2, Virtual pilot</b>

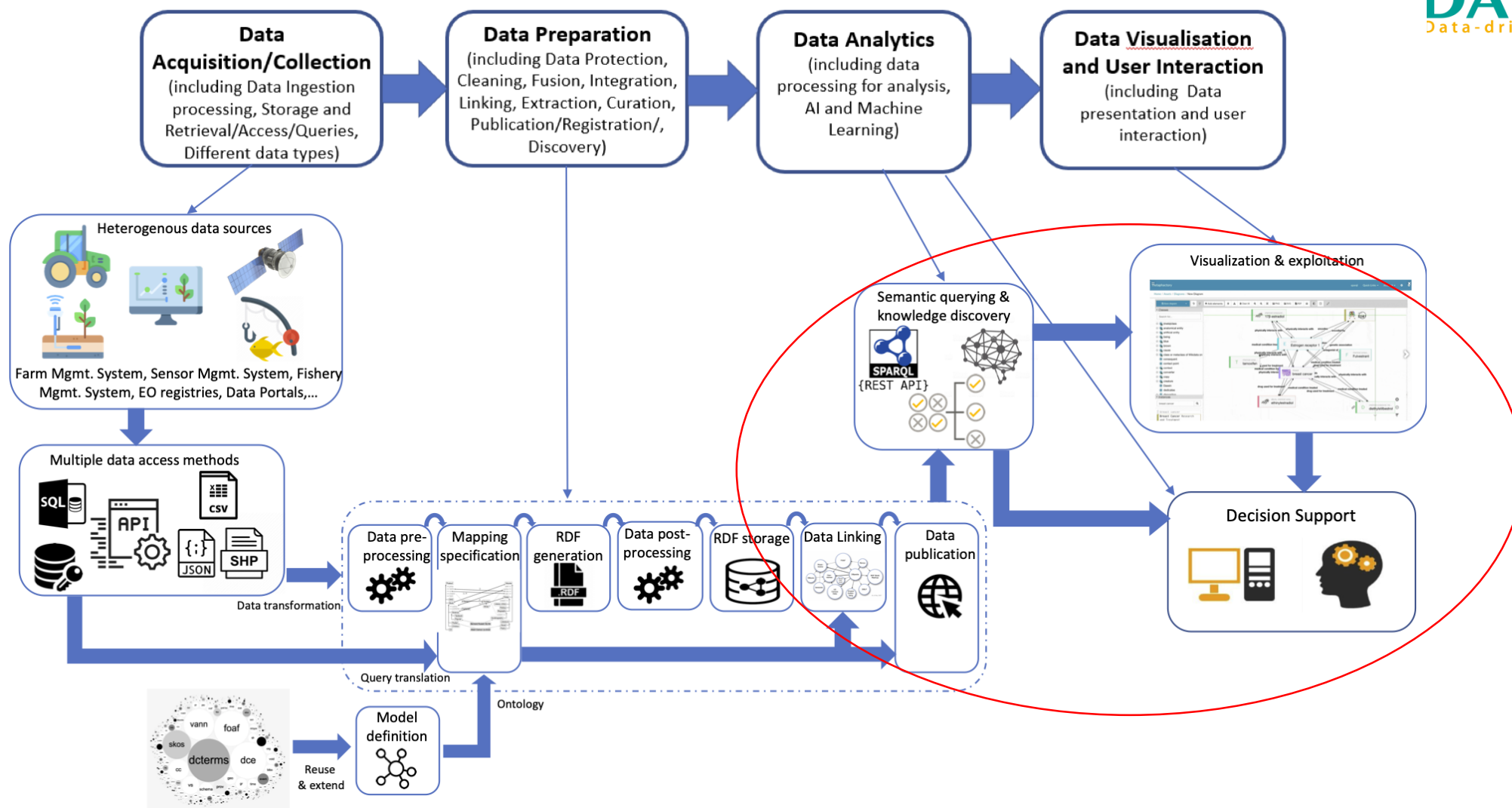


The screenshot shows the DataBio web portal interface. At the top, there is a search bar with the text 'Search the hub' and a 'Explore' button. Below the search bar, there is a search result for 'IoT data real-time processing and decision making'. The result includes a version number '0.0.1', a 'Release' button, and a 'Updated tunti sitten' status. The description states: 'The "Generic pipeline for IoT data real-time processing and decision making" is an example of a pipeline pattern that fits the two aspects of generalization. It has been applied to three pilots in the project from the agriculture and fishery domain, and it can also be applied to other domains as discussed in the summary section. The main characteristic of this generic pipeline is the collection of real-time data coming from IoT devices to generate insights for operational decision making by applying real-time data analytics on the collected data.' Below the description, there is a tag 'GenericPipeline' and a list of categories: 'A1.1', 'C1.1', 'C2.2', 'A1', 'B1', and 'Pipeline'. Another search result is visible below, for 'Linked data integration and publication', which also includes a version number, release button, and update status. Its description mentions: 'In DataBio project and some other agri-food projects Linked Data has been extensively used as a federated layer to support large scale harmonization and integration of a large variety of data collected from various heterogeneous sources and to provide an integrated view on them. The triplestore populated with Linked Data during the course of DataBio project (and few other related projects) resulted in creating a repository of over 1 billion triples, being one of the largest semantic repositories related to agriculture, as recognized by the EC innovation radar naming it the "Arable Farming Data Integrator for Smart Farming". Additionally, projects like DataBio have also helped in deploying different endpoints providing access to the dynamic data sources in their native format as Linked Data by providing a virtual semantic layer on top of them. This action has been realised in DataBio project through the implementation of the instantiations of a "Generic Pipeline for the Publication and Integration of Linked Data", which have been applied in different uses cases related to the bioeconomy sectors. The main goal of these pipelines instances is to define and deploy (semi-) automatic processes to carry out the necessary steps to transform and publish different input datasets for various heterogeneous sources as Linked Data. Hence, they connect different data processing components to carry out the transformation of data into RDF [REF-27] format or the translation of queries to/from SPARQL [REF-28] and the native data access interface, plus their linking, and including also the mapping specifications to process the input datasets.'

