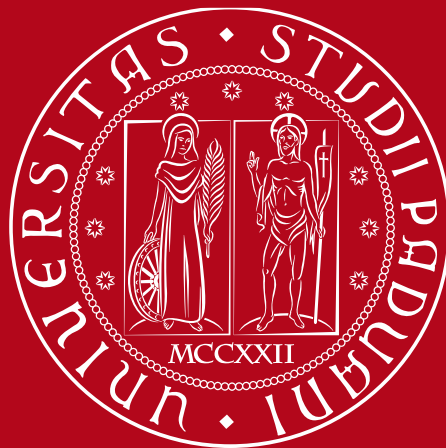DIPARTIMENTO
DI INGEGNERIA
DELL'INFORMAZIONE

UNIVERSITÀ
DEGLI STUDI
DI PADOVA

# Federated Data Analytics for Genomics Data

Student: **Mirco CAZZARO** – mat. **2076745**

Supervisor: Prof. **Gianmaria SILVELLO**

Master Degree Course in **Computer Engineering** – **Web Information and Data Engineering** (WIDE)

A. A. **2023/2024**

# The Challenge

Biomedical data is becoming increasingly complex. The rapid evolution of data storage systems has created a landscape where:

- Diverse data models exist (relational, hierarchical, graph-based);
- Integration of these models is crucial for making data meaningful.

Understanding genetic diseases and advancing personalized treatments relies on integrating and analyzing this diverse data effectively.

# Background

**Methodologies**:

1. Data **Federation**

2. Data **Virtualization**

3. Semantic Data **Integration** (OBDA)

4. Knowledge Graphs, Ontologies & RDF

**Technologies**:

1. SPARQL query language;

2. Data **Virtualization** solutions: <u>Denodo, Dremio, Teiid</u>

3. **OBDA** Solutions: <u>Mastro, Ontop</u>

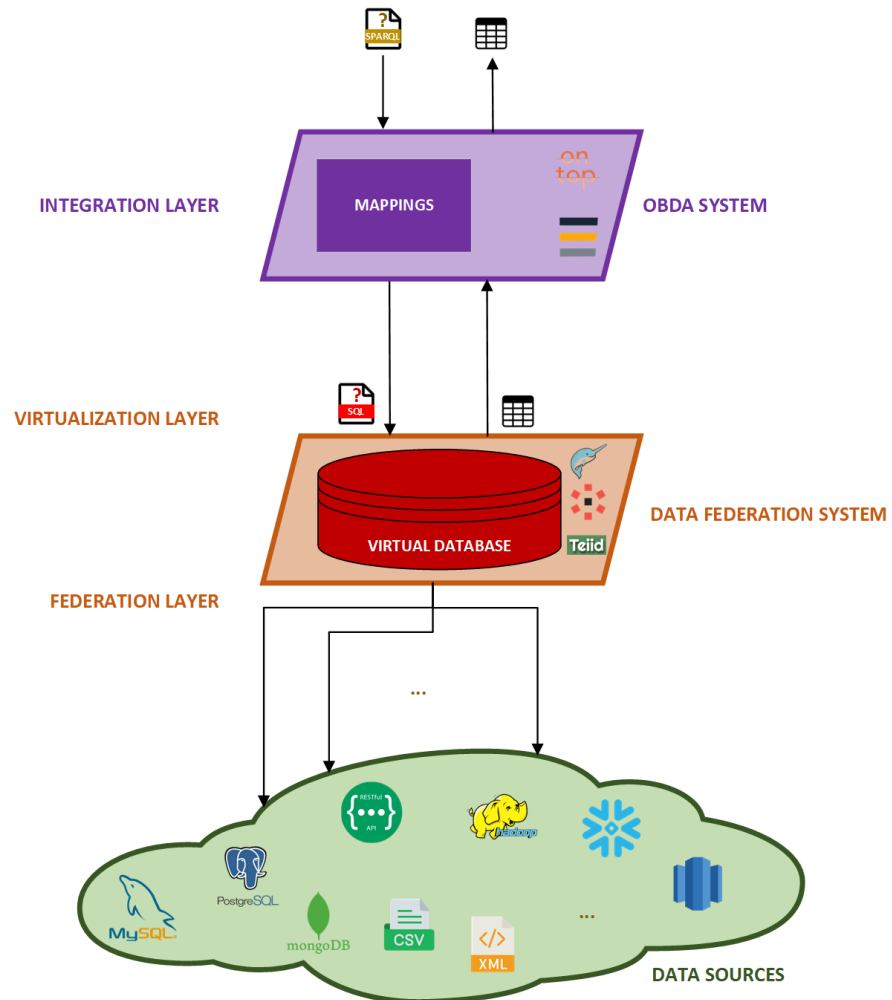| | Support | Free and Open Source | Well Documented | Scalability | Solid Logging Capabilities |
|---|---|---|---|---|---|
| Denodo | ✓ | | ✓ | ✓ | ✓ |
| Teiid | | ✓ | | | |
| Dremio | ✓ | ✓ | ✓ | ✓ | ✓ |

