

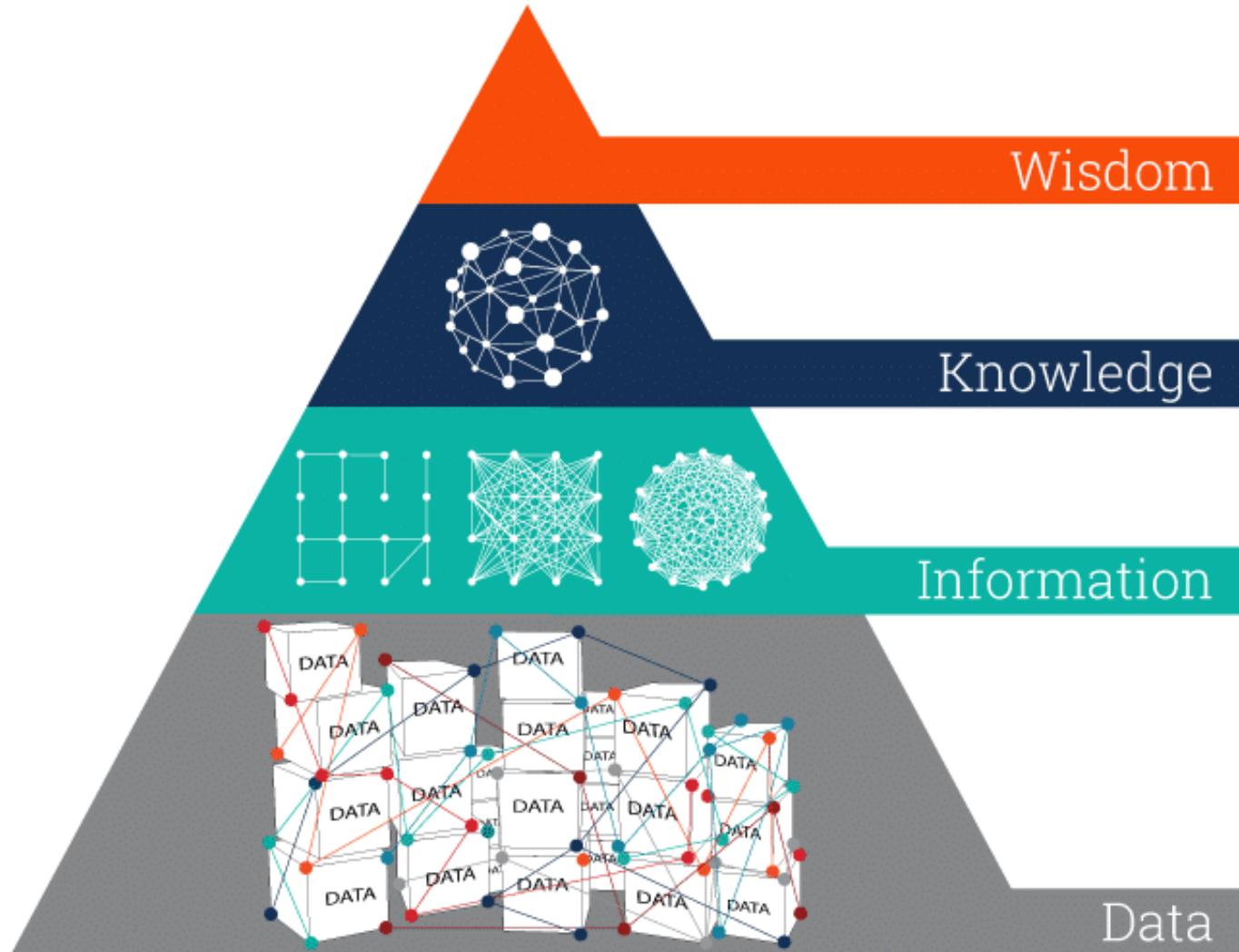
# **Metadata, Data Provenance, Metadata, and looking to the future with Data Mesh**

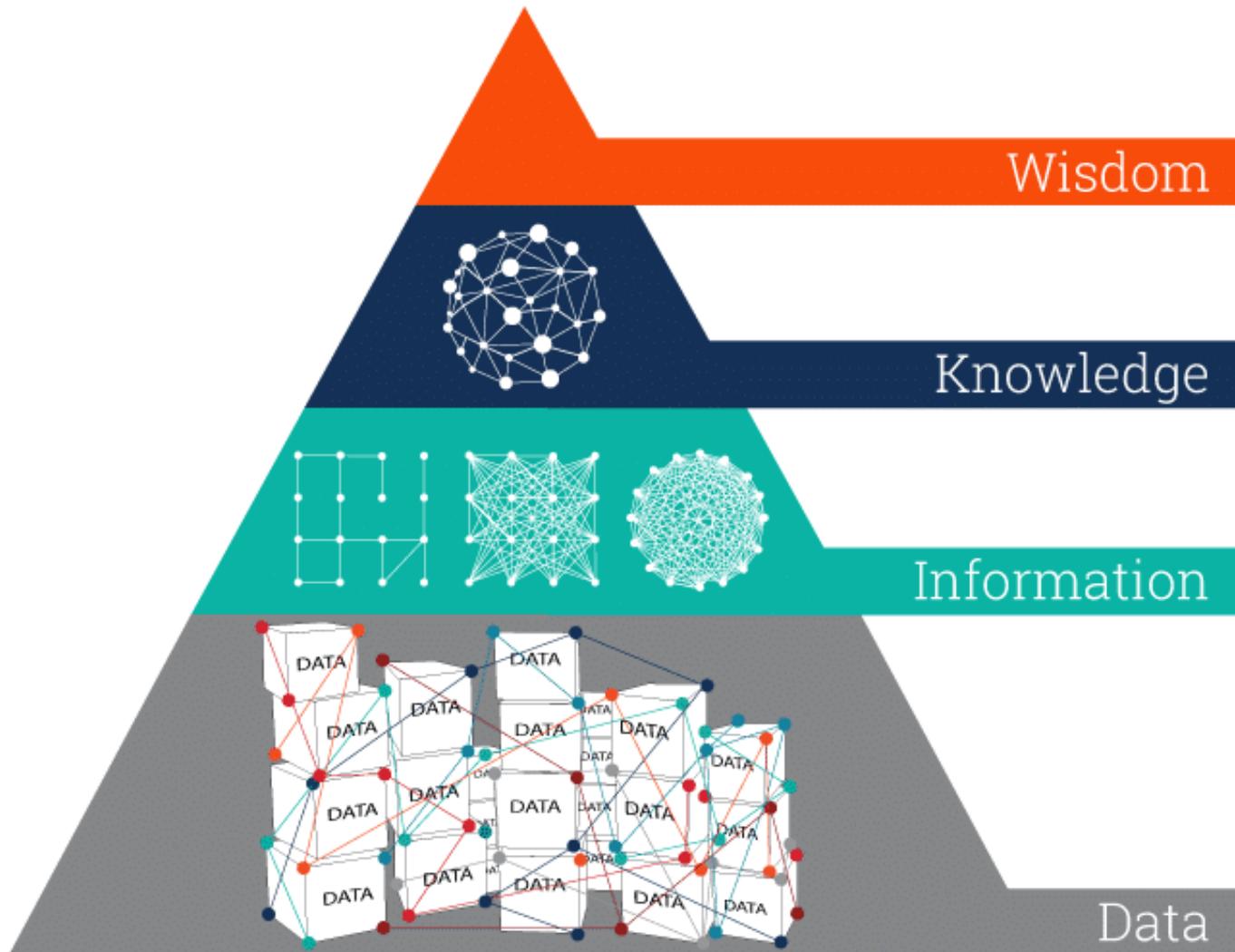
# Agenda



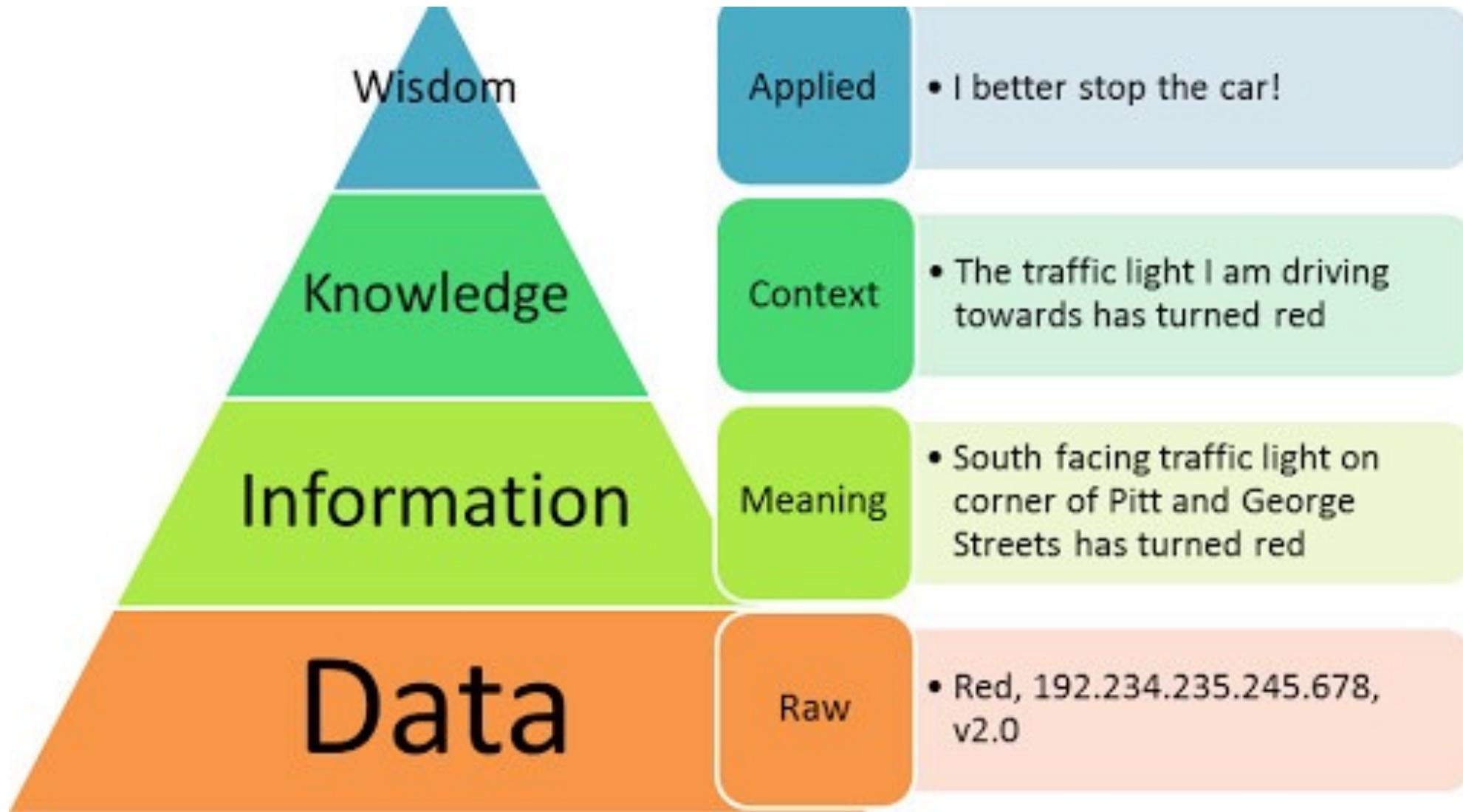
- The DIKW Pyramid!
- Metadata, Metadata, METADATA!
- (Some) Tools of the trade
  - Apache Atlas
  - LinkedIn Datahub
  - Current Research
- Data Mesh

# **Data vs Information vs Knowledge vs Wisdom**





Each step up  
the pyramid  
answers  
questions  
about and  
**adds value**  
to the initial data.