# Linked Data Integration and Publication Pipeline



# Case: Link discovery from knowledge graphs



# Case: Link discovery from knowledge graphs

We applied discovery of RDF spatial links based on topology (Geo-L):

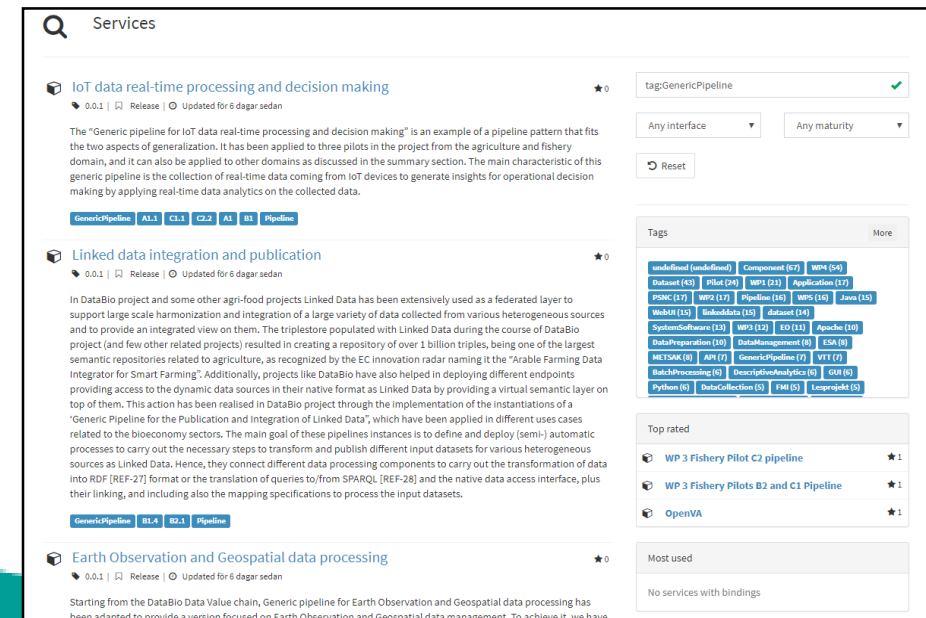
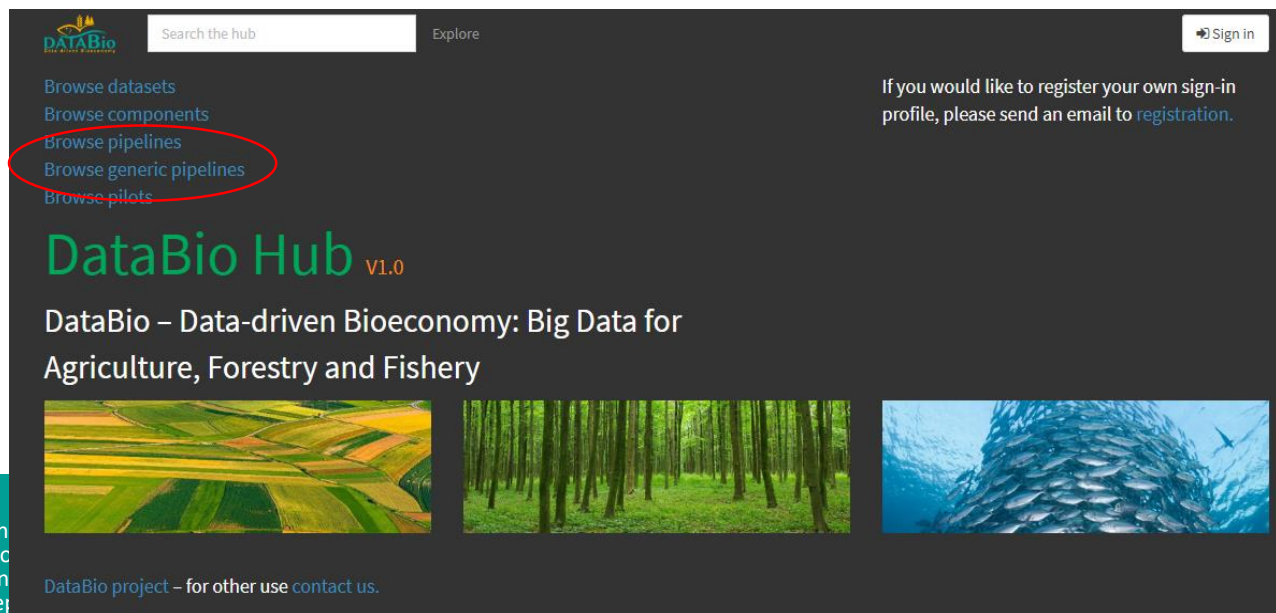
- Identifying fields from *Czech LPIS* data with specific soil type, from *Czech open data*
- Identifying all fields in a specific region which grow the same type of crops like the one grown in a specific field over a given period of time
- Identifying "risky" plots from *Czech LPIS* data which intersect with *buffer zones* around water bodies.