# Proposed Solution

**Federated Data Analytics System**

A system designed to integrate and analyze clinical and genomics data seamlessly.

- Objective: Enable researchers to perform complex queries across multiple datasets without extensive preprocessing.

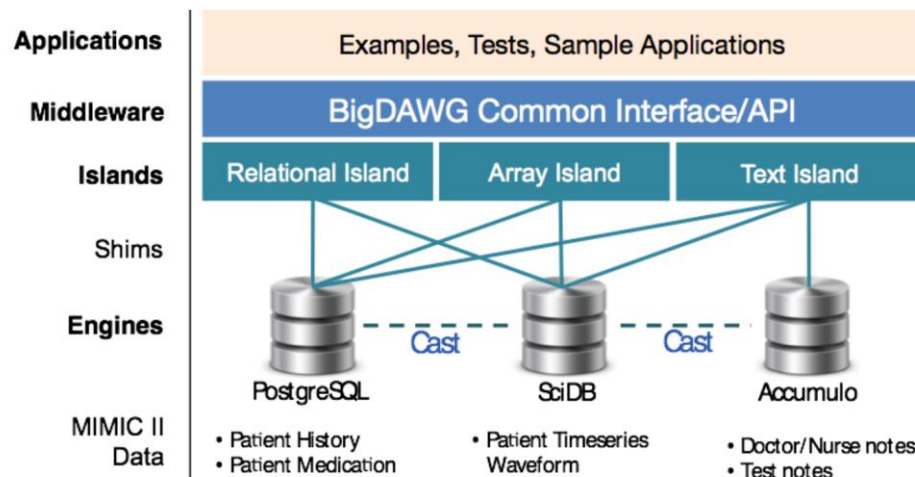- Key Components: <u>OBDA</u>, <u>Data Virtualization</u>, <u>Ontology</u>.

This approach reduces complexities in biomedical data management and accelerates research.
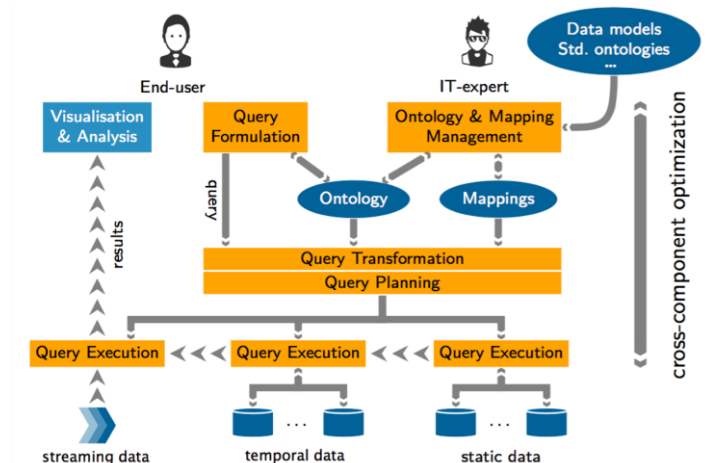
# Related Polystore Solutions

## BigDAWG

- It is structured in 4 layers, from the data stores layer up to the application layer.