# **Metadata, Metadata, METADATA!**

# What is Metadata?



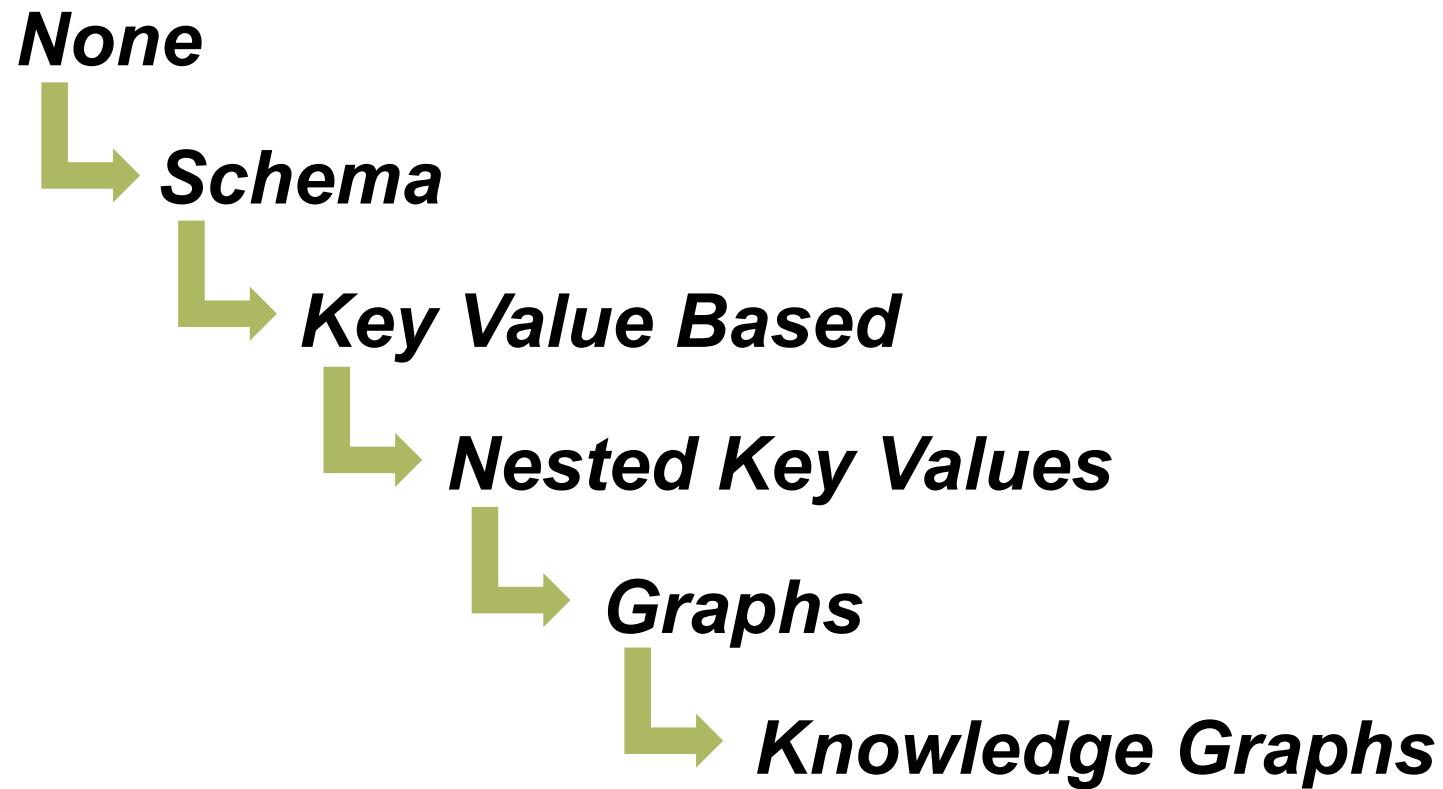
# Why have Metadata?

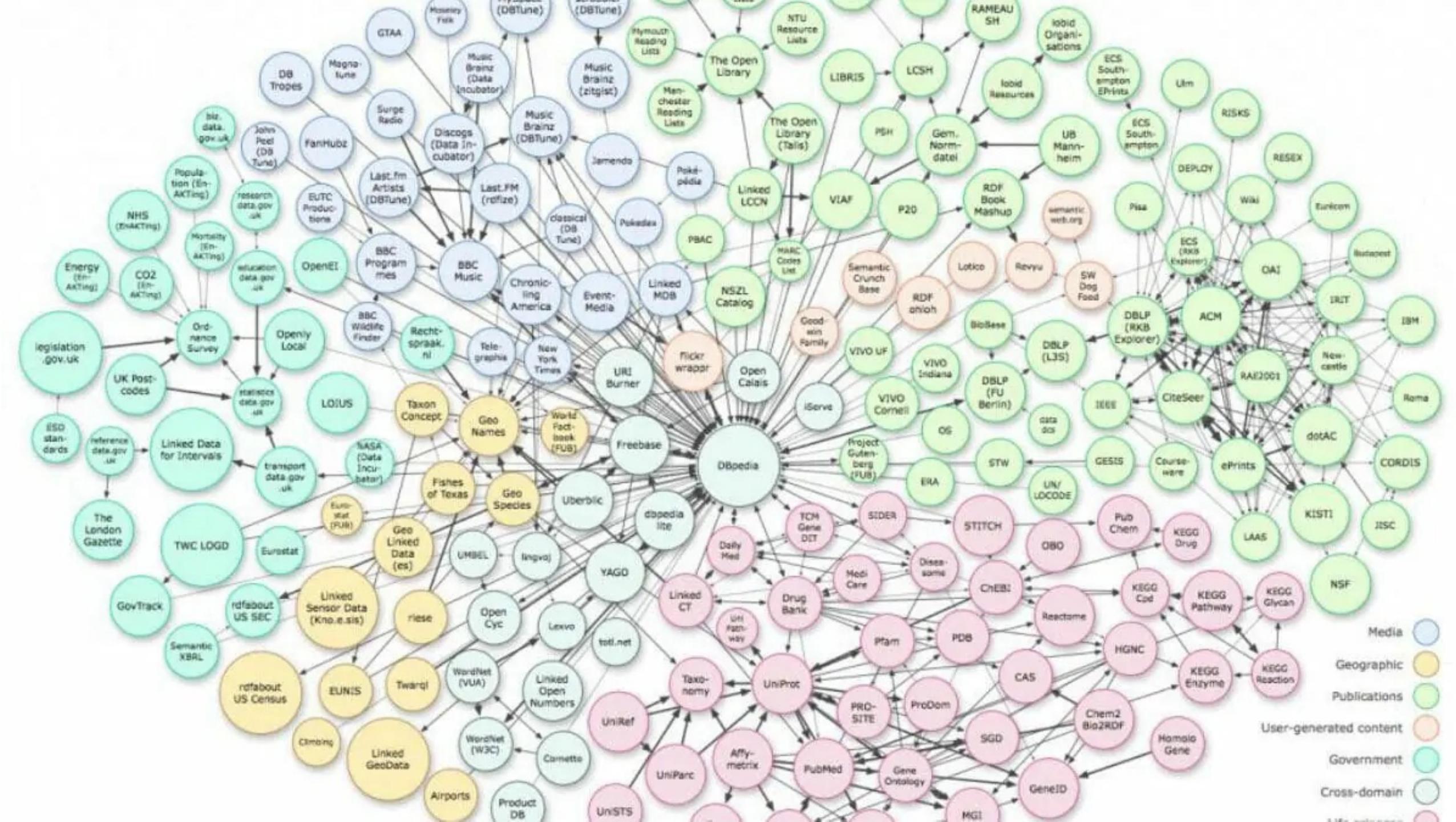


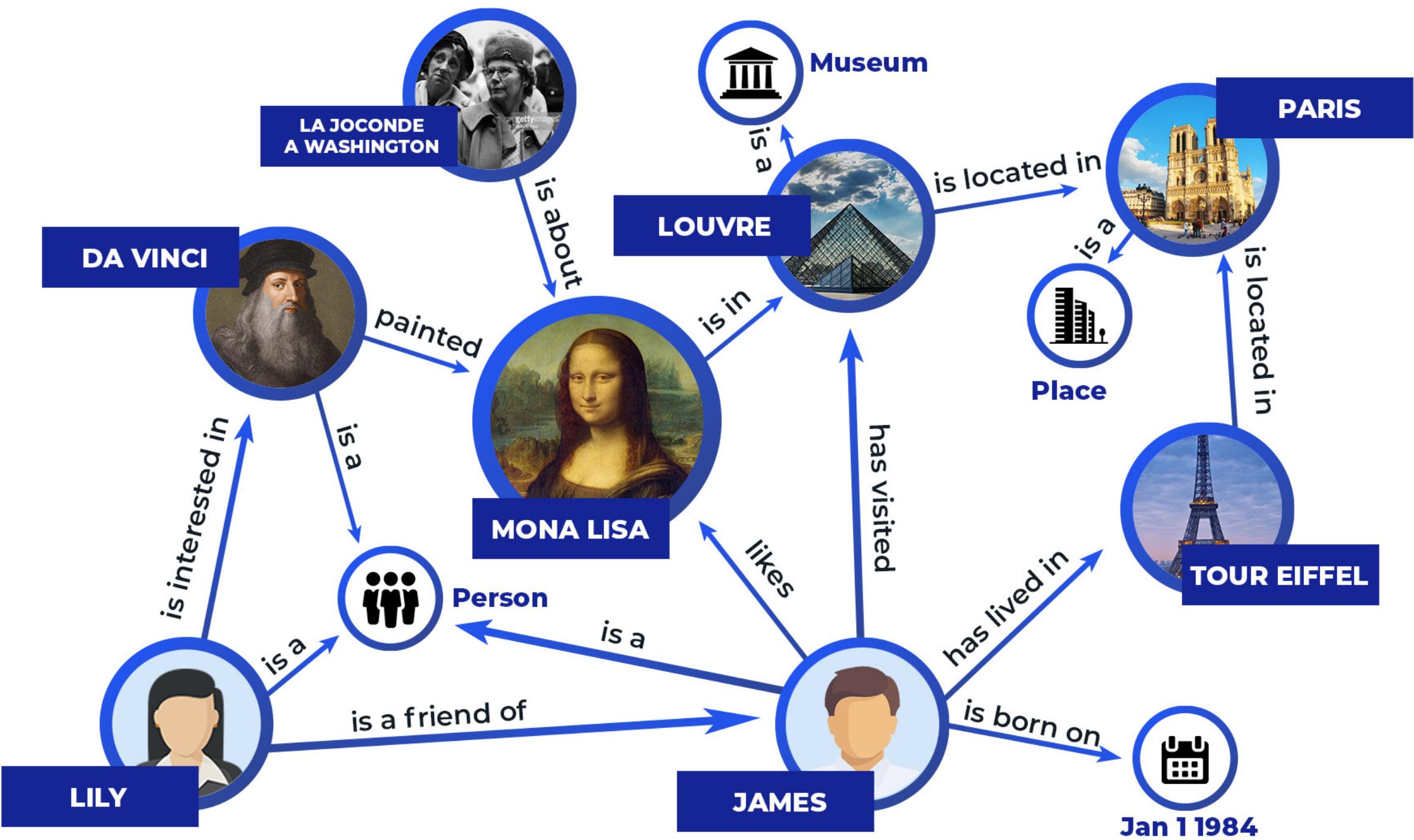
# Types of metadata

Technical <i>Definitional</i>	Operational <i>Descriptive</i>	Business <i>Descriptive</i>	Social <i>Descriptive</i>
Schemas, data types, models, etc.	Process outputs, lineage metadata, ETL, performance metadata, etc.	Data tags, classifications, mappings to business relationships, etc.	Metadata about user-generated content, business user knowledge, etc.

# Typical history of descriptive metadata







# What is Data Provenance?



provenance

/'prəv(ə)nəns/

*noun*

the place of origin or earliest known history of something.

"an orange rug of Iranian provenance"

Similar:

origin

source

place of origin

birthplace

spring

wellspring

fount



- the beginning of something's existence; something's origin.

"they try to understand the whole universe, its provenance and fate"

- a record of ownership of a work of art or an antique, used as a guide to authenticity or quality.

plural noun: **provenances**

"the manuscript has a distinguished provenance"