Figure: Risky overlap area between a crop field and a buffer zone of a water

# Managing project assets: components, pipelines, datasets and reports

- [DataBio Hub](https://DataBioHub.eu) (DataBioHub.eu) is central in the development platform
  - Provides a catalogue of public (and private) digital assets of DataBio
  - Links resources together (project reports, models, docker modules etc)
  - Describes currently 101 **components**, 65 **datasets**, **25 pipelines** (7 generic), 27 **pilots**
- Provides links to Deliverables, Interfaces (SPARQL endpoints, REST...)





- DataBio designed a *common* development platform for 27 pilots in agriculture, forestry and fishery
- Big Data *pipelines* were central in this platform as a tool for developers
- We showed that a few *generic* pipelines can be applied in numerous diverse applications
- *Instantiations* of the generic pipelines were used in the pilots
- All Dtabio results are available in the searchable and cross-linked *DataBio Hub*.

# Thank you for your attention!





# DEEPHEALTH

Deep-Learning and HPC to Boost Biomedical Applications for Health

## Pipelines for Medical Imaging Use Cases & Requirements for Benchmarking

Federico **Boelli**, Jon Ander **Gómez**, Costantino **Grana**, Roberto **Paredes**

Evaluation schemes for Big data and AI Performance  
of high Business impact

EBDVF 2020, November 3-5, 2020



*The project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 825111.*

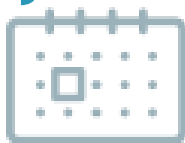


# About DeepHealth

## Aim & Goals

- Put HPC computing power at the service of biomedical applications with DL needs and apply DL techniques on large and complex image biomedical datasets to support new and more efficient ways of diagnosis, monitoring and treatment of diseases.
- Facilitate the daily work and increase the productivity of medical personnel and IT professionals in terms of image processing and the use and training of predictive models without the need of combining numerous tools.
- Offer a unified framework adapted to exploit underlying heterogeneous HPC and Cloud architectures supporting state-of-the-art and next-generation Deep Learning (AI) and Computer Vision algorithms to enhance European-based medical software platforms.

## Key facts



Duration: 36 months  
Starting date: Jan 2019



Budget 14.642.366 €  
EU funding 12.774.824 €



22 partners from 9 countries:  
Research centers, Health organizations,  
large industries and SMEs

### Research Organisations



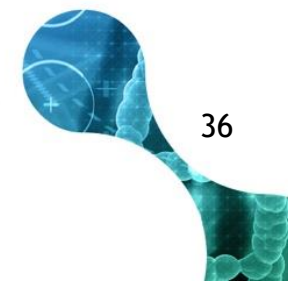
### Health Organisations



### Large Industries



### SMEs



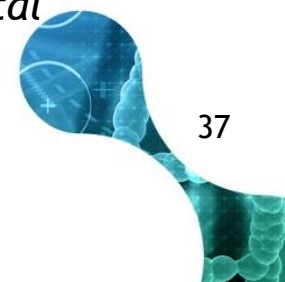
# Developments & Expected Results

- **The DeepHealth toolkit**

- Free and open-source software: 2 libraries + front-end.
  - *EDDLL: The European Distributed Deep Learning Library*
  - *ECVL: the European Computer Vision Library*
- Ready to run algorithms on Hybrid HPC + Cloud architectures with heterogeneous hardware (Distributed versions of the training algorithms)
- Ready to be integrated into end-user software platforms or applications



- **HPC infrastructure** for an efficient execution of the training algorithms which are computationally intensive by making use of heterogeneous hardware in a transparent way
- Seven enhanced **biomedical and AI software platforms** provided by EVERIS, PHILIPS, THALES, UNITO, WINGS, CRS4 and CEA that integrate the DeepHealth libraries to improve their potential
- Proposal for a structure for anonymised and pseudonymised data lakes
- **Validation** in 14 use cases (*neurological diseases, tumor detection and early cancer prediction, digital pathology and automated image annotation*).

