



The vision,  
the results,  
the future

December 14, 2023



Information Day on a new era of STI policy making with AI

# Our vision: background



- Open Research Data
- Using AI to obtain evidence for policy-making
- Appealing visualization of results
- Catalogue of web services:
  - For agenda setting
  - For proposal management

# IntelComp's Goal



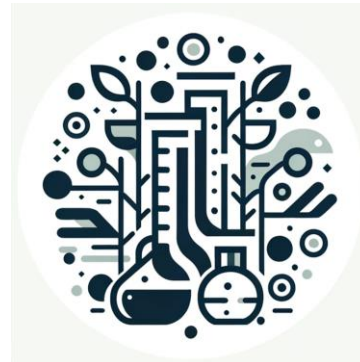
Use AI tools to provide policy makers with evidence based on data



*Industrialize the approach*



*Expert in the loop*

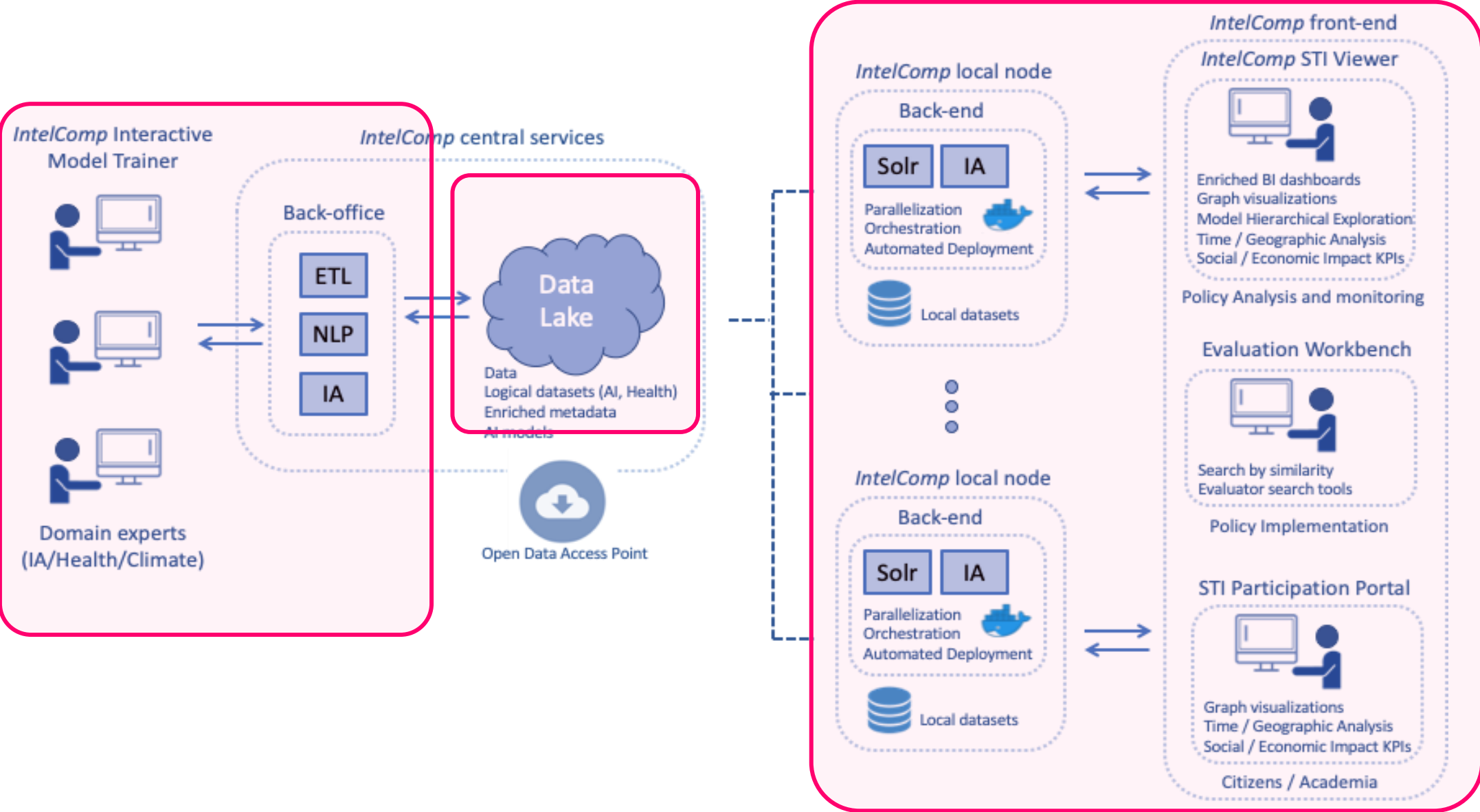


*Living Labs*



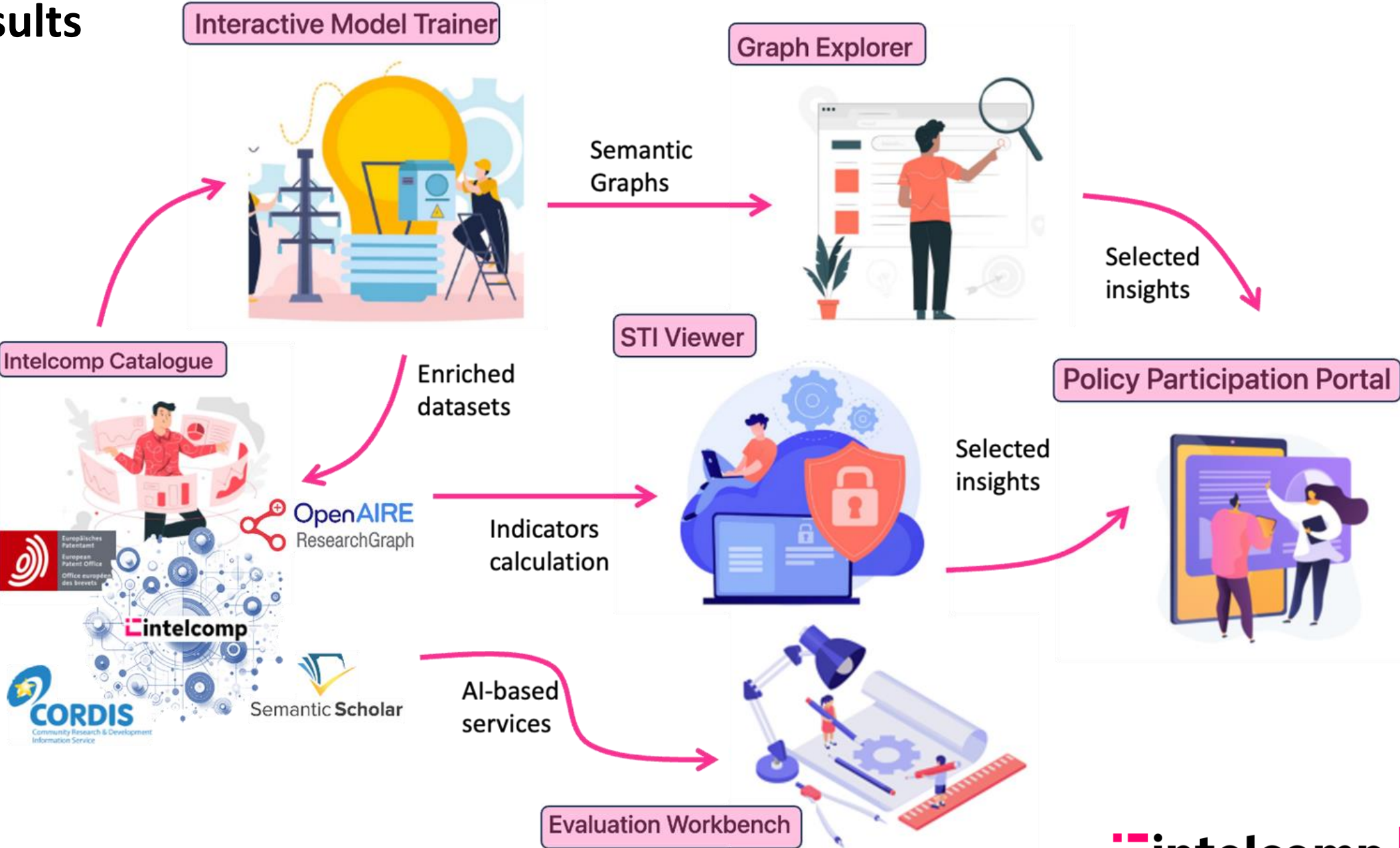
*Integrated Platform*

# IntelComp's Concept & Original Design

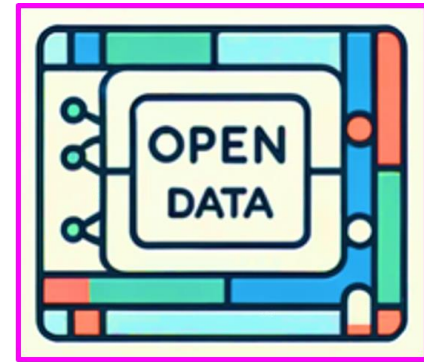