# Data Provenance in Data Science



Lineage



Traceability



Compliance



Dependency Analysis

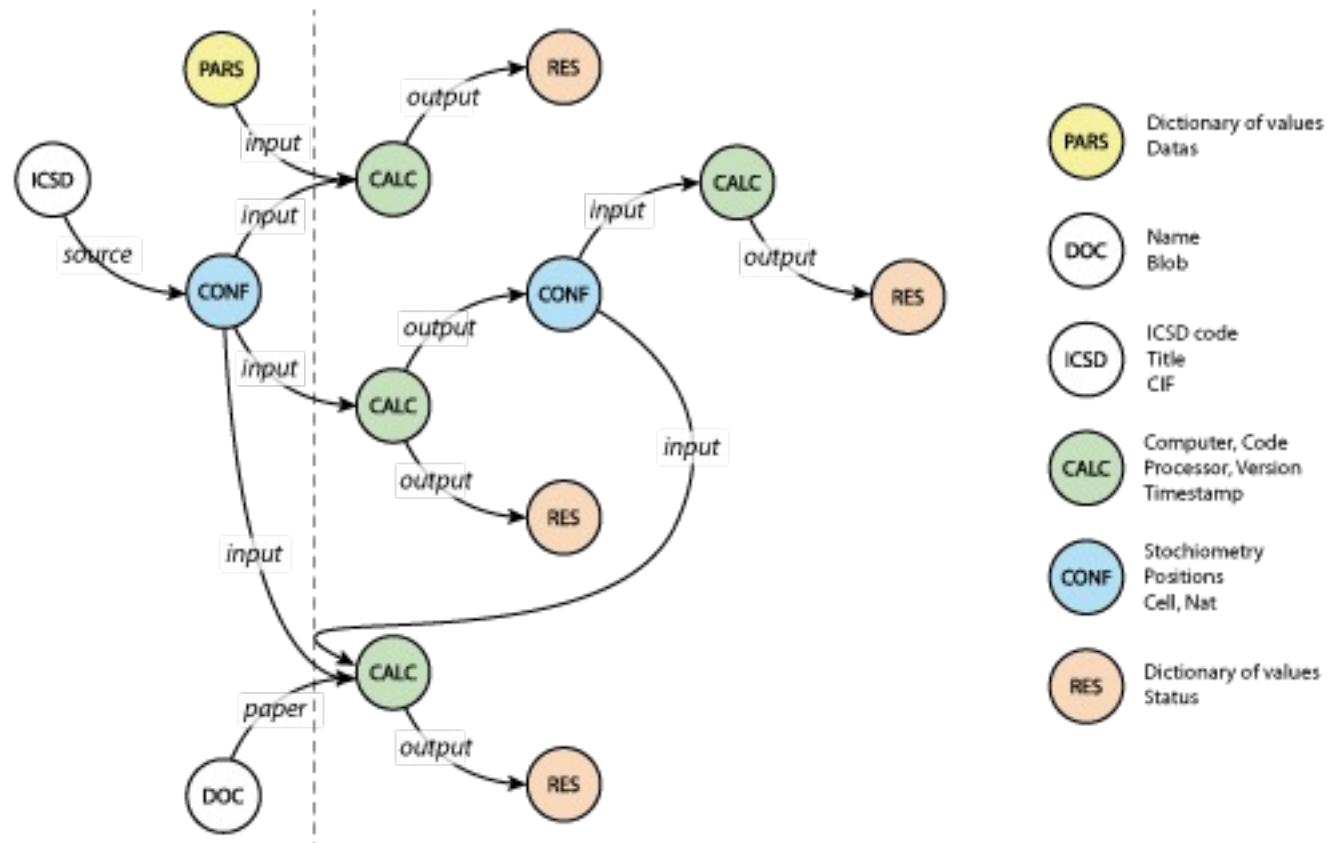


Data and AI Reproducibility

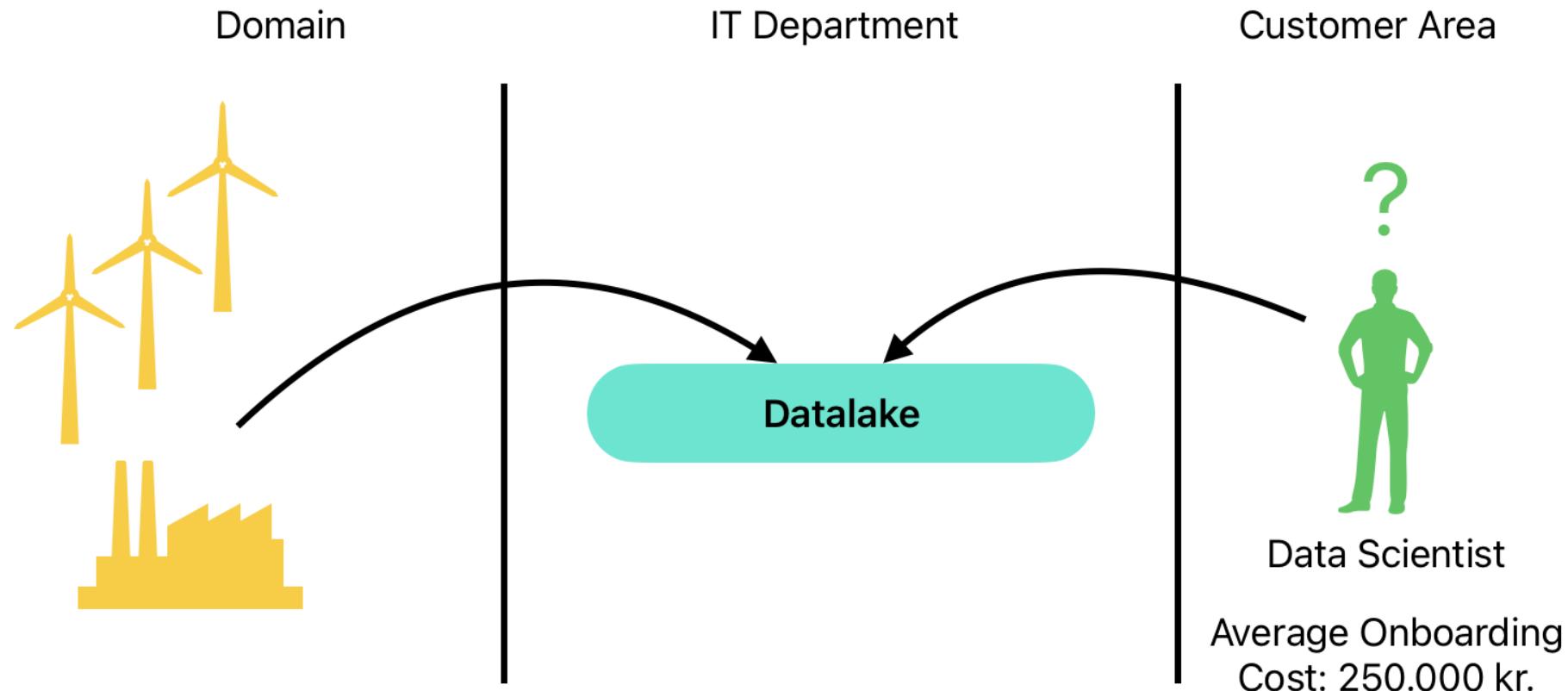


AI Explainability

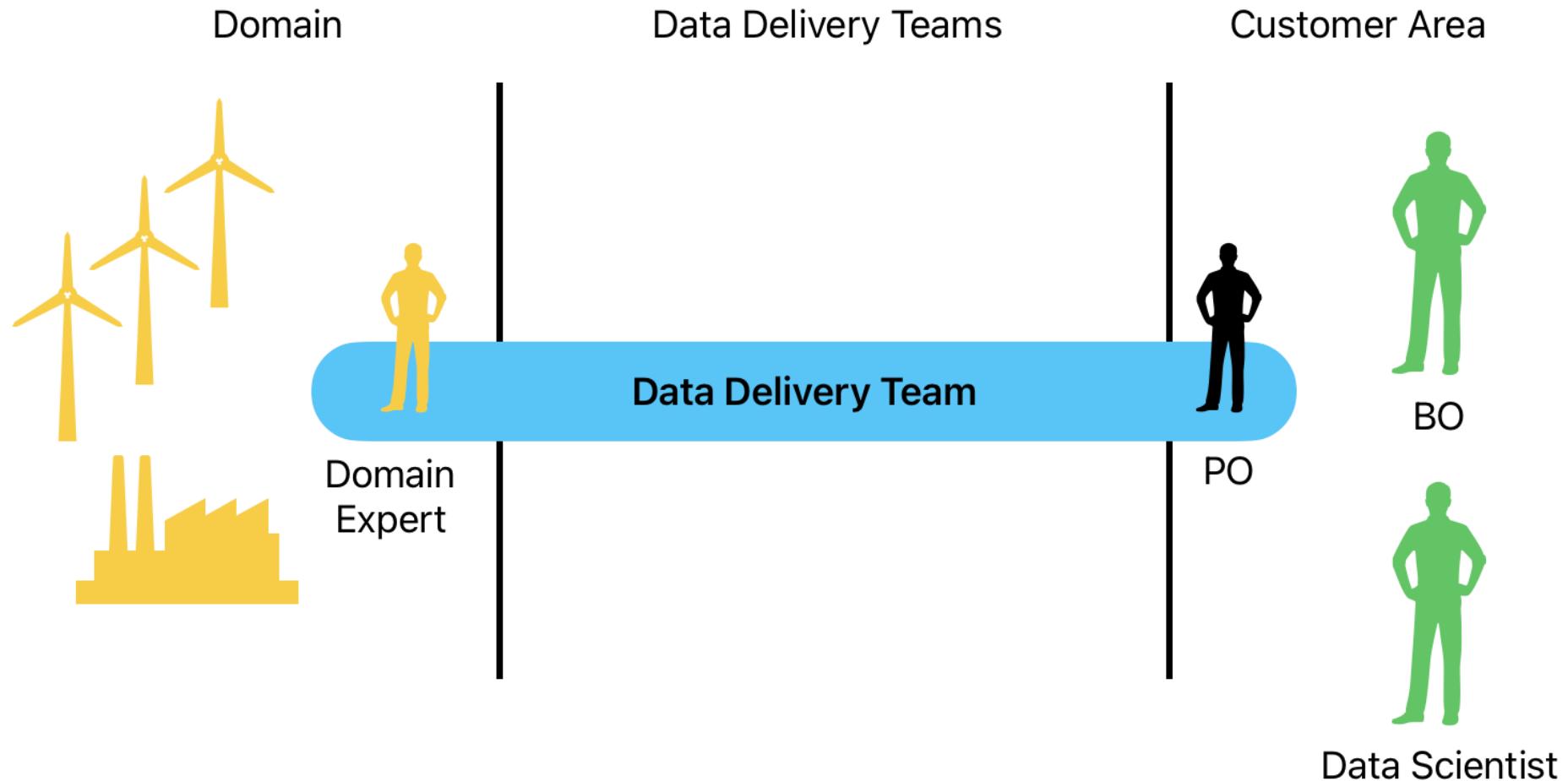
# Lineage, Tracability, Reproducibility



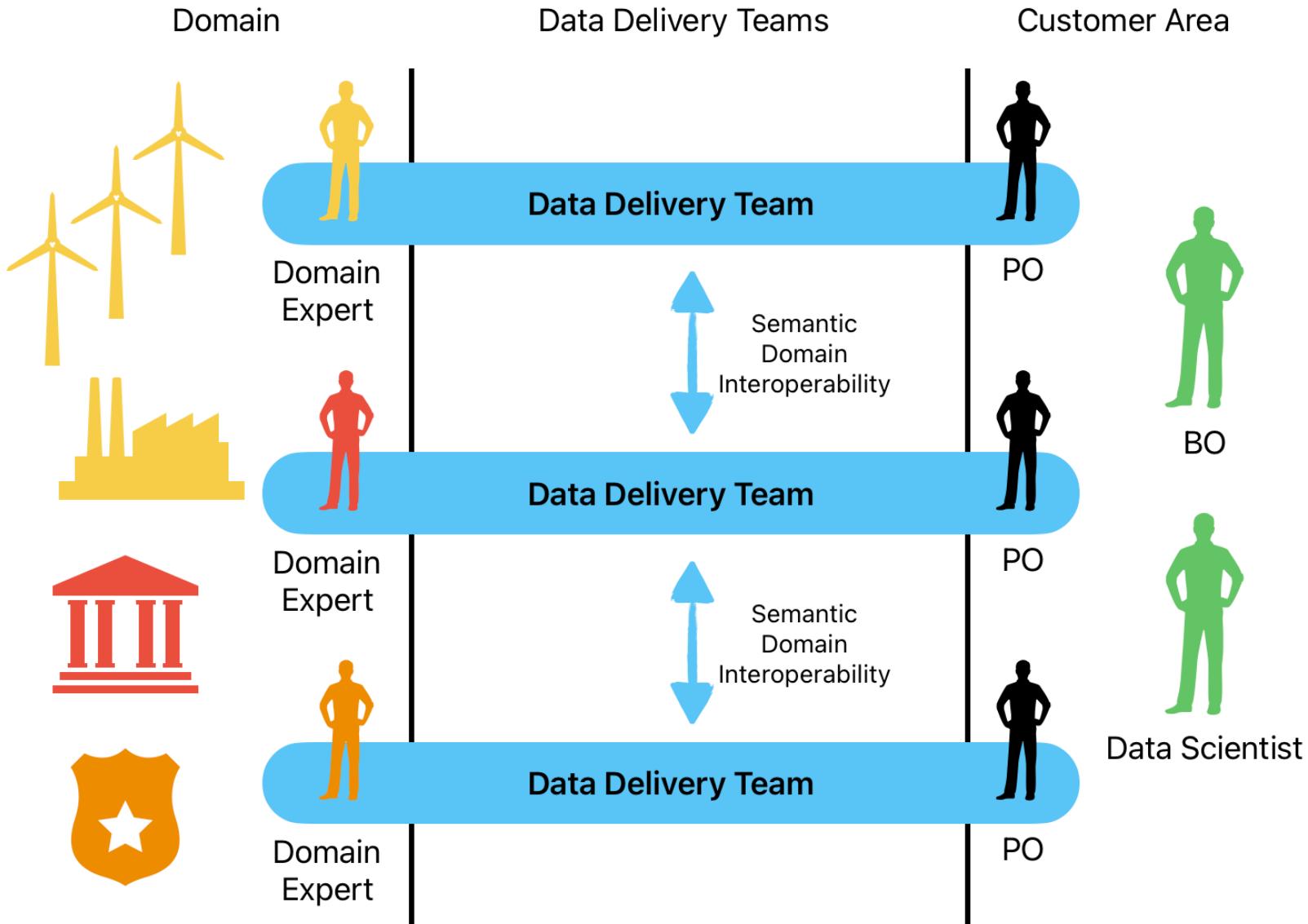
# Current State – Domain Knowledge is Impossible to access