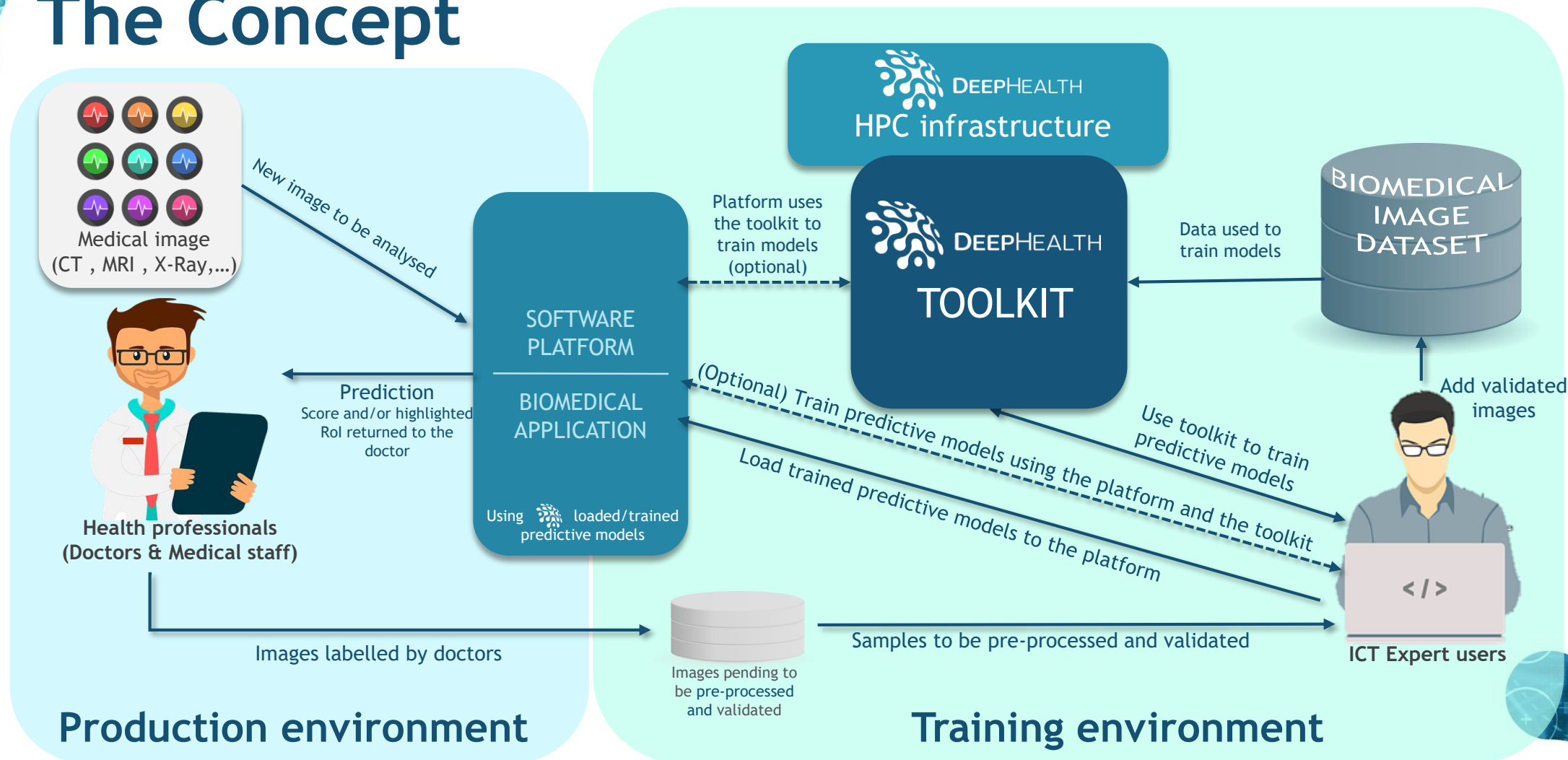






DEEPHEALTH

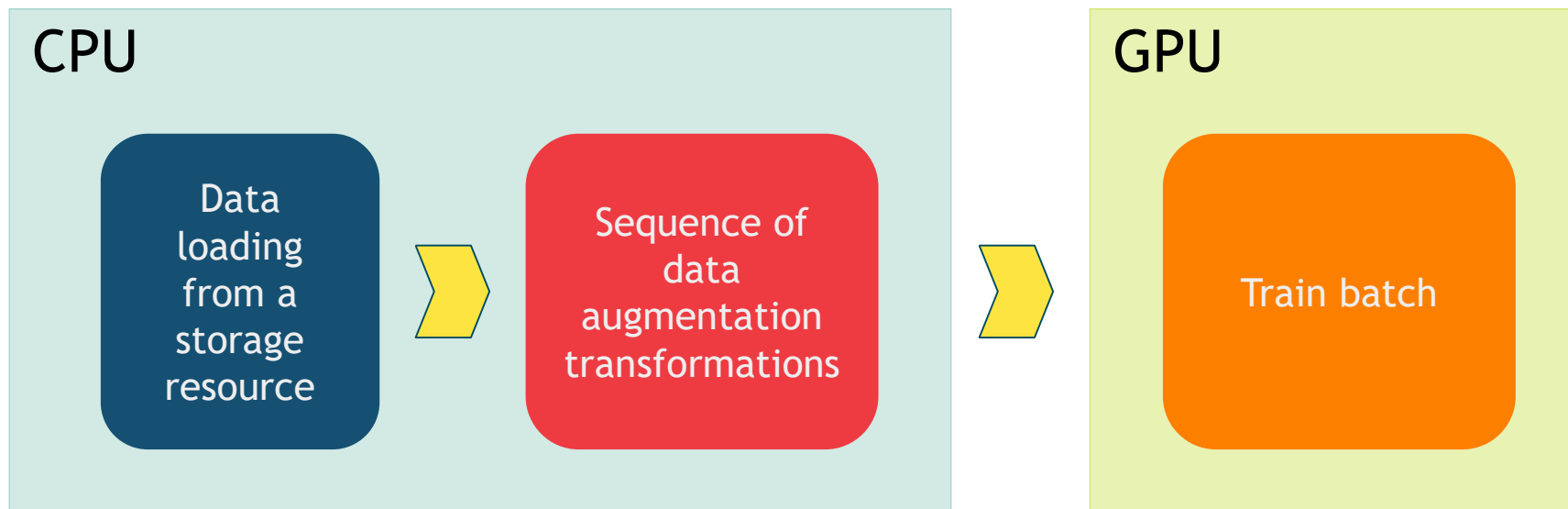
# The Concept



# Pipeline Definition

What do we understand by pipeline in this context?