- A crucial role is played by the "**Island layer**", a layer of autonomous components that abstracts data sources.



## Optique

- Initially conceived as an EU funded project, now implemented and maintained by Siemens.

- 3-layers architecture, **OBDA based**, optimized for streaming data. Provides a visual query tool to ease its use.

# EU Project HEREDITARY

**Objective:**

Transforming our understanding of brain diseases through integrated multimodal data analysis.
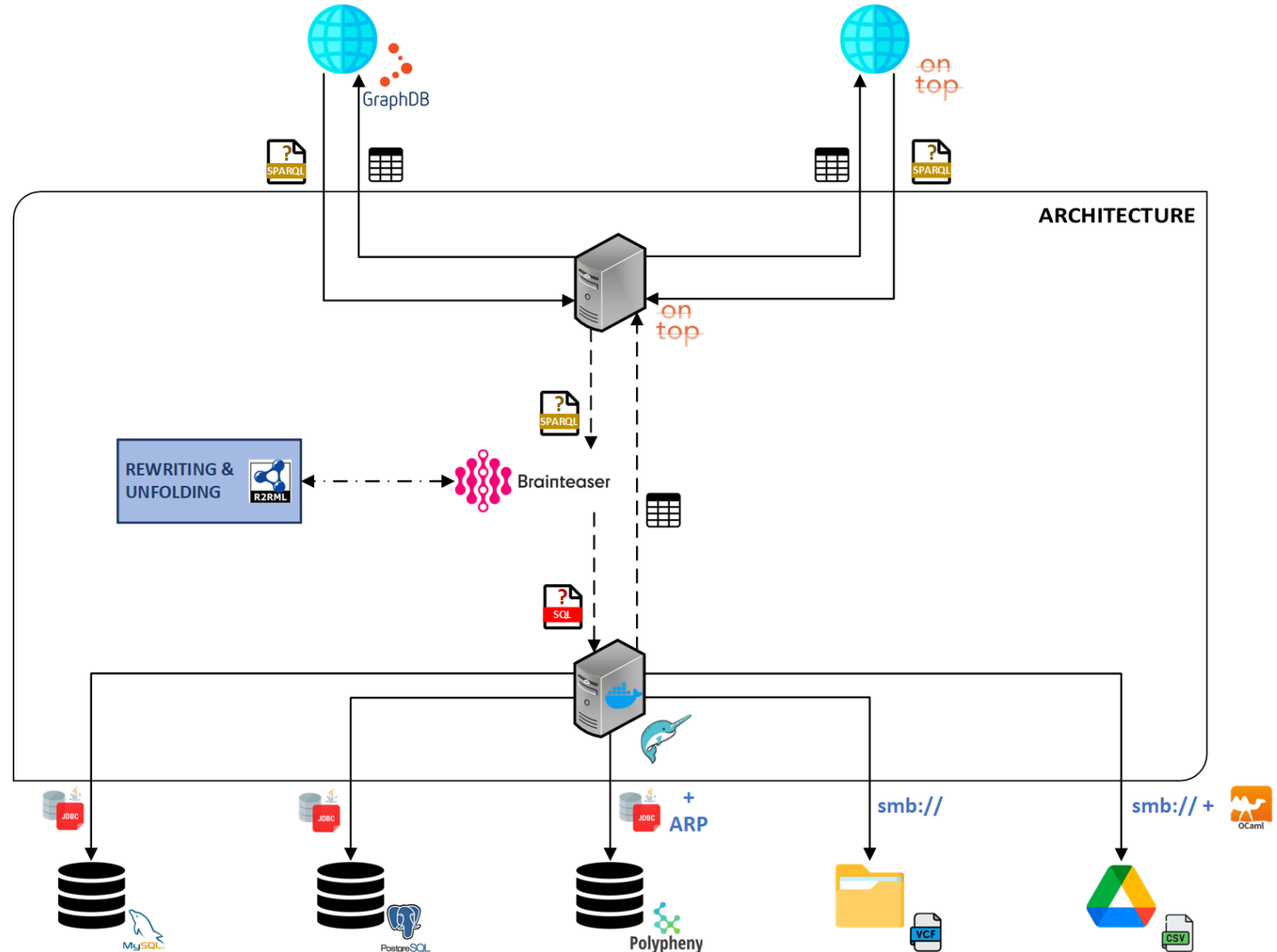
**Focus Areas:**