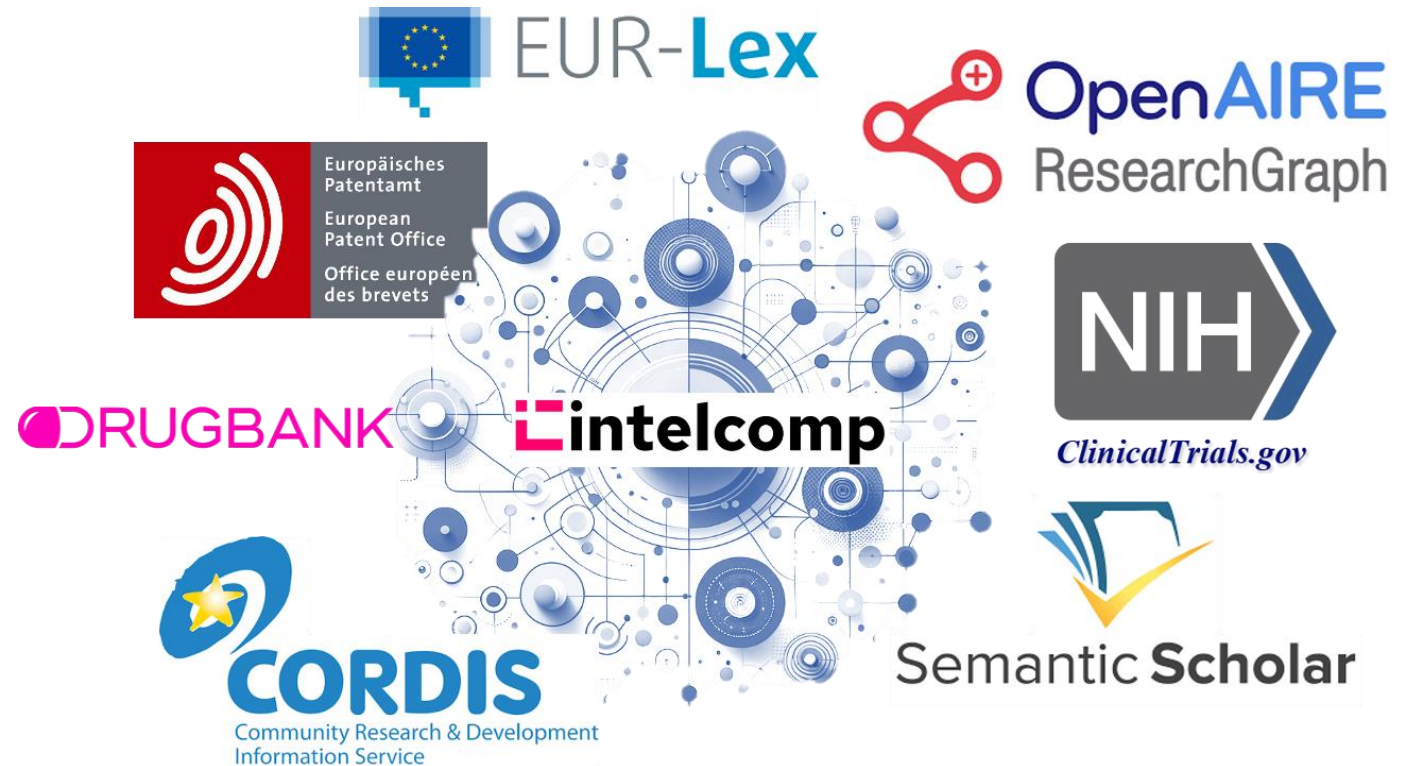
# The results



# Data Catalogue



- IntelComp's data lake collects several STI-related open data sets
  - Metadata heterogeneity
  - Variable curation quality
  - Very large datasets
- IntelComp's approach:
  - Unstructured datalake
  - Distributed storage (HDFS, parquet)
  - English as an anchor language
  - Text-based representations