# Data Delivery Teams Create Access to Domain Knowledge



# And Allows Us to Bridge Domains



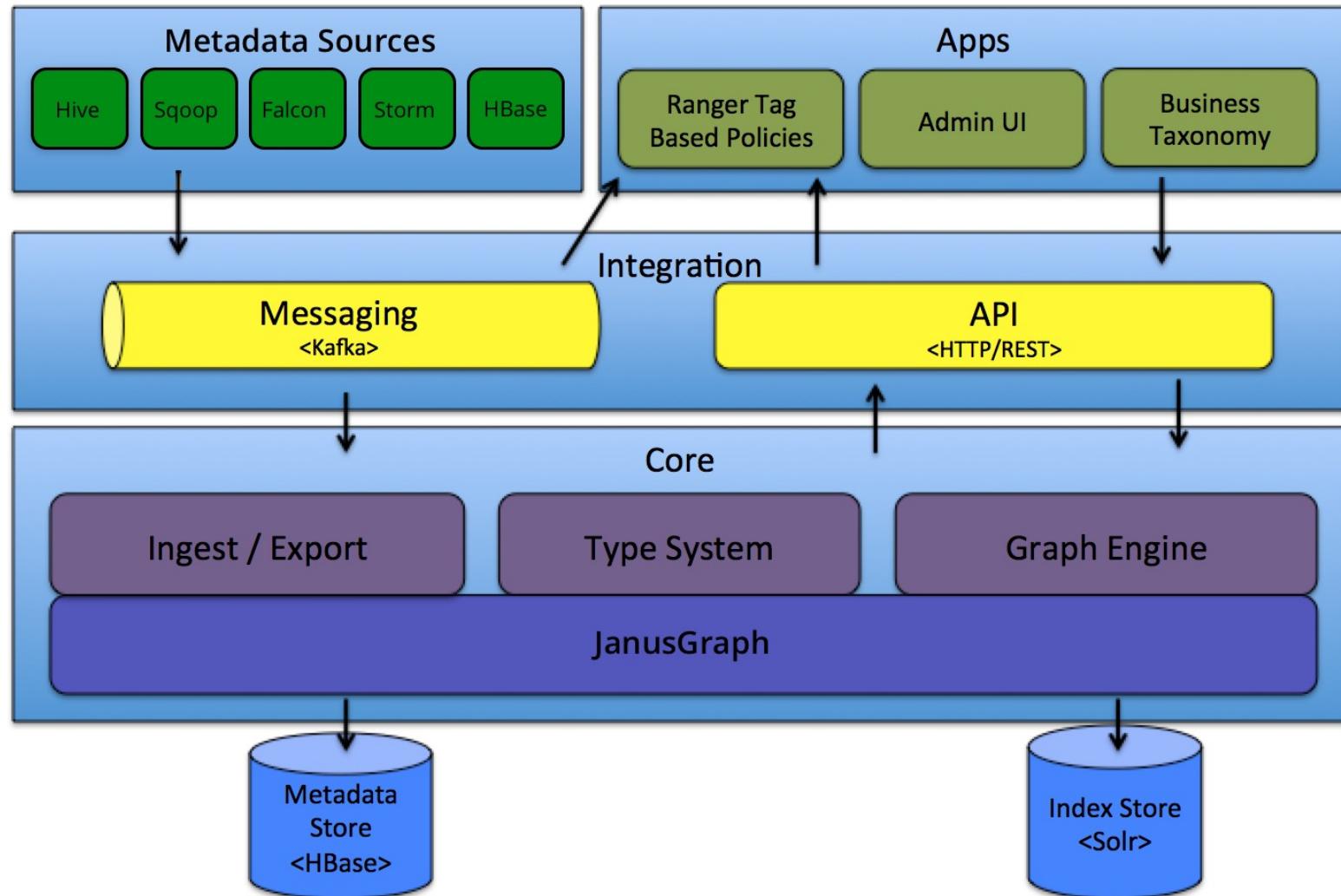
# Metadata Implementations: Apache Atlas



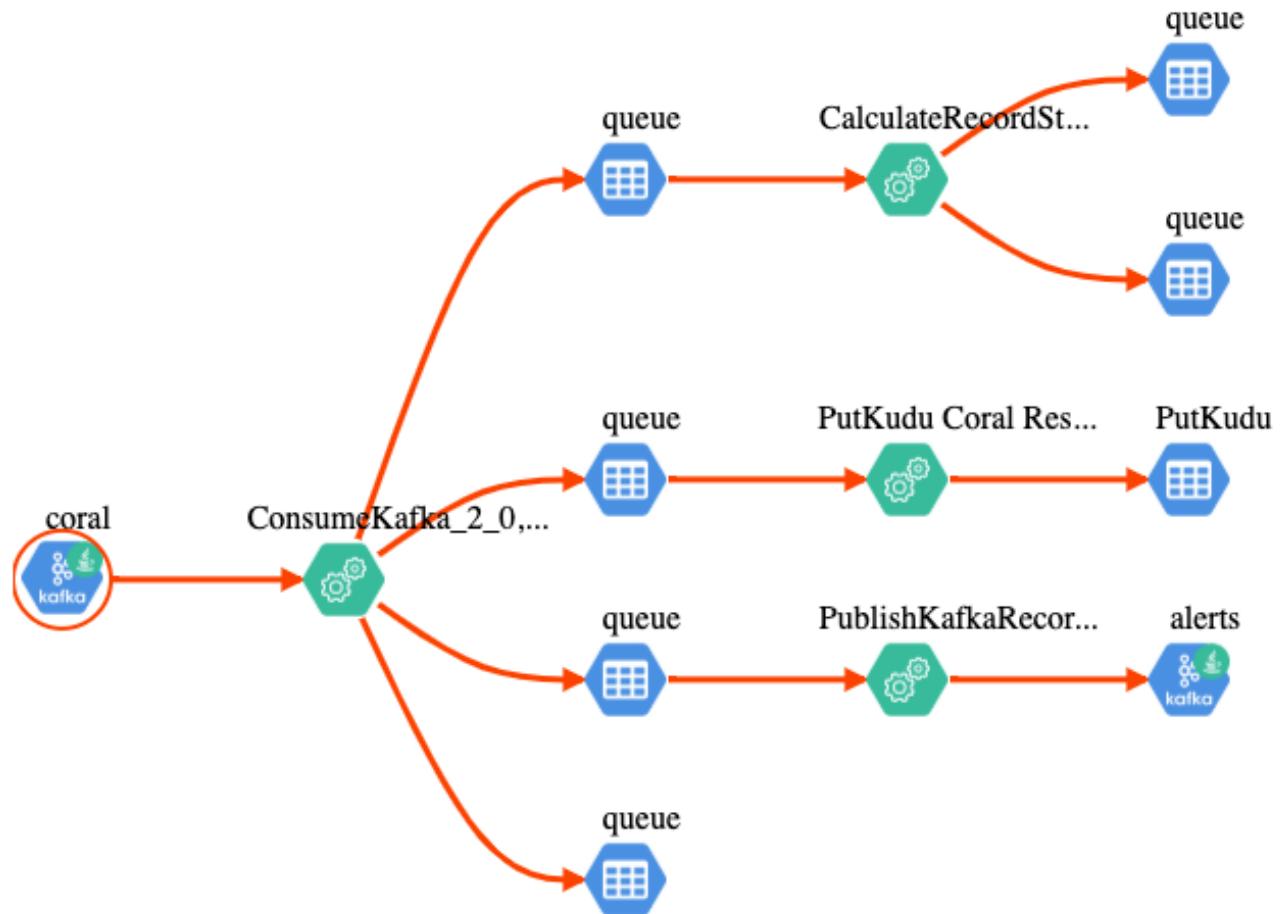
# Apache Atlas

- Features
  - Types (Schema)
  - Classification
  - Lineage
  - Search/Discovery
  - Governance (Security and Data Masking)
- Older metadata alternative
- Creates semantic connections between terms
- Relatively simple compared to newer alternatives
- Integrates well with the Apache ecosystem

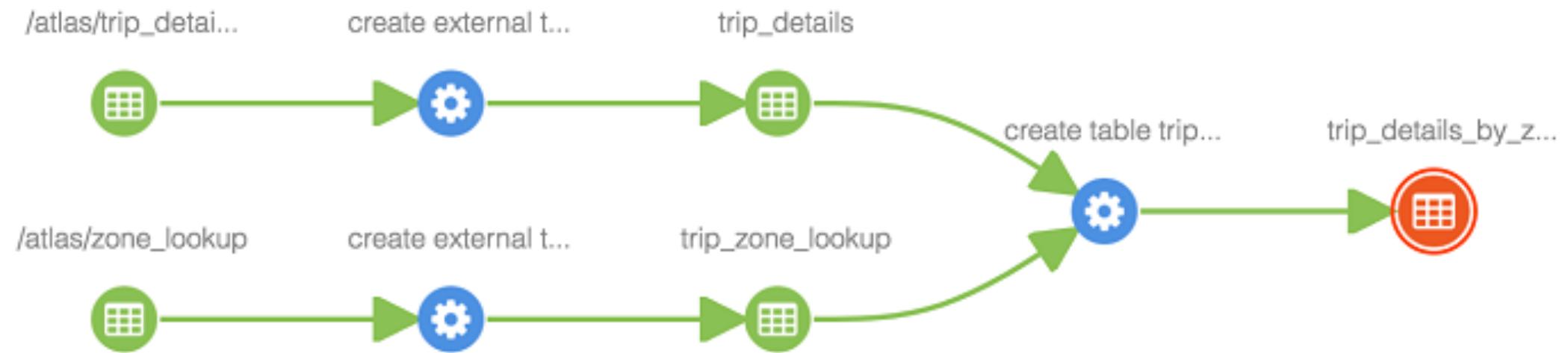
# Atlas Architecture



# Atlas Lineage – Realtime Data (Kafka)



# Atlas Lineage – Historical Data (Hive)

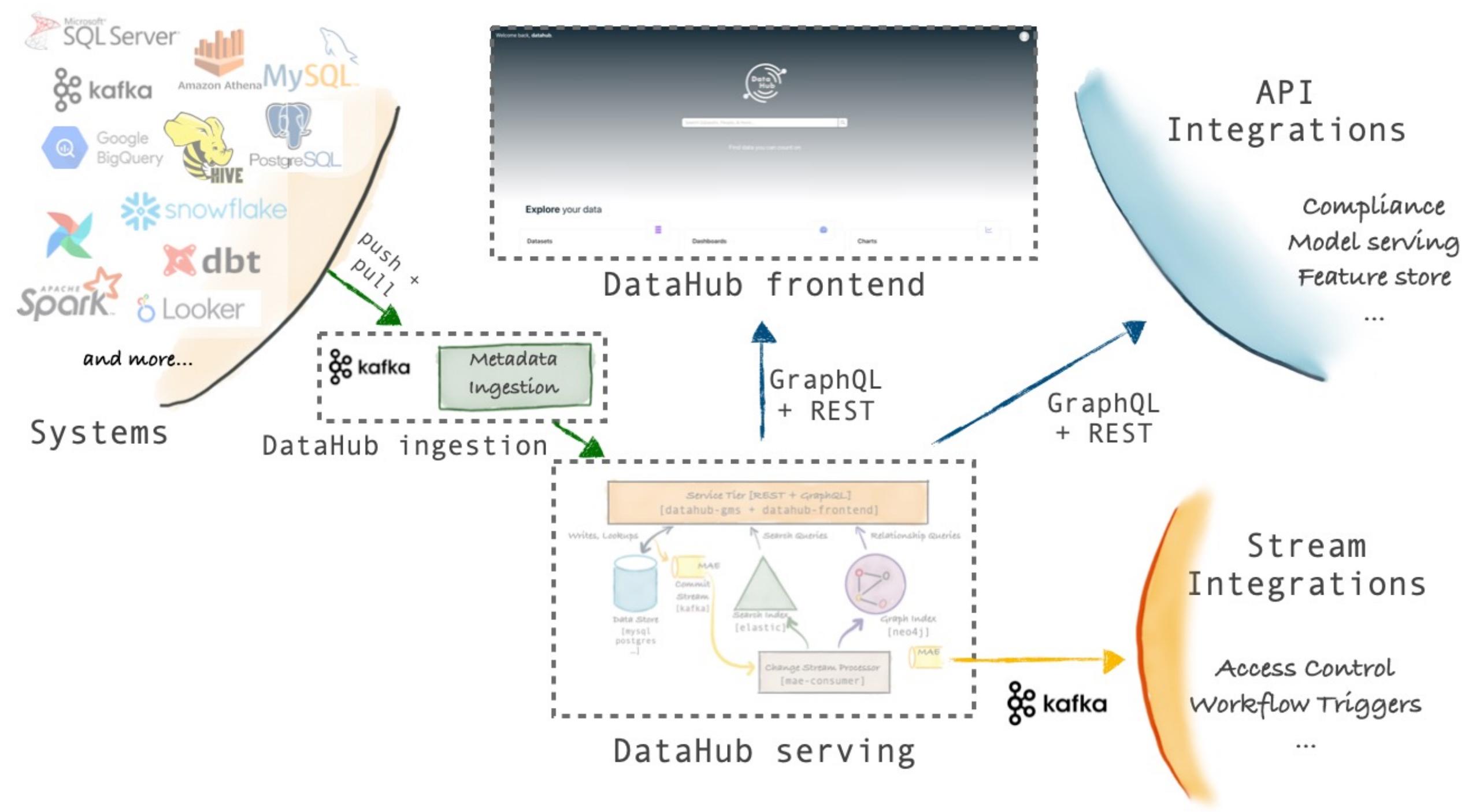


# **Metadata Implementations: LinkedIn Datahub**



## LinkedIn Datahub

- Features
  - Search and Discovery
  - Access Control
  - Data Lineage
  - Compliance
  - Data management and automated metadata generation
  - AI Explainability and reproducibility
- Extremely extensive compatibility with databases and data lakes
- Not bound to a specific ecosystem