1. Integrating genomic and clinical data from various European stakeholders;

2. Addressing the challenges of data heterogeneity;

3. Ensuring privacy compliance and data security.

**(one of the) Outcomes:**

A unified DBMS that supports advanced querying capabilities across different data sources and that allows to gather insightful analytics.

# System Architecture



-  **Ontop**: Semantic Data Integration Layer.

-  **Dremio**: Data Federation and Virtualization Layer.

-  **BRAINTEASER Ontology**: Provides the vocabulary structure for querying.

# System Architecture – Data Sources

- **MySQL**: Stores structured clinical data.

- **PostgreSQL**: Manages additional clinical datasets.

- **Polypheny**: Supports various data models, enhancing flexibility.

- **NAS Shared Folders**: Hosts genomics data in VCF files.

- **Google Drive**: Integrates a cloud-based storage use case, hosting CSV files.

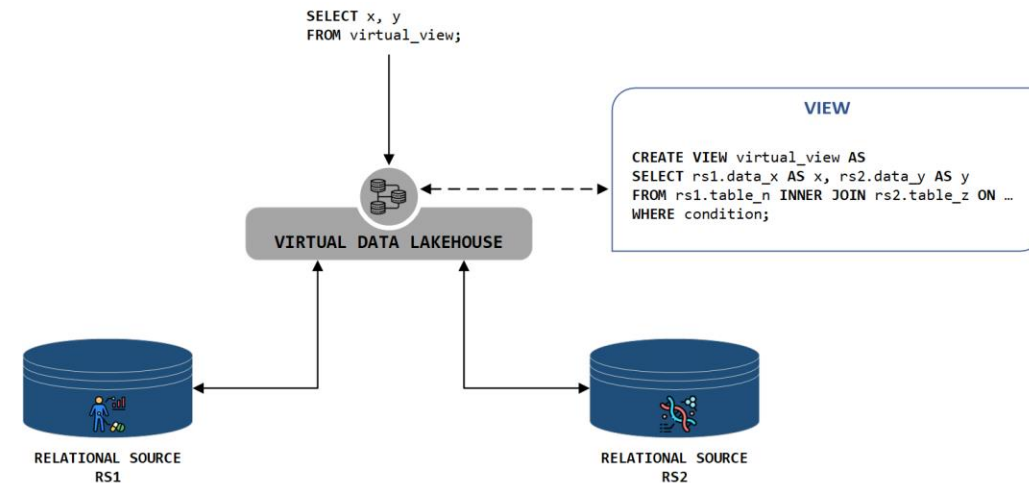Each data source is connected via Dremio's federation layer, enabling a unified view across diverse data repositories.

# Federation and Virtualization

**Virtualization Layer:**

Creates virtual views without materializing data, optimizing
performance and enabling complex queries across the
federated system.

**Federation Layer:**

Integrates data across various sources, ensuring seamless
access and querying.