# Technologies

- Scalable NLP and NMT
  - GPU processing
  - Spark NLP for parallel processing
- Text-based AI
  - Supervised classification
  - Zero-shot classification
  - Document domain selection
  - Topic Modeling
- LLMs
  - Transformer-based document representations
  - GPT models to improve interpretability



# AI and NLP service catalogue

## System for subcorpus generation

- Transformers
- Expert in the loop approach
- Relevance feedback

## Semantic Similarity & Graph Analysis

### Interactive Model Trainer

- Find documents with similar semantics
- Graph-based impact analysis

The goal is better experience

with the STI Viewer and the EWB

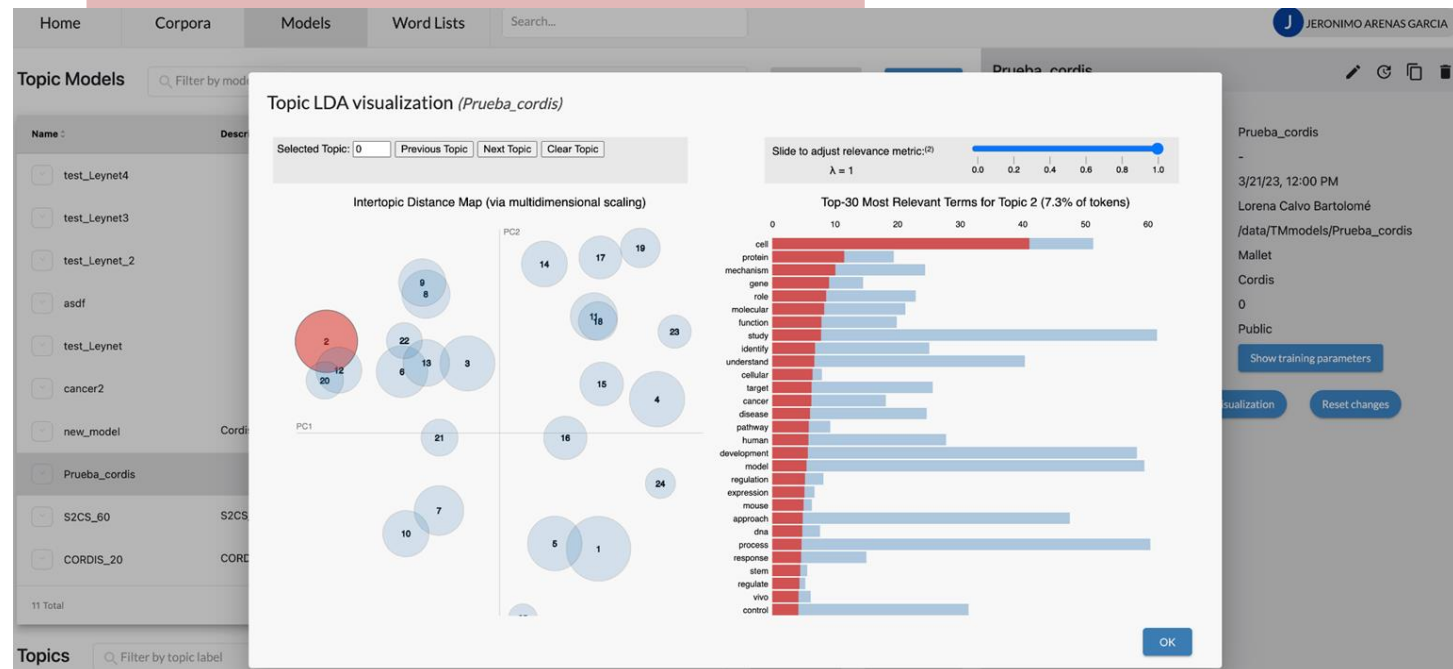
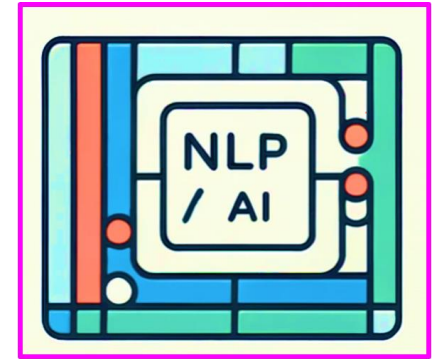
## Classification service

- Maximize alignment with domain experts' intuition
- Transformers
- Zero-shot classification