## Metadata Store

Metadata storage, indexing, serving



## Metadata Models

Entities, Aspects, & Relationships



## Ingestion Framework



## GraphQL API

Metadata Graph Queries & Mutations

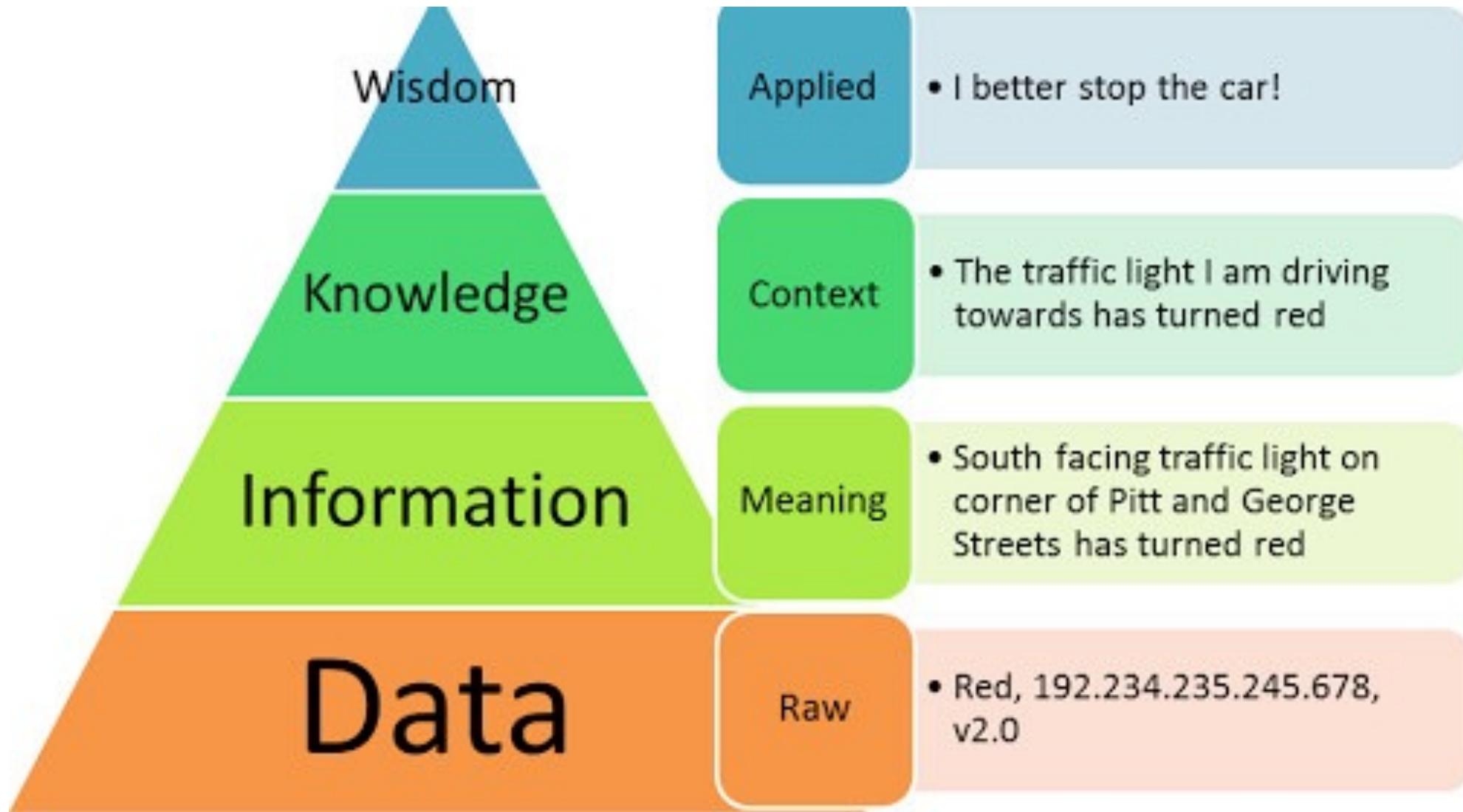


## User Interface

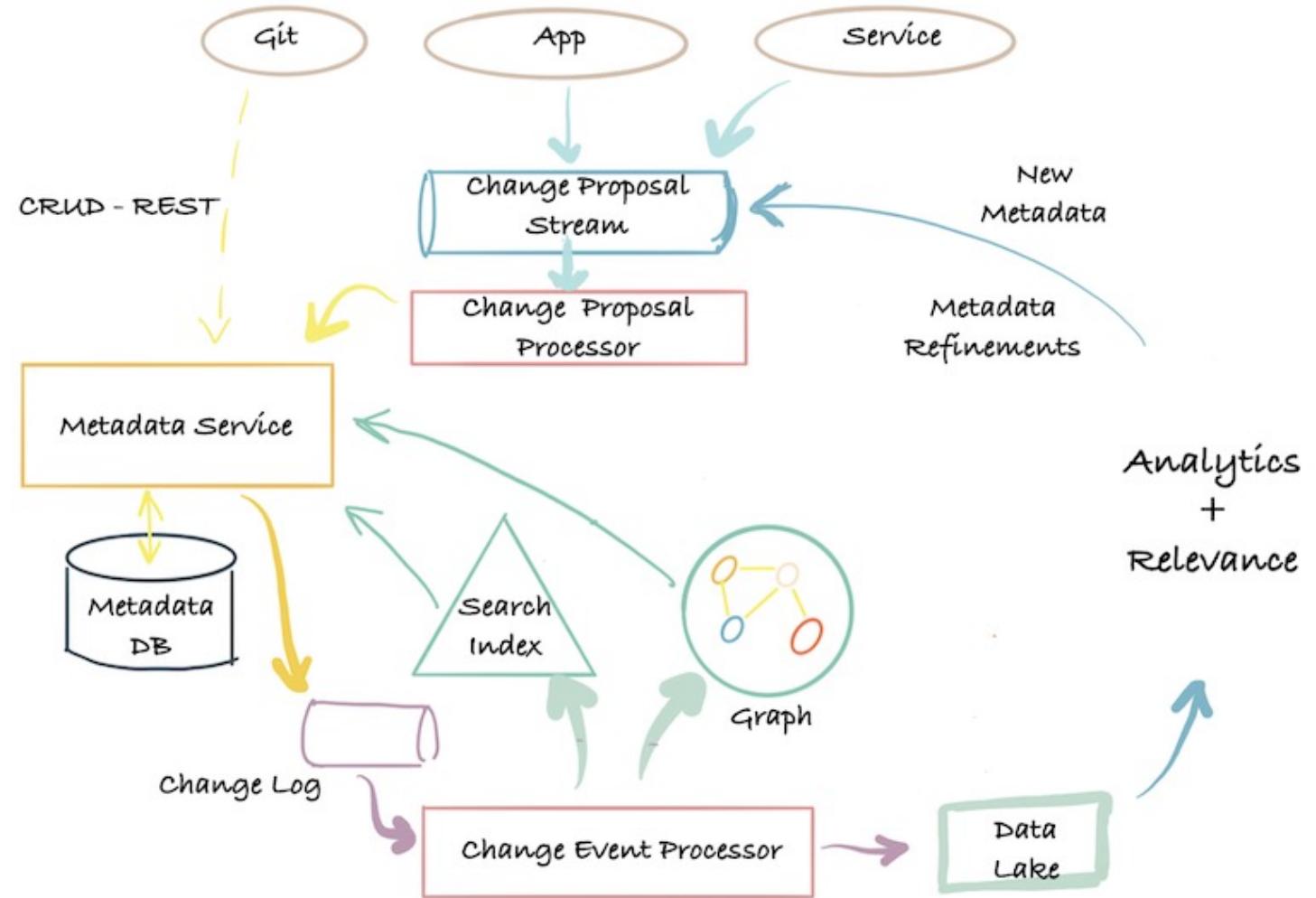
Discovery, Governance, Observability

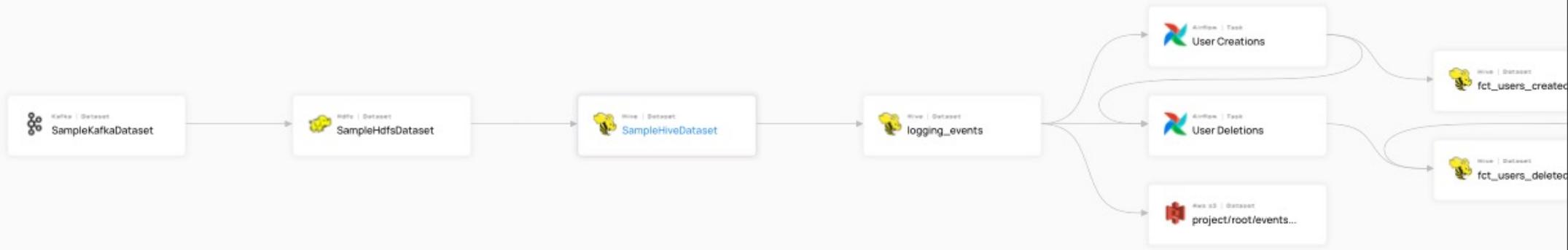


## What Meta Data limitations are there?



# Dynamic Metadata Changes





Type ▲ ▼ Name ▲ ▼

Last Status ▲ ▼



BigQuery



Redshift



Snowflake



Kafka



Looker



LookML



Tableau



MySQL



Postgres



MongoDB



Azure AD

okta

Okta



Glue



Oracle



Hive



Superset



Athena



SQL Server



ClickHouse



Trino



Druid



Metabase



MariaDB



Power BI



Mode



Custom

 Hive

View all results for Hive

Datasets

SampleHiveDataset

### Domains



### Explore your Metadata

Datasets

7

Dashboards

1

Charts

2

Pipelines

1

Glossary Terms

5

Feature Tables

5

ML Models

1

### Platforms



Hive

6



Feast

5



Looker

3



Airflow

3



Kafka

1



HDFS

1



AWS S3

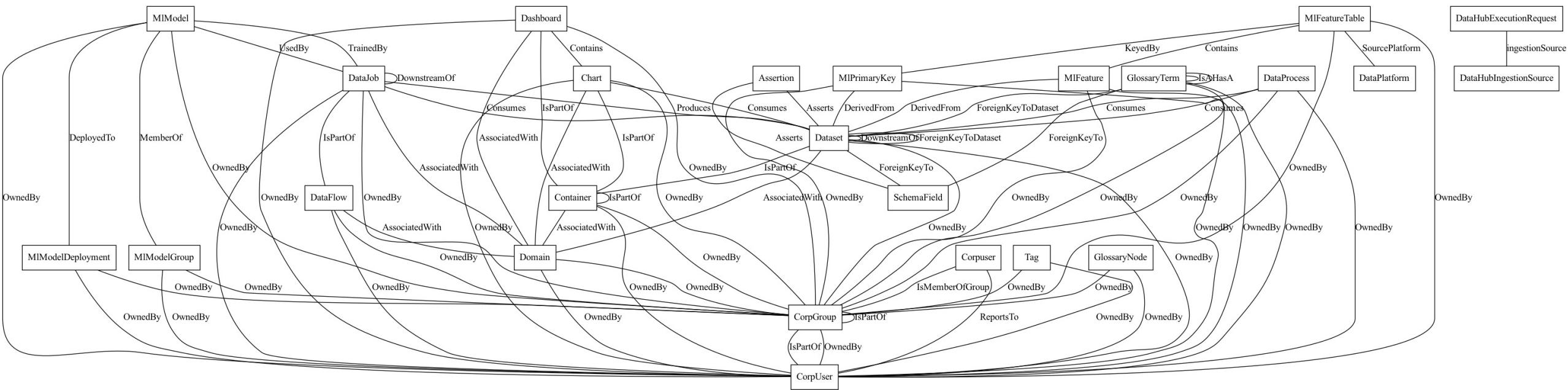
1

# The Core Entities

DataHub's "core" Entity types model the Data Assets that comprise the Modern Data Stack. They include

1. **Data Platform**: A type of Data "Platform". That is, an external system that is involved in processing, storing, or visualizing Data Assets. Examples include MySQL, Snowflake, Redshift, and S3.
2. **Dataset**: A collection of data. Tables, Views, Streams, Document Collections, and Files are all modeled as "Datasets" on DataHub. Datasets can have tags, owners, links, glossary terms, and descriptions attached to them. They can also have specific sub-types, such as "View", "Collection", "Stream", "Explore", and more. Examples include Postgres Tables, MongoDB Collections, or S3 files.
3. **Chart**: A single data visualization derived from a Dataset. A single Chart can be a part of multiple Dashboards. Charts can have tags, owners, links, glossary terms, and descriptions attached to them. Examples include a Superset or Looker Chart.
4. **Dashboard**: A collection of Charts for visualization. Dashboards can have tags, owners, links, glossary terms, and descriptions attached to them. Examples include a Superset or Mode Dashboard.
5. **Data Job** (Task) : An executable job that processes data assets, where "processing" implies consuming data, producing data, or both. Data Jobs can have tags, owners, links, glossary terms, and descriptions attached to them. They must belong to a single Data Flow. Examples include an Airflow Task.
6. **Data Flow** (Pipeline) : An executable collection of Data Jobs with dependencies among them, or a DAG. Data Jobs can have tags, owners, links, glossary terms, and descriptions attached to them. Examples include an Airflow DAG.

See the [Metadata Modeling/Entities](#) section on the left to explore the entire model.