A **pipeline** is a set of operations **sequentially applied** to a data block (subset of samples)

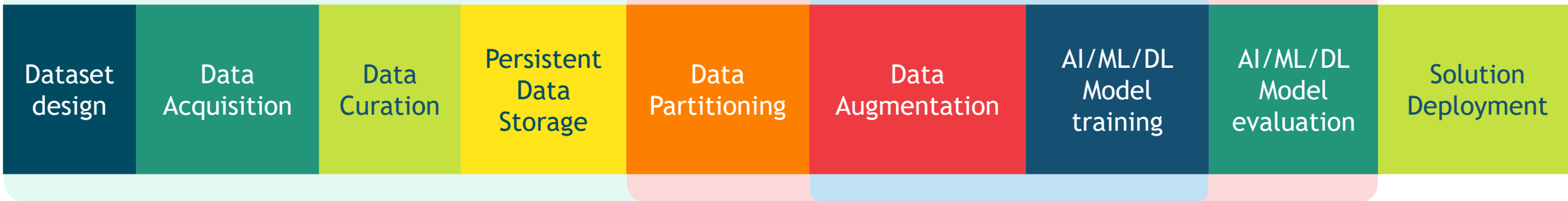


# Data & Model Pipelines

## Data Pipeline

## AI/ML/DL Model Pipeline

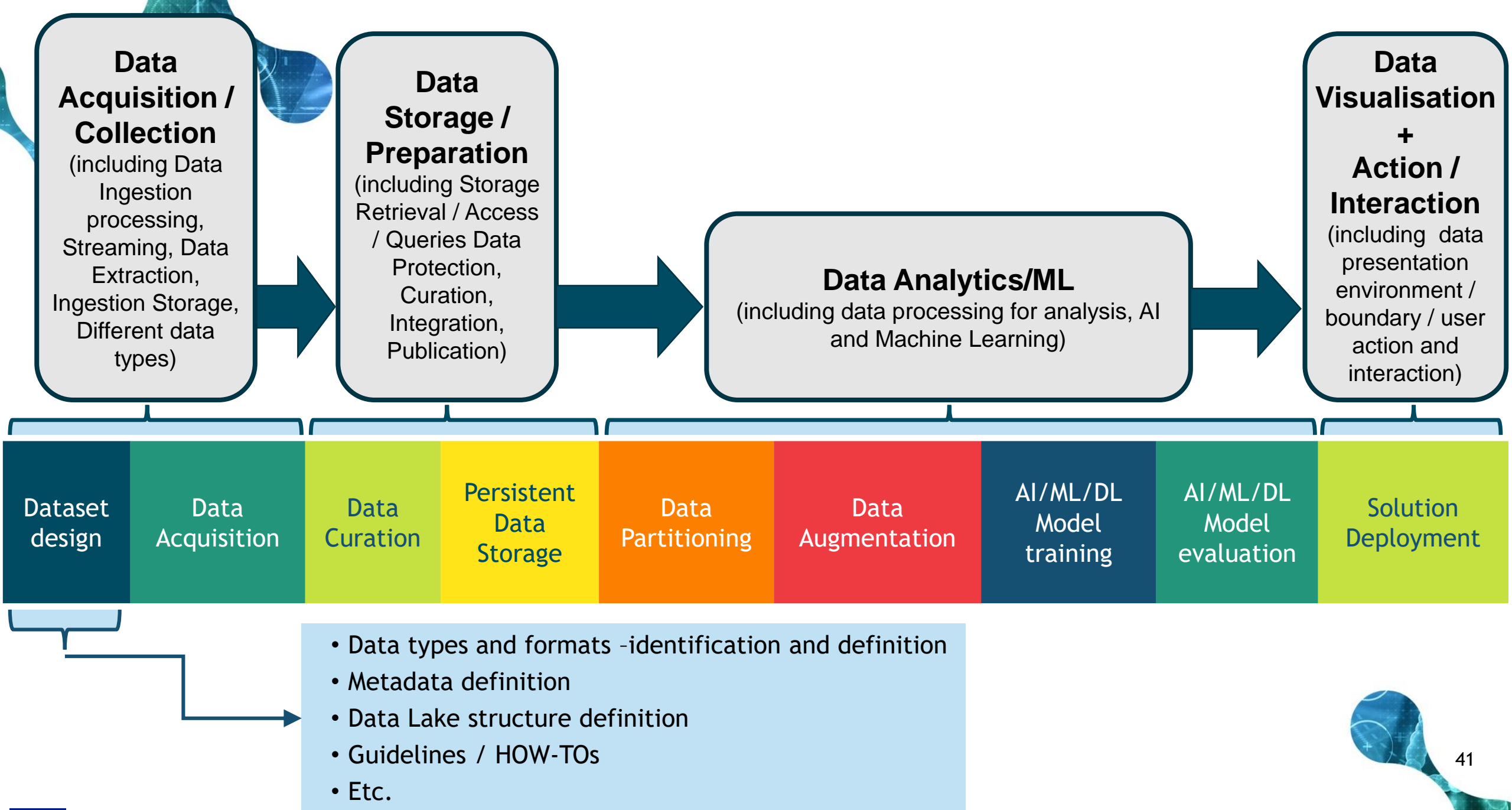