Dremio serves as the core engine for both layers,
providing a robust platform for handling large-scale data
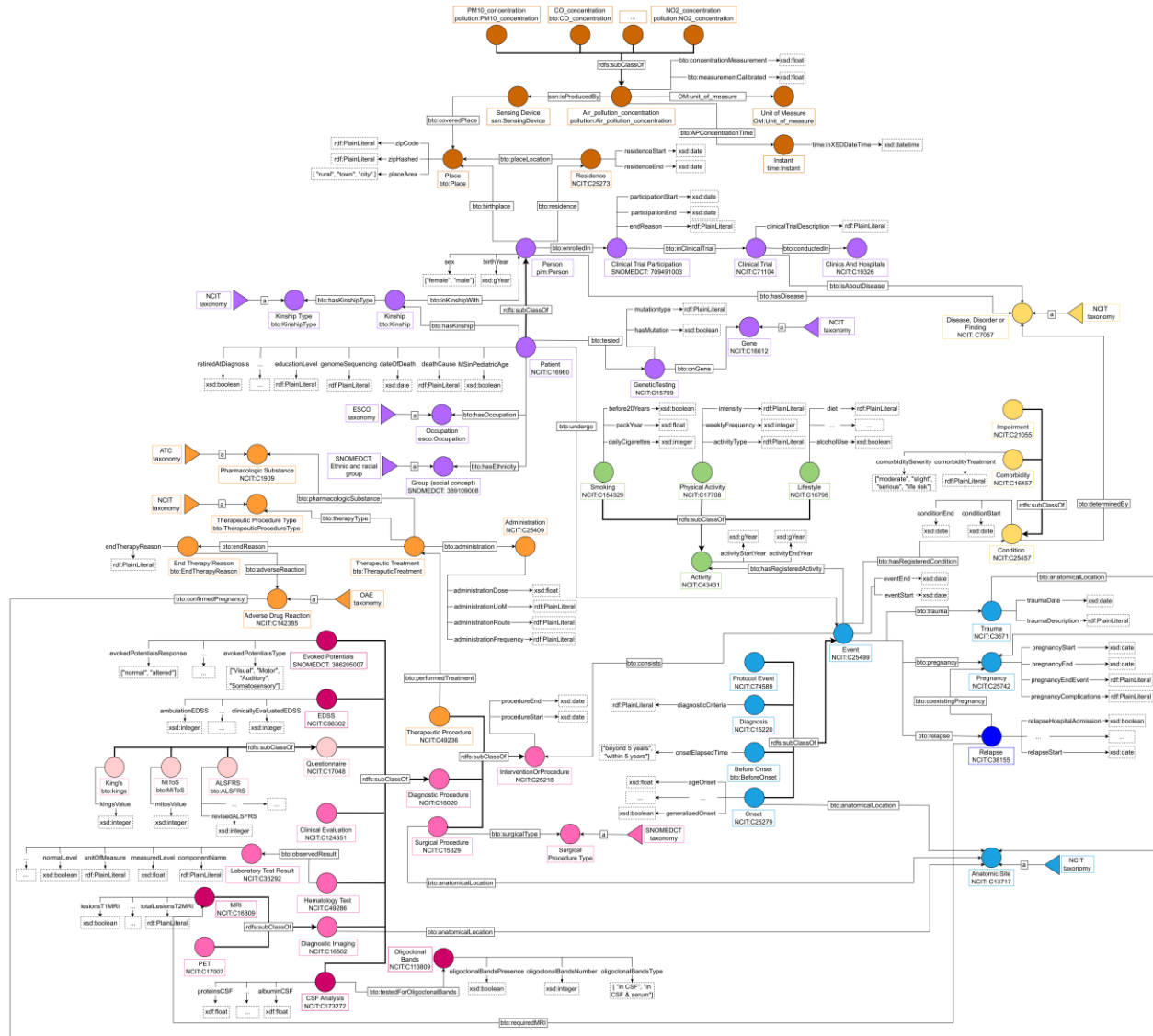integration tasks.

# Ontology Integration

The **BRAINTEASER** Ontology is central to the system

**Features:**

- Models relationships between clinical and genomics data;
- Enhances query capabilities through semantic enrichment;
- Ensures interoperability and scalability.

# Use Case: ALSFRS Data Querying

Use Case: Querying ALSFRS (Amyotrophic Lateral Sclerosis Functional Rating Scale) data across multiple sources.