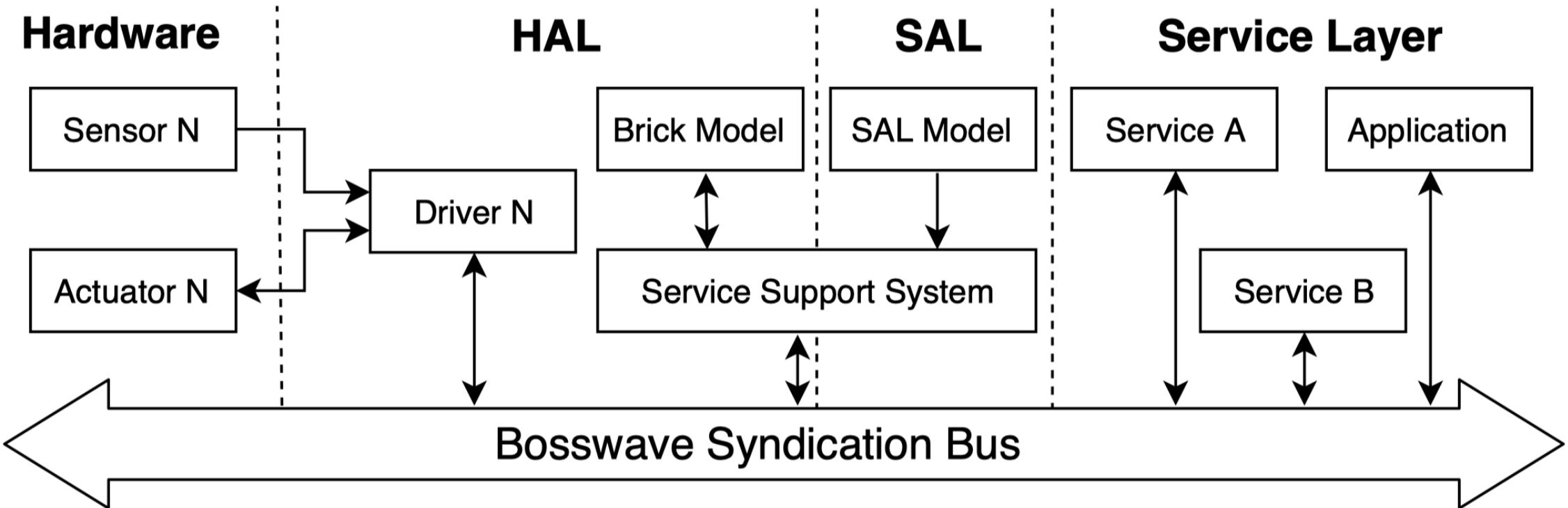
**Atlas, Datahub, and others, only uses graphs for  
Lineage, Dependencies, Ownership and Data  
Flow, not to argument the metadata context itself.**

**For a short explanation of DataHub, watch this:**

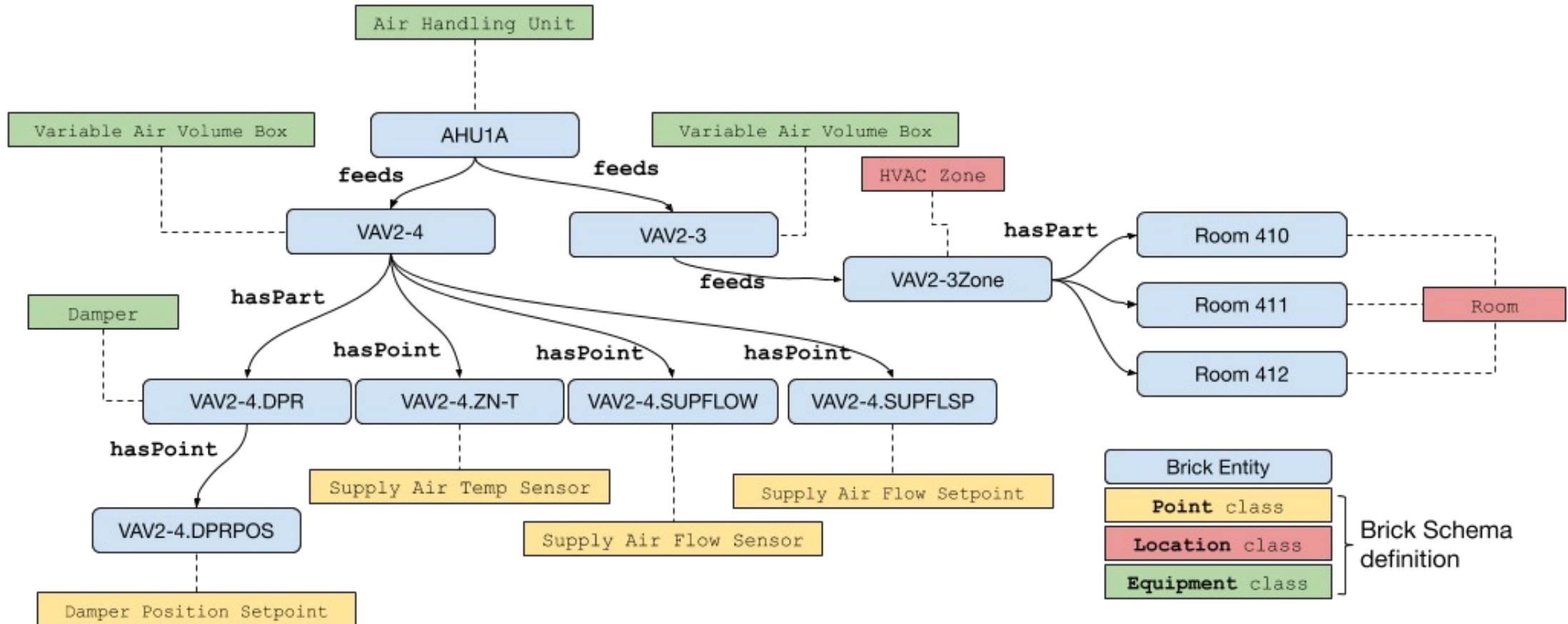
**<https://www.youtube.com/watch?v=VY57iRdG-Us>**

# **Metadata Research: Brick and SAL (Knowledge Graphs)**

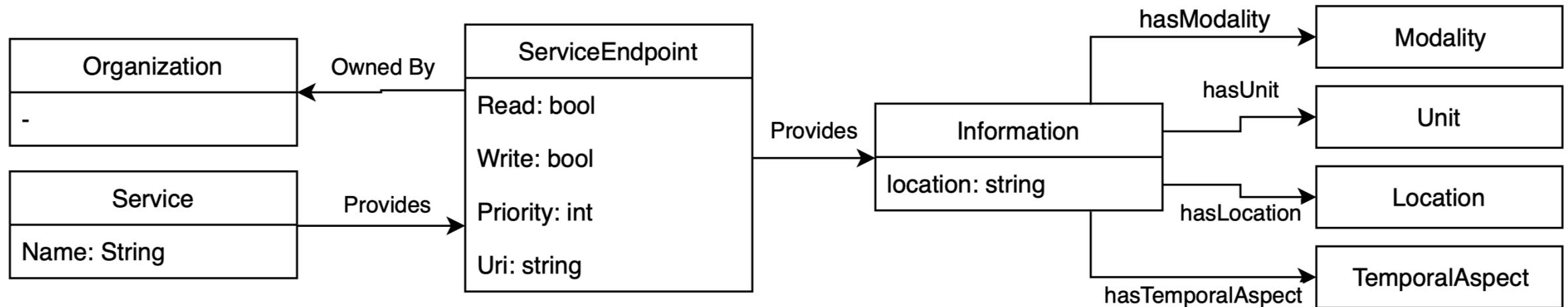
# Why are knowledge graphs important?



# Brick instantiation example



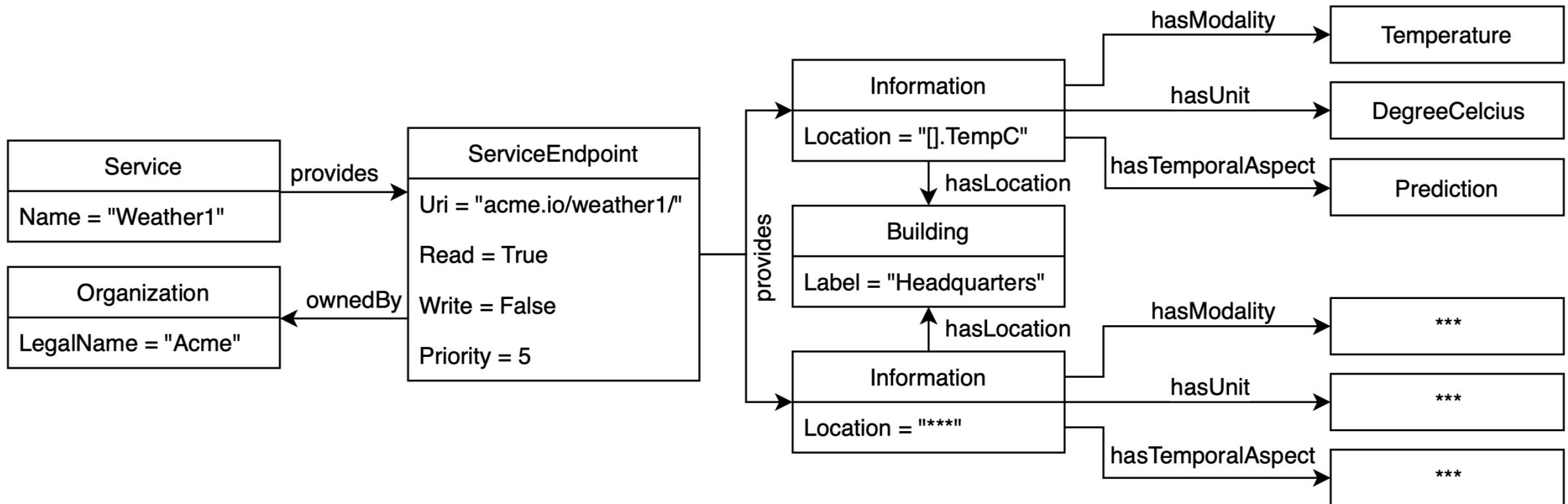
# SAL Ontology Classes



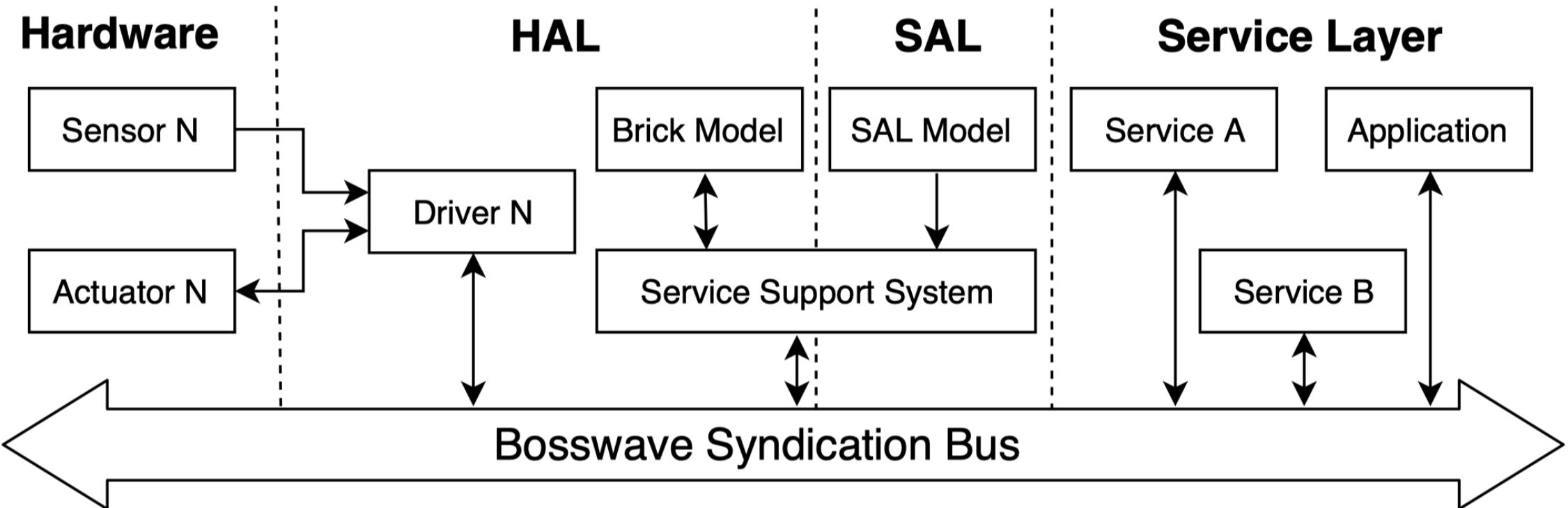
# SAL Types

Modality	Unit	TemporalAspect
Angle, CO <sub>2</sub> , Presence, Flow, Illuminance, Power, Pressure, Rain, Humidity, Temperature, Wind, AbsoluteTime, RelativeTime, PowerFlexibility, Performance, Energy, Certainty, Time	Boolean, Count, CubicMeters, CubicMetersPerHour, DegreeCelsius, DegreeFahrenheit, Degrees, GigaByte, KiloByte, Hours, Hertz, Joules, JoulesPerCubicMeter, Kelvin, KiloJoulesPerSquareMeter, KiloMeters, KiloWatts, KiloWattHours, Lux, CubicMeters, CubicMetersPerHour, CubicMetersPerSecond, MilliAmperes, Minutes, MilliMeters, MilliSeconds, MetersPerSecond, MilliVolts, MilliWatts, MilliWattHours, Pascal, Percent, PartsPerMillion, RotationsPerMinute, Volts, Watts, Unitless, Time, DateTime, Date	Prediction RealTime Archival