## Training Pipeline

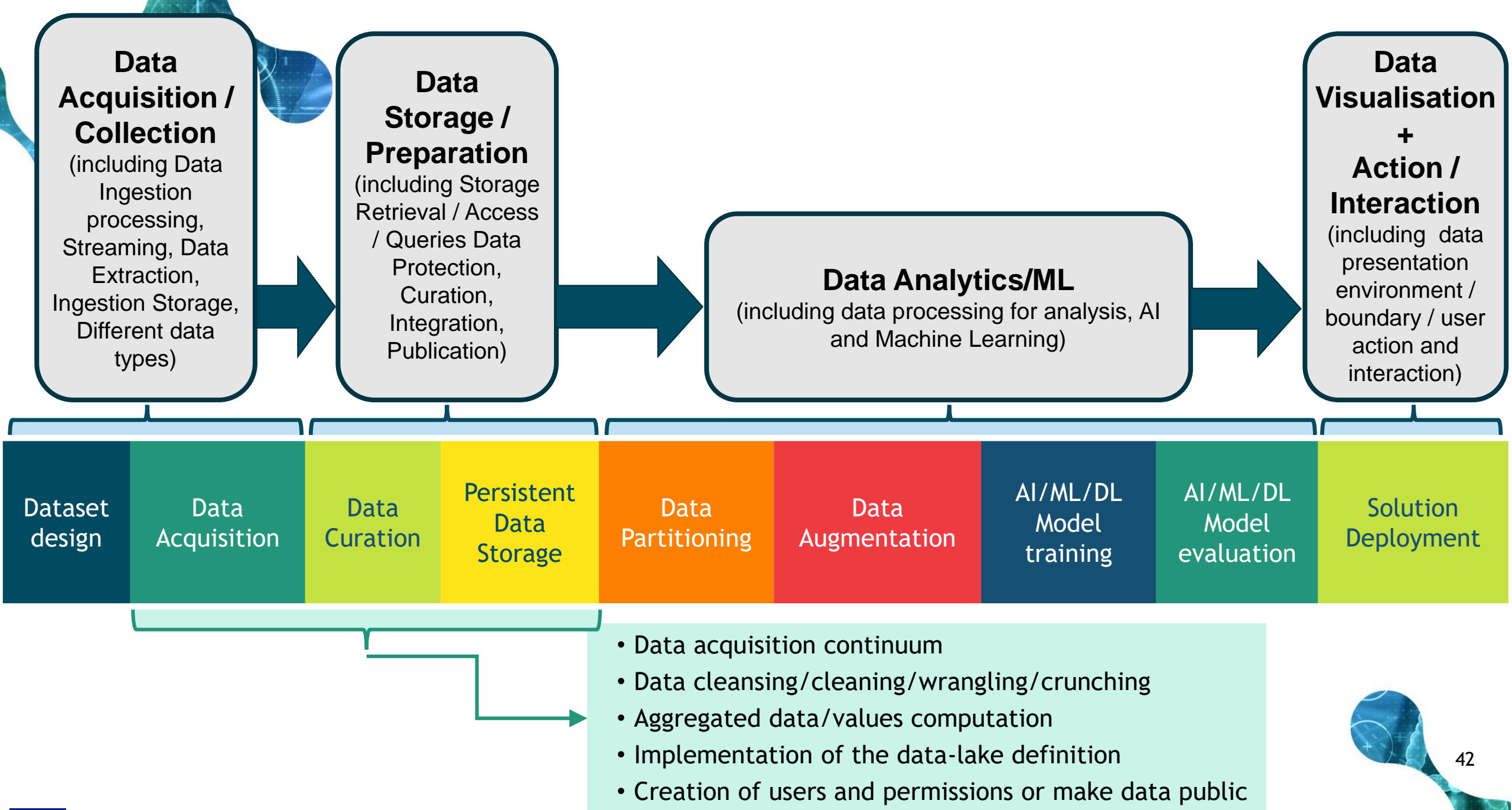


*Training Pipeline is at the core of the Model Pipeline which in turn is considered part of the Data Pipeline*

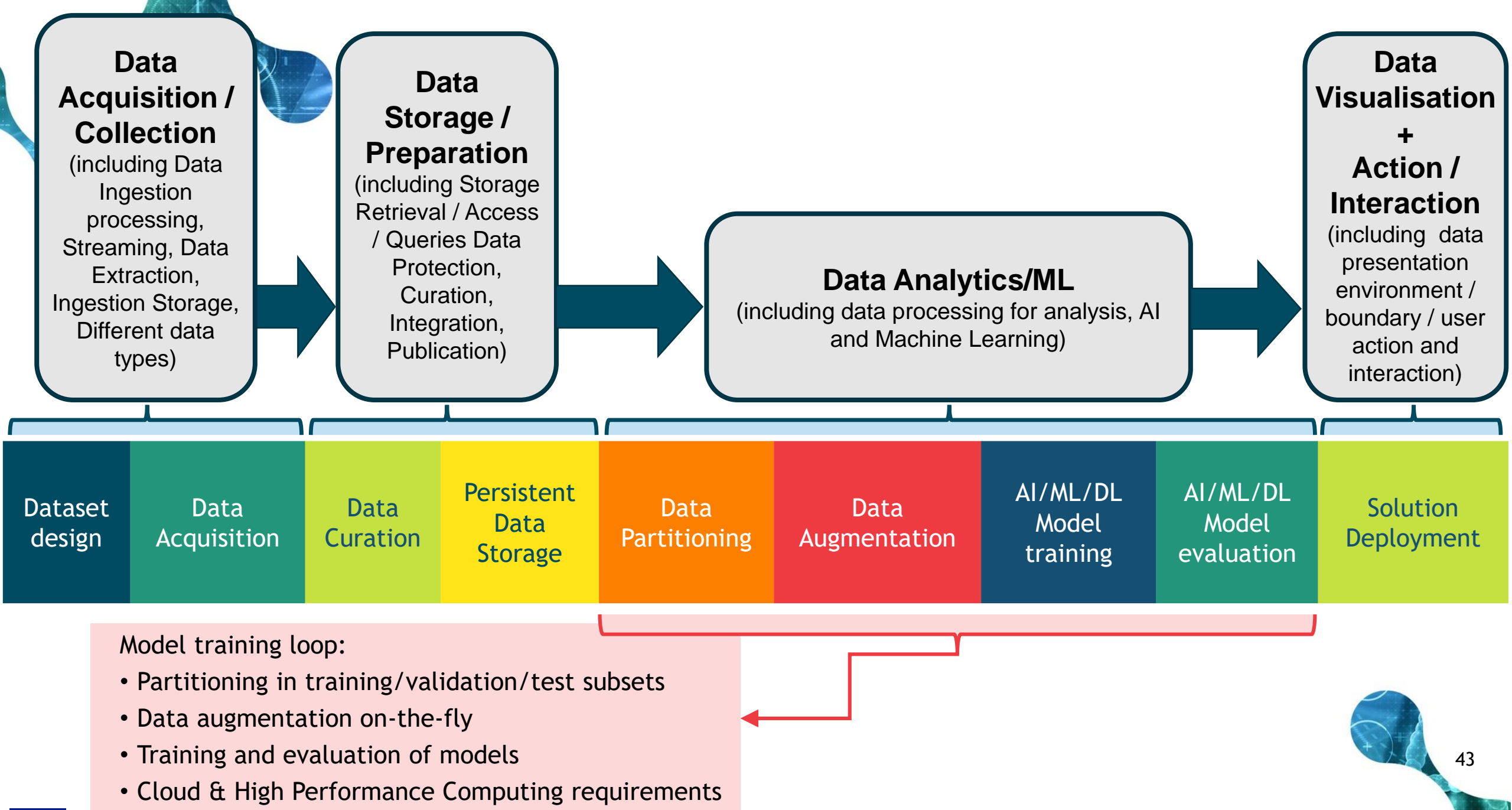
Both pipelines are suitable for business and research applications