- **SPARQL Query**: Retrieves comprehensive patient data;

- **Process**:

    - Query Rewriting;

    - Unfolding into SQL;

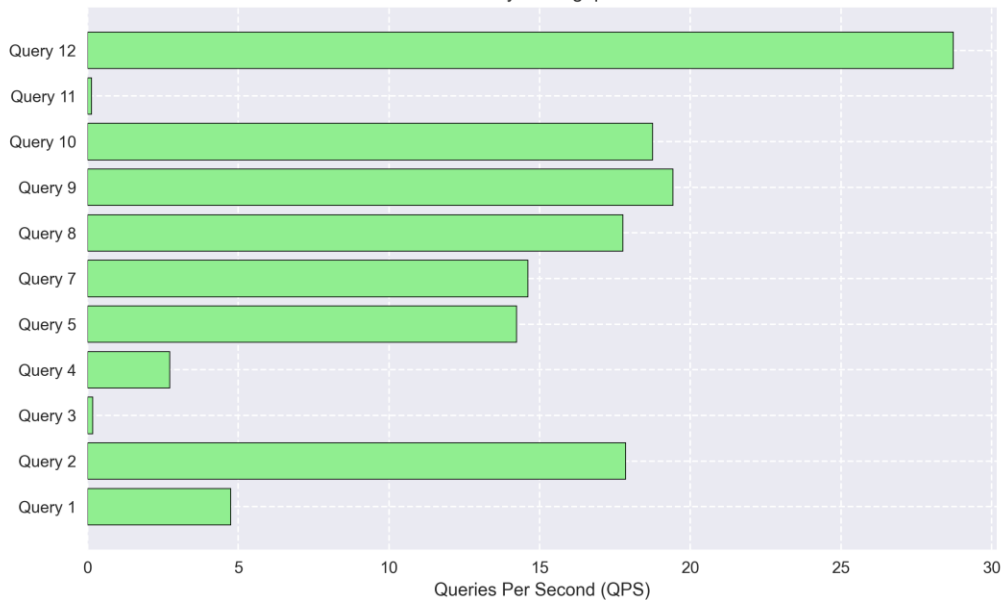    - Execution in Dremio.