- Improve recommendations of the EWB

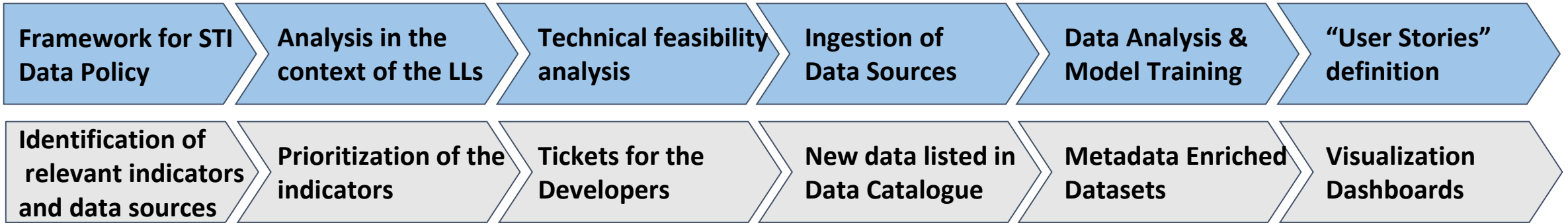
## Topic Modeling Service

- Mallet
- Neural solutions (prodLDA, CTM)
- Expert assisted Curation tools
- Graphical tools for topic model exploration

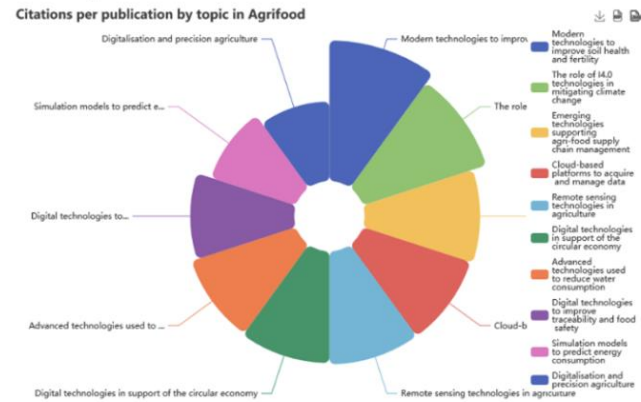




# Indicators calculation and visualization: STI Viewer



Scientific Impact Characteristics



**Description/Source:** The graph shows the average number of citations per scientific publication for different topics in the Agrifood domain. The assignment of topics is taken from Abbate, Centobelli, Cerchione (2023) (<https://doi.org/10.1016/j.techres.2023.102222>) which we enhanced and mapped to topics in our Fields of Science (FoS) publication classification system (<http://doi.org/10.21203/rs.3.rs-2875191/v1>). Data Sources: OpenAIRE Graph using the metadata of the scientific publication to the SDGs that are related to it. Data Sources: OpenAIRE Graph using the metadata of the scientific publication to the SDGs that are related to it. Data Sources: OpenAIRE Graph using the metadata of the scientific publication to the SDGs that are related to it.