The whole **Data Pipeline** is applicable to any sector. Our project is focused on the Health sector









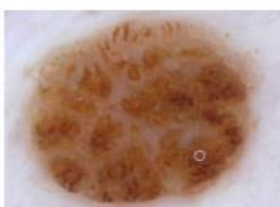


# Skin Lesion Detection and Classification

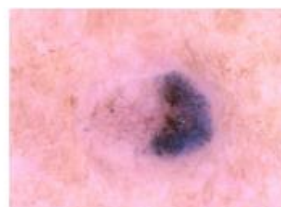
- Use case n° 12 of the DeepHealth project is based on the International Skin Imaging Collaboration dataset
- Aims: identification (**segmentation**) and diagnosis (**classification**) of skin lesion images among different classes



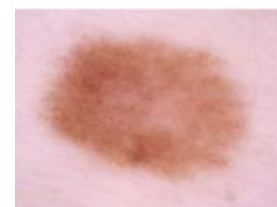
*Actinic Keratosis*



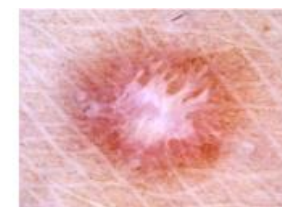
*Benign Keratosis*



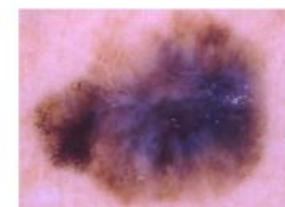
*Basal Cell Carcinoma*



*Nevus*



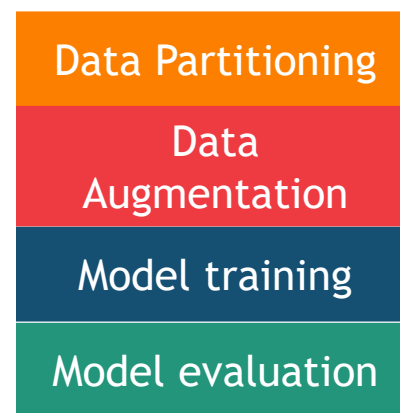
*Dermatofibroma*



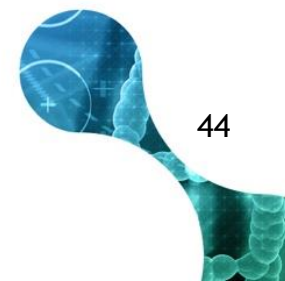
*Melanoma*



- Retrospective acquisition
- 23.906 annotated images
- Publicly available on the ISIC archive website
- jpeg data format



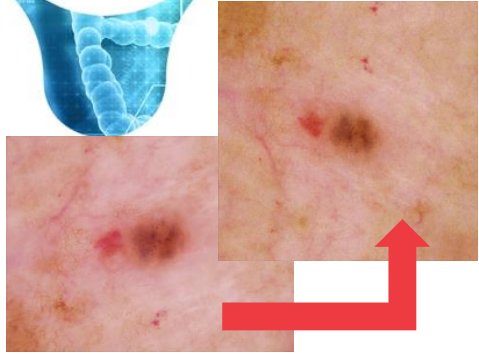
- Training 19.000
- Validation 906
- Test 4.000





# Skin Lesion Detection and Classification

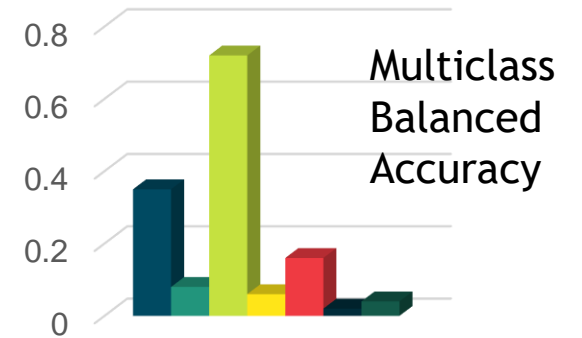
Performed using the DeepHealth toolkit. Models are already available in the front-end.