# SAL Instantiation Example



# How are the knowledge graphs applied?



# Looking to the future: Data Mesh

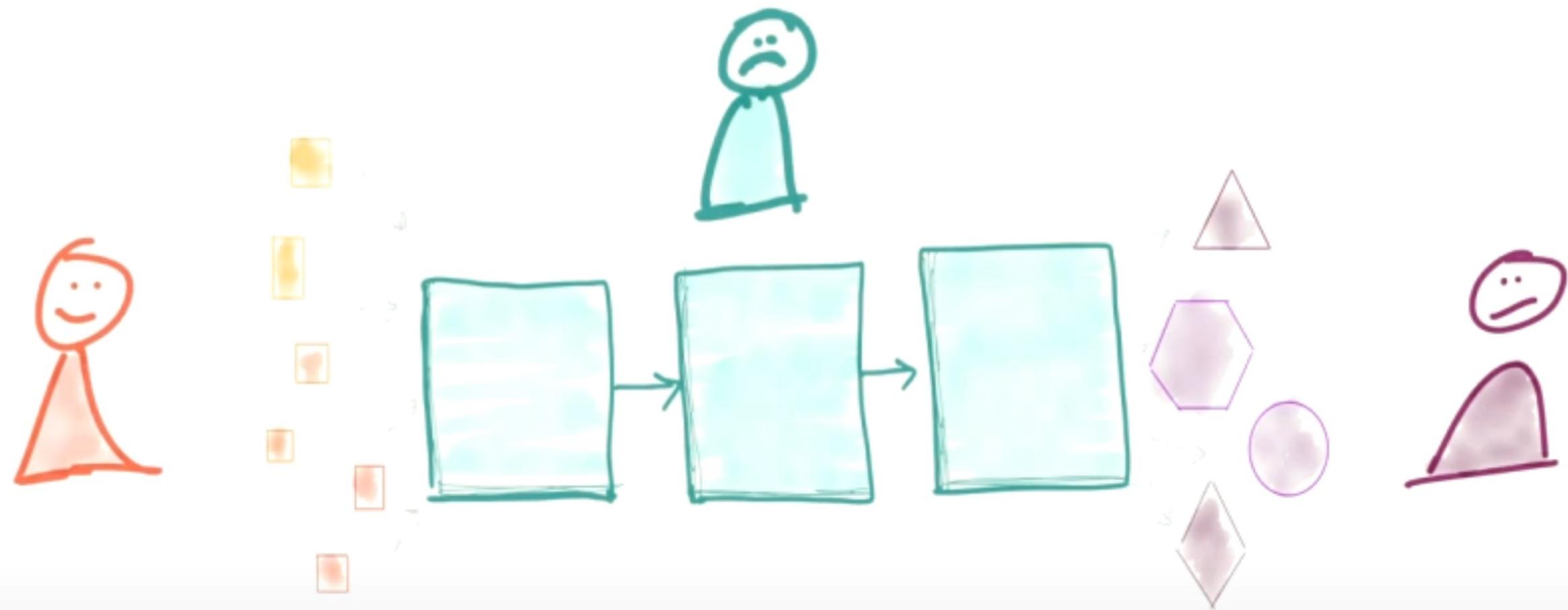
**This section will give you a quick look into the future, and ensure you have a concept of what Data Mesh means if you are confronted with the term in the industry.**

**Slides on loan from Zhamak Dehghani**

**For a comprehensive walkthrough, watch:**

**<https://www.youtube.com/watch?v=52MCFe4v0UU>**

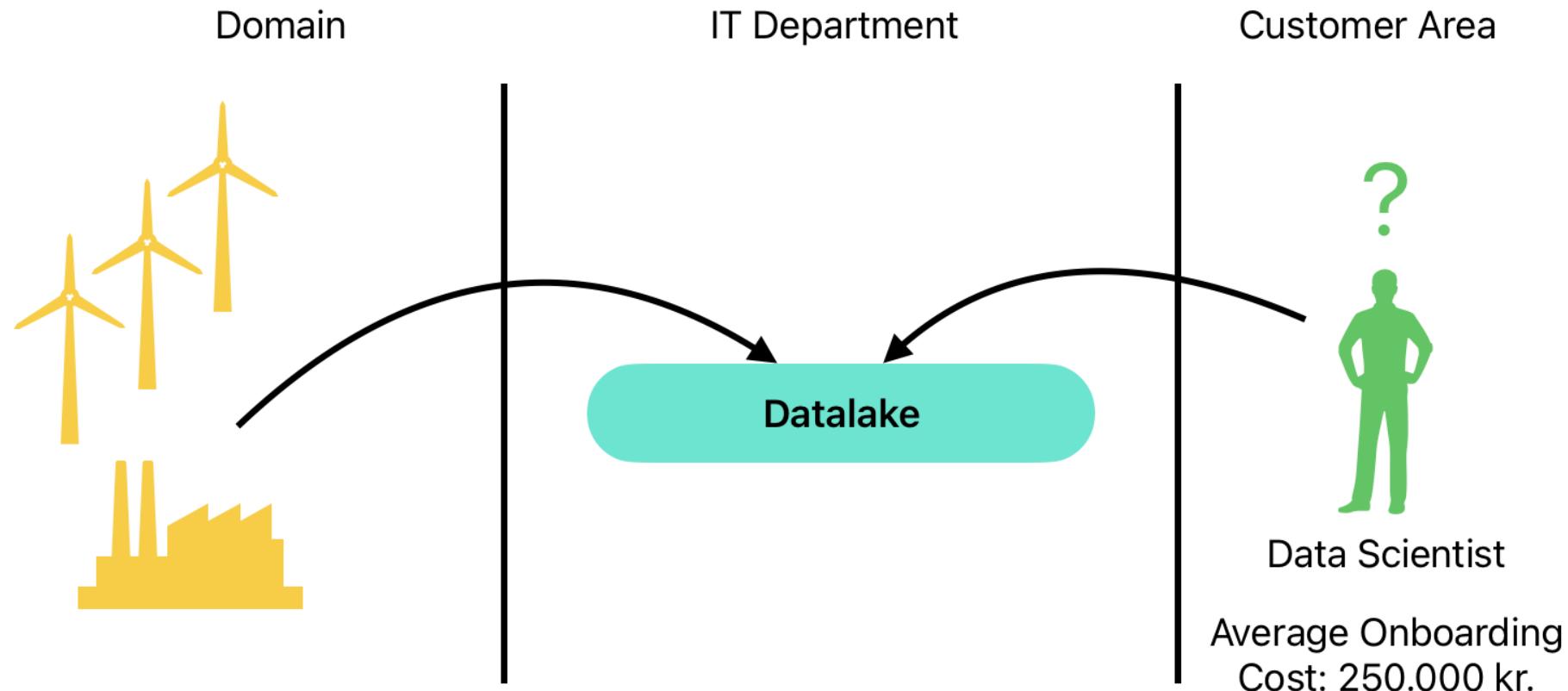
## Data Platform Engineers



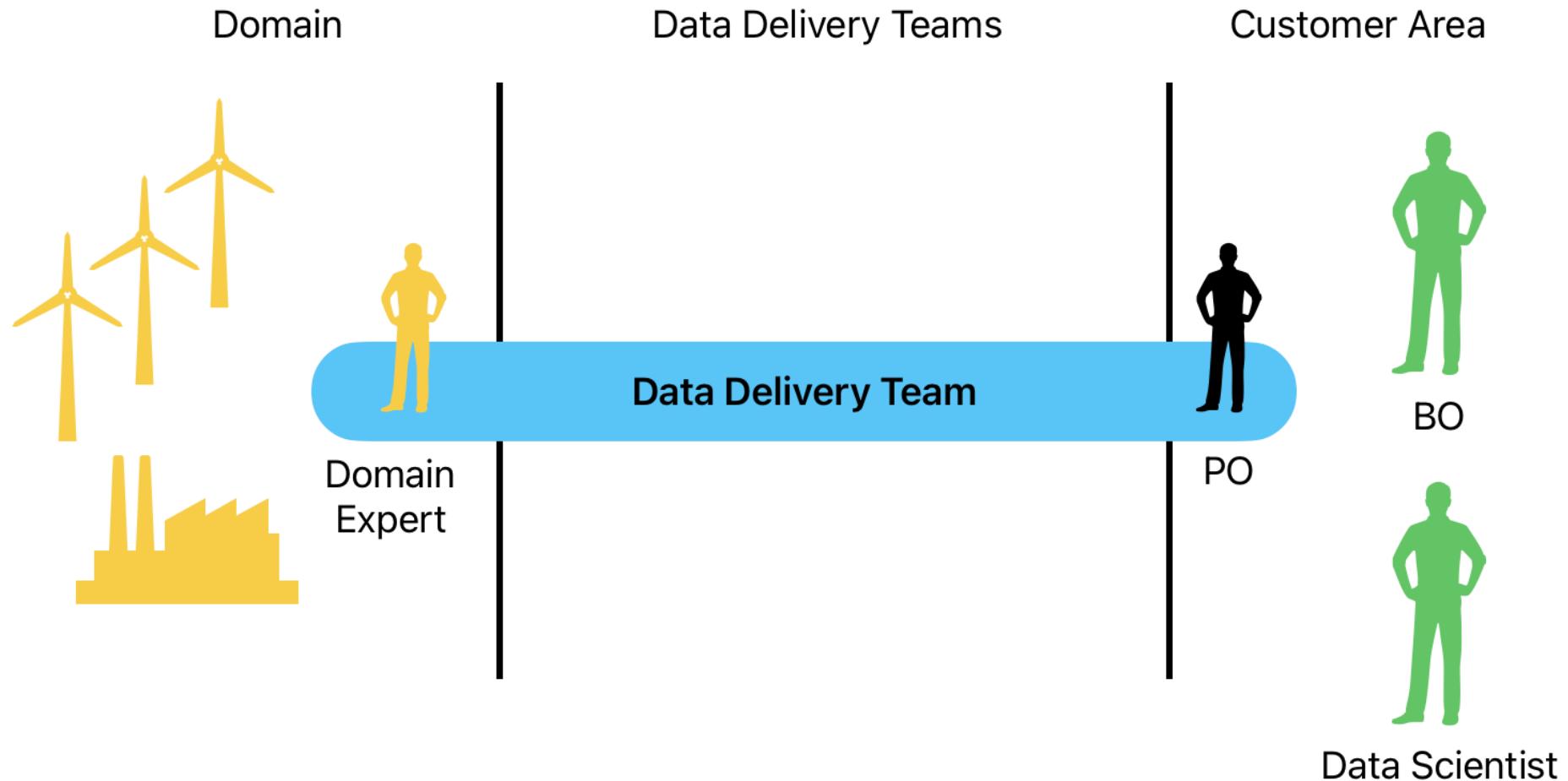
Domains' Operational systems Teams

Data Scientists, BI teams, ...

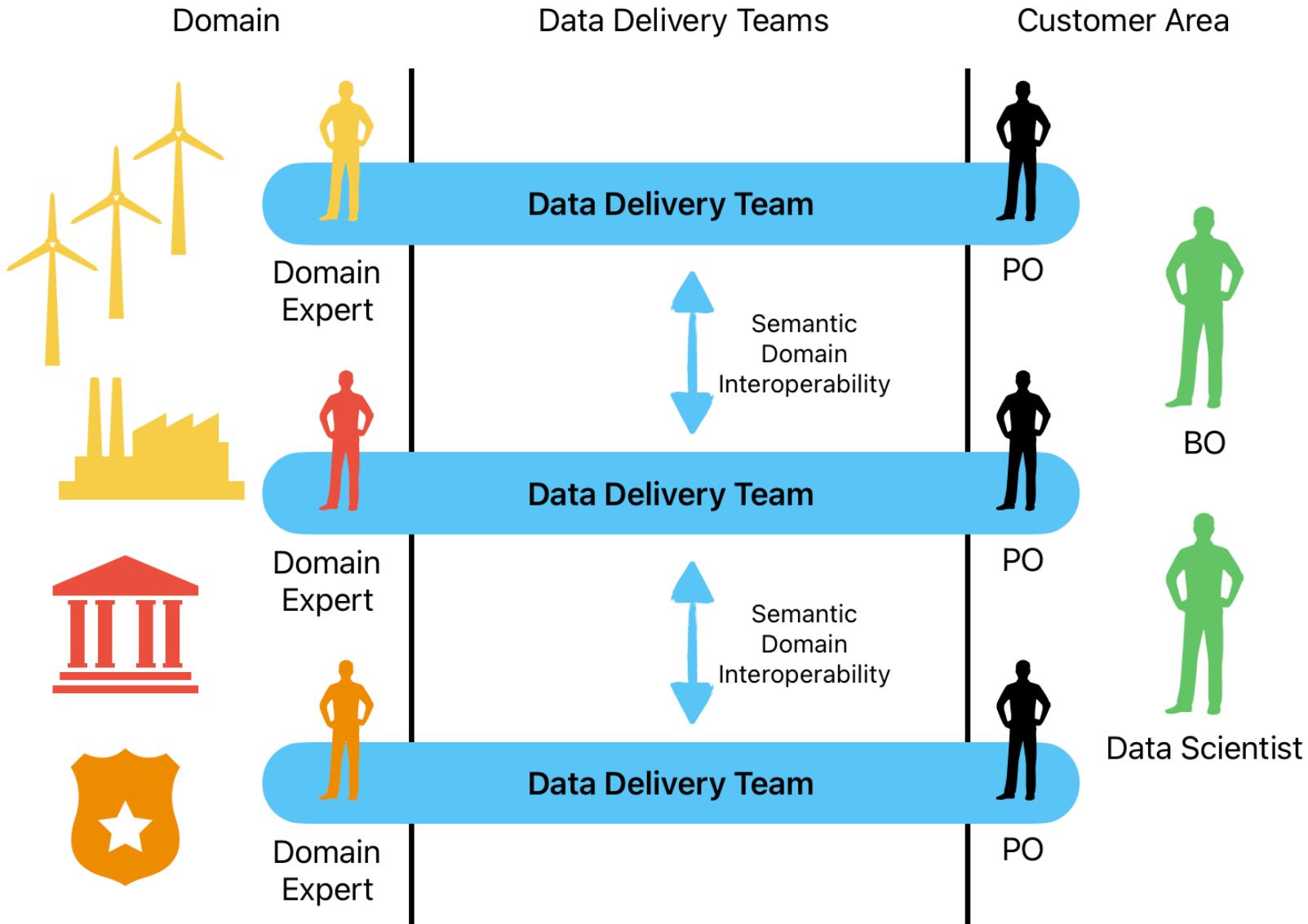
# Current State – Domain Knowledge is Impossible to access