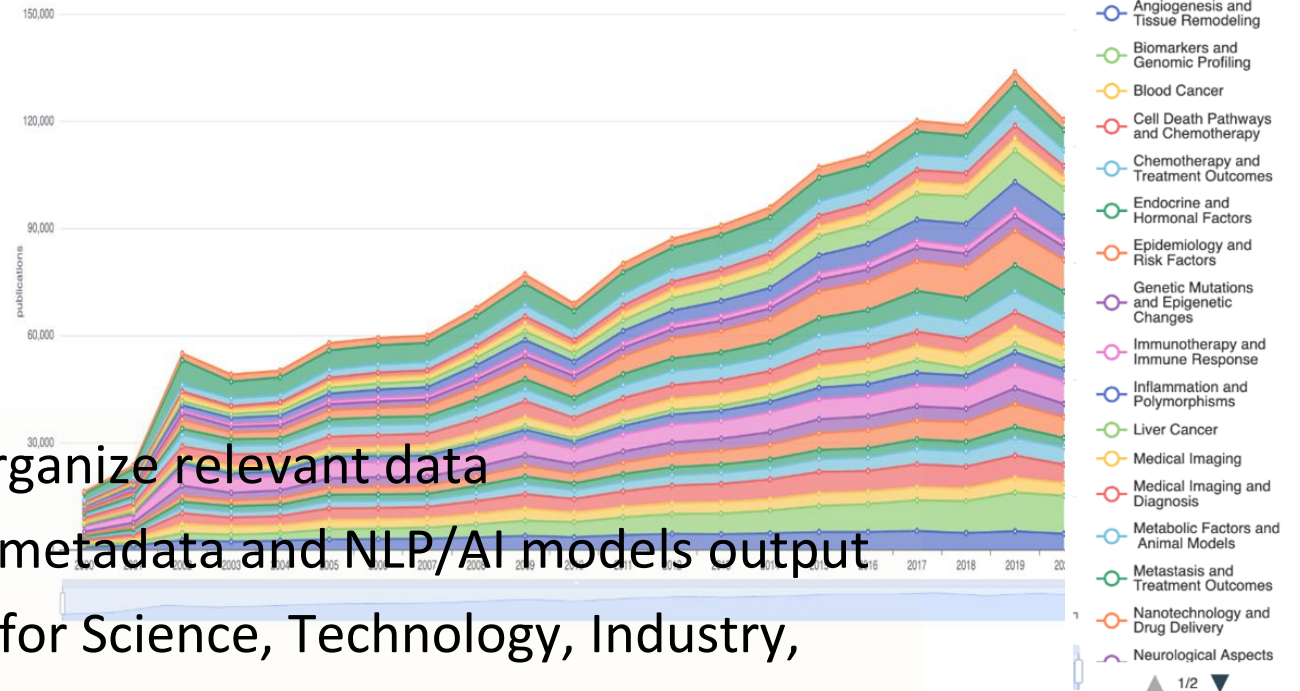
Scientific Impact Characteristics

Citations per publication by SDG



**Description/Source:** The graph shows the average number of citations per scientific publication in Agrifood, by UN Sustainable Development Goal (SDG, <https://sdgs.un.org/goals>). Our SDG classification system uses deep learning models to map the scientific publication to the SDGs that are related to it. Data Sources: OpenAIRE Graph using the metadata of the scientific publication to the SDGs that are related to it. Data Sources: OpenAIRE Graph using the metadata of the scientific publication to the SDGs that are related to it. Data Sources: OpenAIRE Graph using the metadata of the scientific publication to the SDGs that are related to it.

Share of publications in Cancer topics over time



User-defined storylines to organize relevant data

Visualizations based on raw metadata and NLP/AI models output

Used in the three living labs for Science, Technology, Industry,

regulations indicators

# Graph explorer



- Nodes represent projects / publications, etc ...
- Location is associated with semantic closeness
- Hierarchical exploration
- Use metadata to:
  - Filter nodes
  - Coloring criteria
- Export to Participation Portal

# Evaluation of Proposals: Evaluation Workbench

Corpus: hfri | Model: hfri-30

Model Overview | Similarities | Document Search |  All Topics |  Relevant Topics

Option for document similarity search

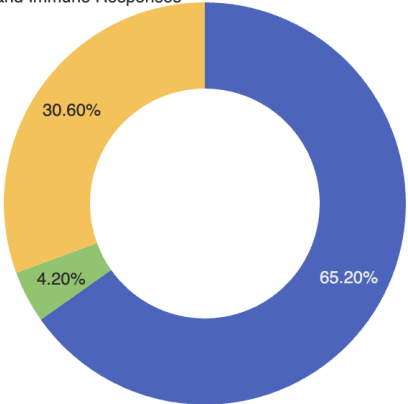
Similarity of an inputted document and others in the corpus

Document Text to search similars to

Find Similar Documents

Taxonomy: fos

Cell Biology and Molecular Mechanisms
  Urban Sustainability and Business Innovation
  Clinical Studies and Immune Responses



Classify

Deep Learning has been exploited in different directions in the context of AI. Recent developments in NLP have lead to Large Language Models that excel in conversational capabilities

Classify

Results:

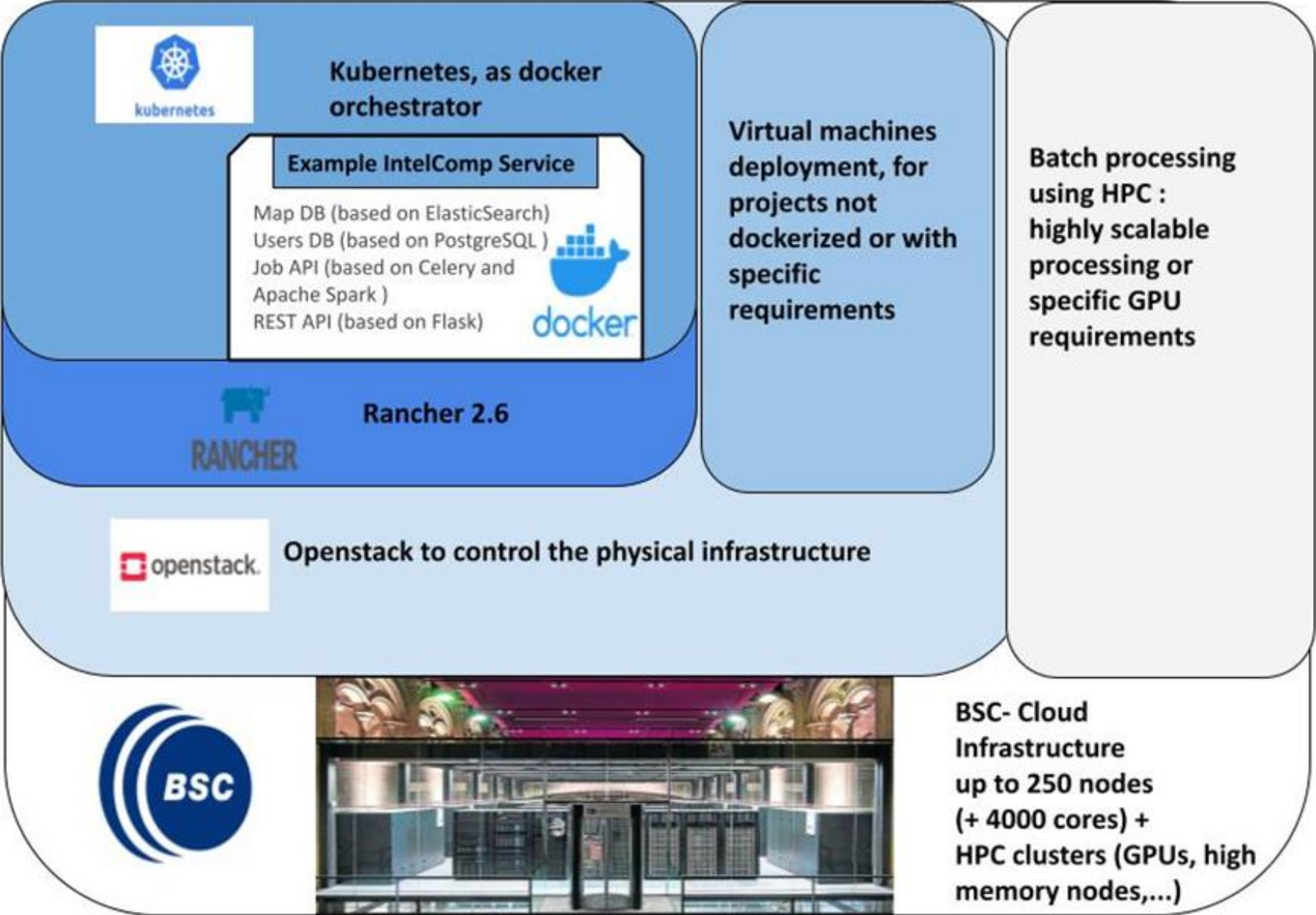
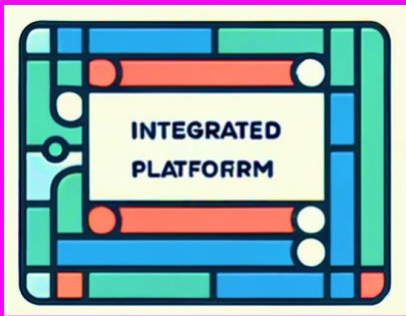
Code	Punctuation
Computer Science	0.96092784

Punctuation

Powered by Topic Models and efficient services provided by Python and Solr backend:

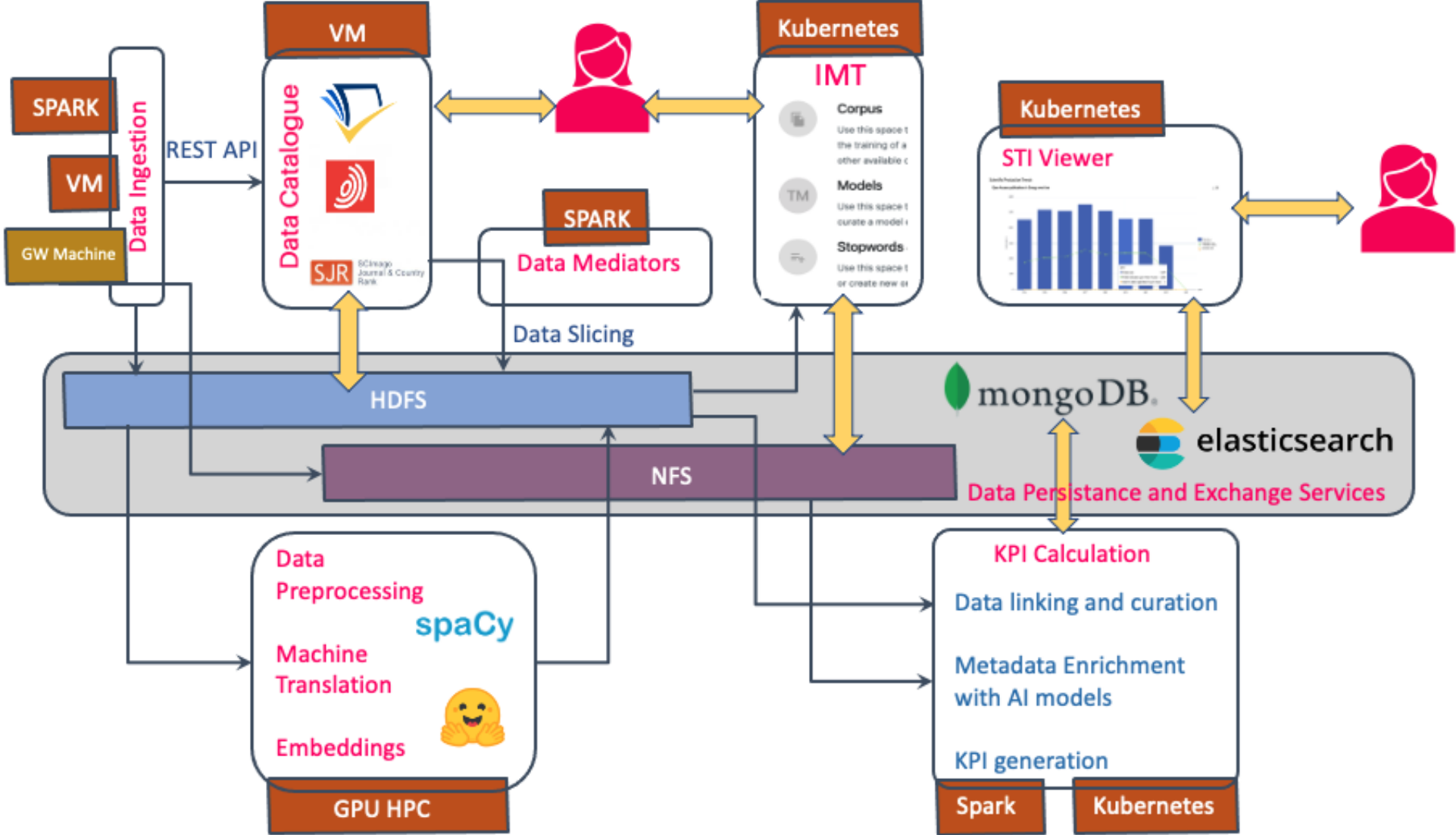
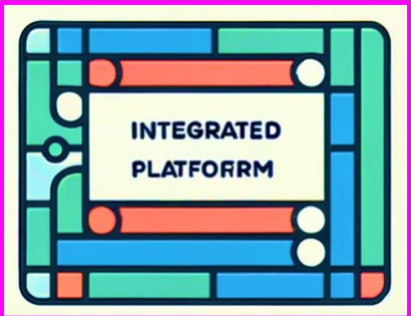
1. Topic analysis of projects / experts
2. Search by similarity of projects / experts
3. Taxonomical classification of project abstracts

# Infrastructure





# Service Deployment





# The future

- Platform infrastructure active in BSC for 2 years
- Ensuring and simplifying high-value workflows
- Data maintenance and update
- Continue fine-tuning of web services
- More intuitive interaction with the platform, LLM-like





intelcomp

The vision,  
the results,  
the future

<https://apps.intelcomp.bsc.es>