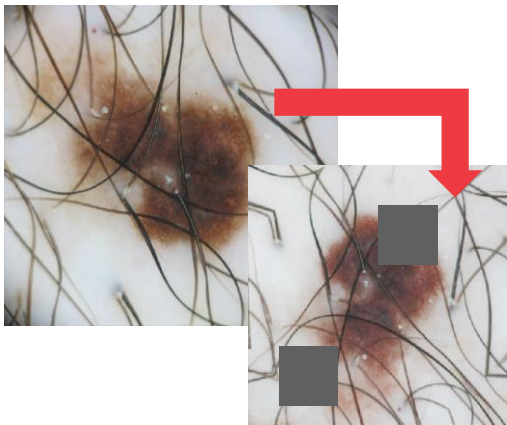
Data  
Augmentation



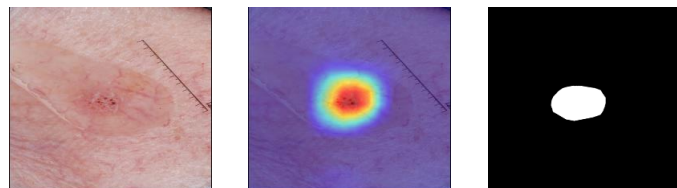
AI/ML/DL  
Model  
training



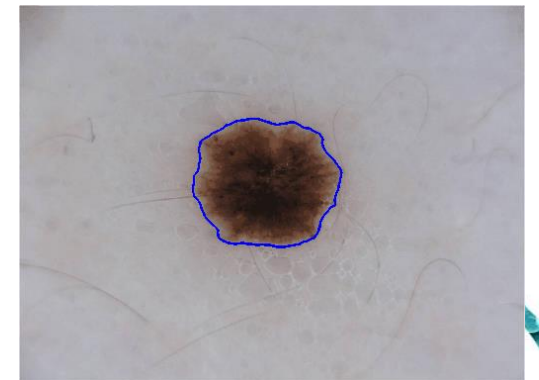
AI/ML/DL  
Model  
evaluation



**Explainability** plays a fundamental role in this context. Ensuring **Confidence Calibration** and providing a **Visual Explanation** of the models is essential to support clinicians.



Jaccard Index  
(Intersection over Union)





# Needs & Requirements

Evaluate datasets in terms of

1. **Findability** - where should a data scientist search for the dataset?
2. **Availability** - how long does a data scientist need to start the initial exploratory data analysis?
3. **Interoperability** - how long does a data scientist need to start training AI/ML/DL models with a dataset?
4. **Reusability** - are previously obtained results with a dataset public and available to other researchers / data scientists?
5. **Privacy / Anonymisation** - can the dataset be made public without compromising the identity of individuals?
6. **Quality** - is the dataset biased or unbalanced? What procedure has been followed to validate annotations?



# Needs & Requirements

Evaluate Deep Learning libraries in terms of

1. **Speed-up** - is distributed learning really efficient?
2. **Convergence** - does the distributed learning reach the same model accuracy in less time?
3. **Usability** - how long does a developer need to use the libraries effectively?
4. **Integrability** - how difficult is it to integrate the libraries as part of solutions to deploy?
5. **KPIs**: time-of-training-models (**totm**), performance/power/accuracy trade-off, etc.
6. **Others** - can you help us to evaluate other aspects?





# Needs & Requirements

## Evaluate Software Platforms in terms of

1. **Usability** - how long does a domain application expert need to manage the software tool effectively?
2. **Completeness** - does the application platform provide all the algorithms/procedures/functions to allow domain application experts to easily define the sequences of steps to implement the data and/or model pipelines?
3. **Compatibility** - how many data formats does the platform admits to import/export data and models from/to other frameworks?
4. **KPIs**: time-to-model-in-production (**ttmip**), time-of-pre-processing-images (**toppi**), etc.
5. **Others** - can you help us to evaluate other aspects?



DEEPHEALTH

# Questions?

Federico Bolelli [federico.bolelli@unimore.it](mailto:federico.bolelli@unimore.it)  
Jon Ander Gómez [jon@prhlt.upv.es](mailto:jon@prhlt.upv.es)  
Costantino Grana [costantino.grana@unimore.it](mailto:costantino.grana@unimore.it)  
Roberto Paredes [rparedes@prhlt.upv.es](mailto:rparedes@prhlt.upv.es)

<https://deephealth-project.eu>



@DeepHealthEU



The project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 825111.

# Contacts



[www.databench.eu](http://www.databench.eu)



[info@databench.eu](mailto:info@databench.eu)



[@DataBench\\_eu](https://twitter.com/DataBench_eu)



[DataBench](https://www.facebook.com/DataBench)



[DataBench Project](#)



[DataBench](#)



[DataBench Project](#)



DataBench



This project has received funding from the European Horizon 2020 Programme for research, technological development and demonstration under grant agreement n° 780966

EUROPEAN  
**BIGDATA**  
**VALUE** FORUM

**BERLIN + VIRTUAL**  
**3-5 NOVEMBER 2020**