# Data Delivery Teams Create Access to Domain Knowledge



# And Allows Us to Bridge Domains



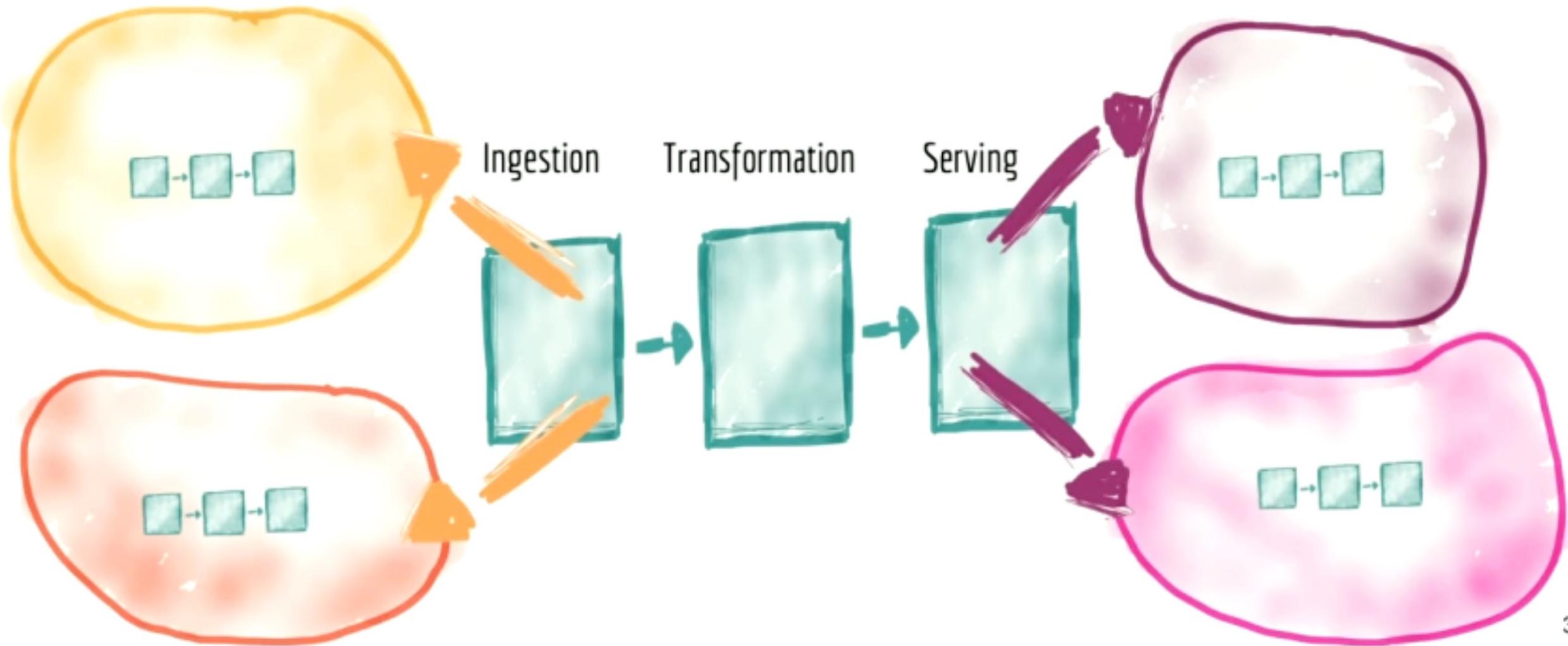
# DATA MESH PRINCIPLES



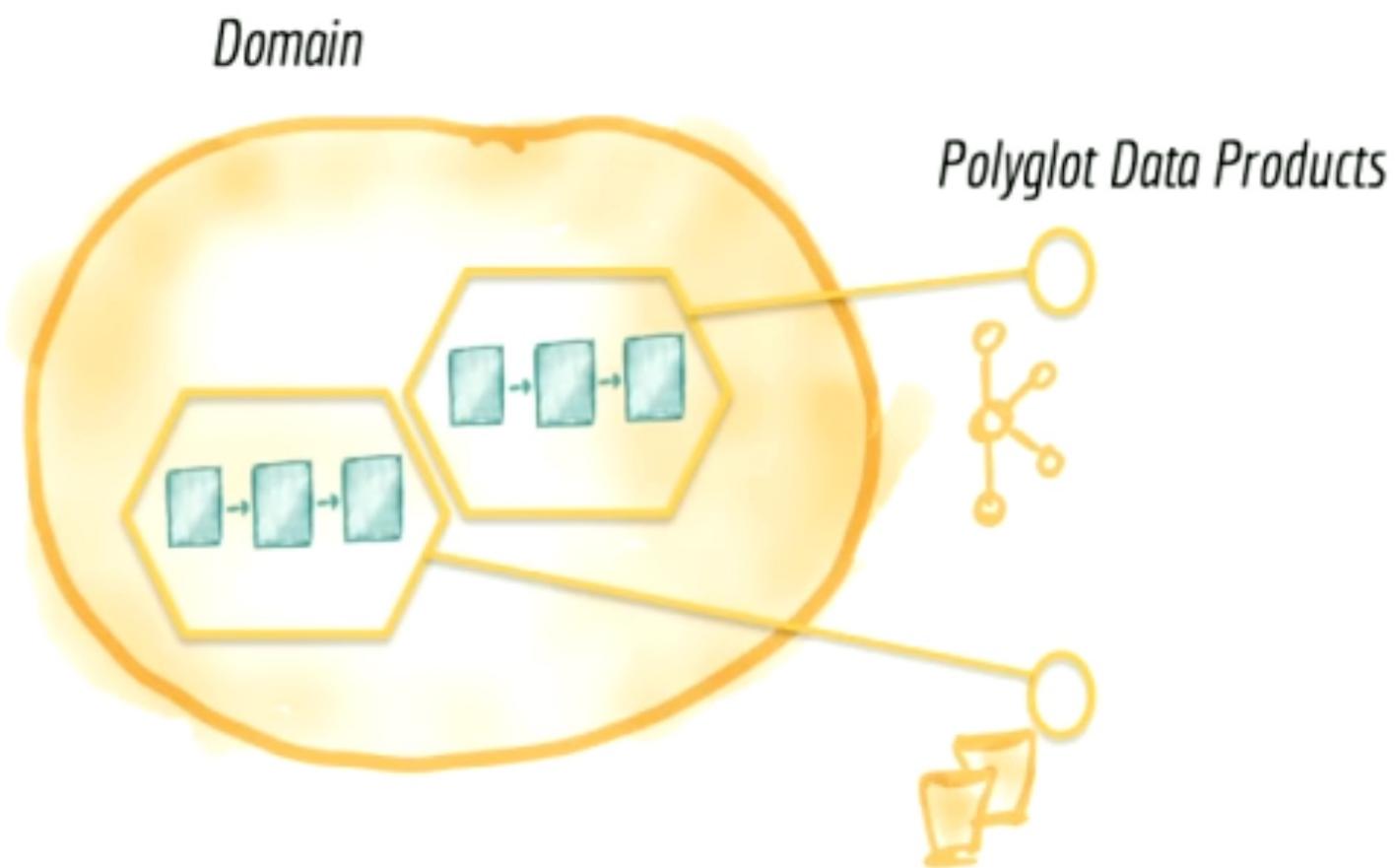
# DISTRIBUTED PIPELINES IN DOMAINS

*More cleansing, integrity checks here*

*More aggregations, ML modelling here*



# DOMAIN DATA AS A PRODUCT



# DOMAIN DATA AS A PRODUCT

Aka Data Products



SHARED | DISCOVERABLE



SELF-DESCRIBING



ADDRESSABLE



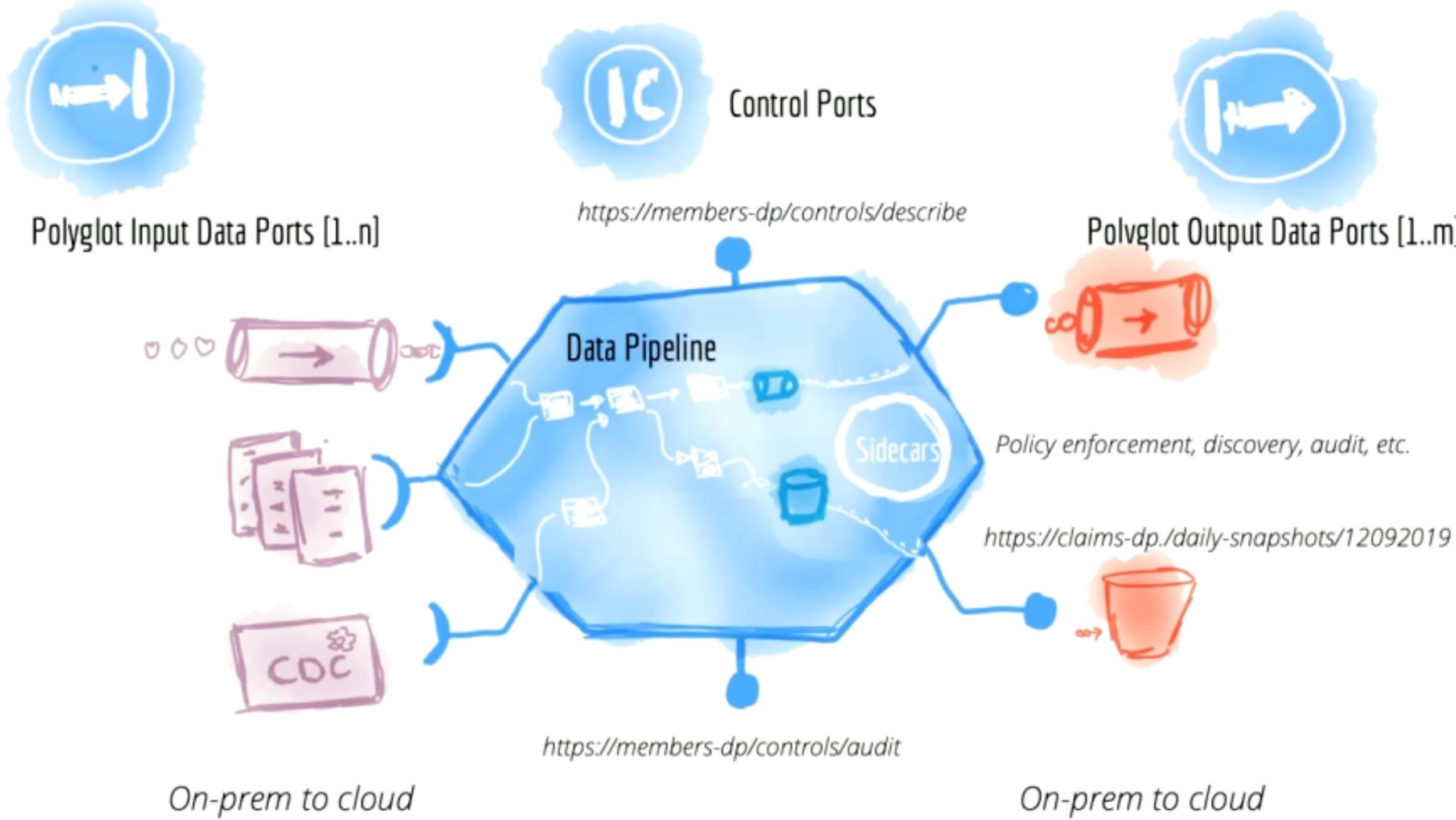
INTER OPERABLE  
(GOVERNED  
BY GLOBAL STANDARDS)



TRUSTWORTHY  
(DEFINED & MONITORED SLOs)



SECURE  
(GOVERNED  
BY GLOBAL ACCESS CONTROL)



# DATA MESH PARADIGM SHIFT

**FROM**

Centralized ownership

Monolithic

Pipeline first class concern

Data as a by-product

Siloed data engineering team

**TO**

Decentralized ownership

Distributed

Domain data first class concern

Data as a product

Cross-functional data domain teams

# ADOPT A NEW LANGUAGE

**FROM**

**TO**

Ingesting

Serving

Extracting & Loading

Discovering & Consuming

Flowing data through centralized Pipelines

Publishing output data ports

Centralized Data Lake | Warehouse | Platform

Ecosystem of Data Products

# Exercises!