Original SPARQL Query

Unfolded SQL Query

# Benchmark 1: BSBM

**BSBM Benchmark**

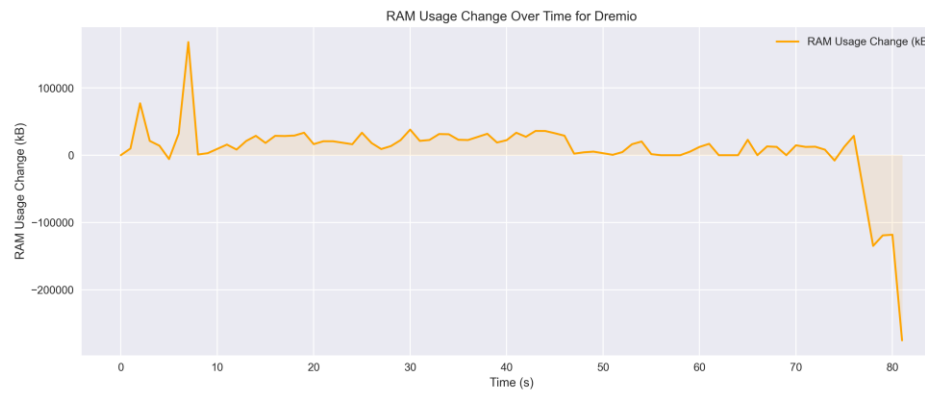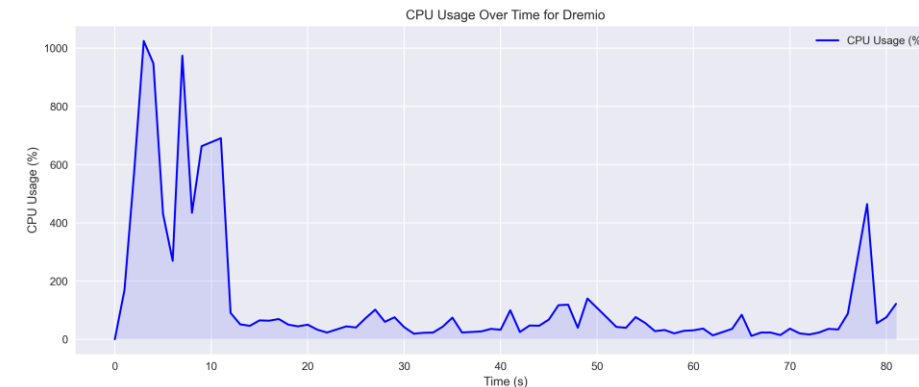Query performances evaluation over a synthetic dataset across diverse sources.
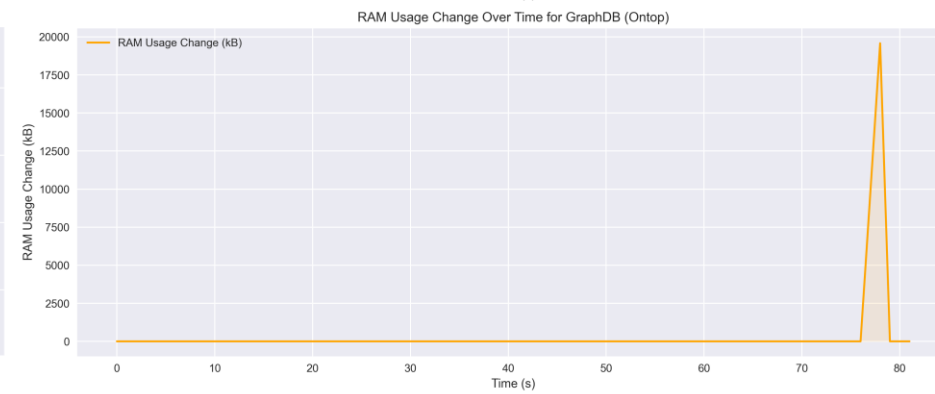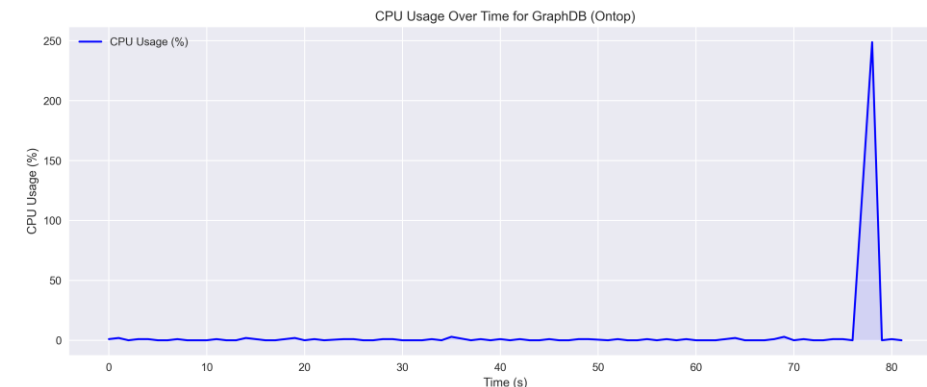
# Benchmark 2: SEASHELL

**SEASHELL Benchmark**

Real-world clinical data highlights system efficiency in handling complex queries.



Dremio

GraphDB

# Conclusions and Future Works

The federated data analytics system developed in this thesis offers a powerful tool for integrating and analyzing heterogeneous biomedical data.

**Future Directions:**

1. **Optimization**: Improve architecture performances, considering SEASHELL monitoring data as an entry point;

2. **Privacy**: Strengthen data security and compliance with regulations, such as GDPR;

3. **Usability**: Simplify both the architecture usage and deployment, so to reach a larger audience.

This system lays the groundwork for future advancements in biomedical data research, contributing to better healthcare outcomes.

# Thank You!

**Do you have any question?**