

Utilizing Data Lakes for Managing Multidisciplinary Research Data

Mark Greiner

Max-Planck Institute for Chemical Energy Conversion

2022-06-27

Problem Space

Description of domain

Problem statement

Presentation
Outline

Data Governance

Architectures

Technologies

Future directions

Solution Space

1. Description of the Domain



Problem space

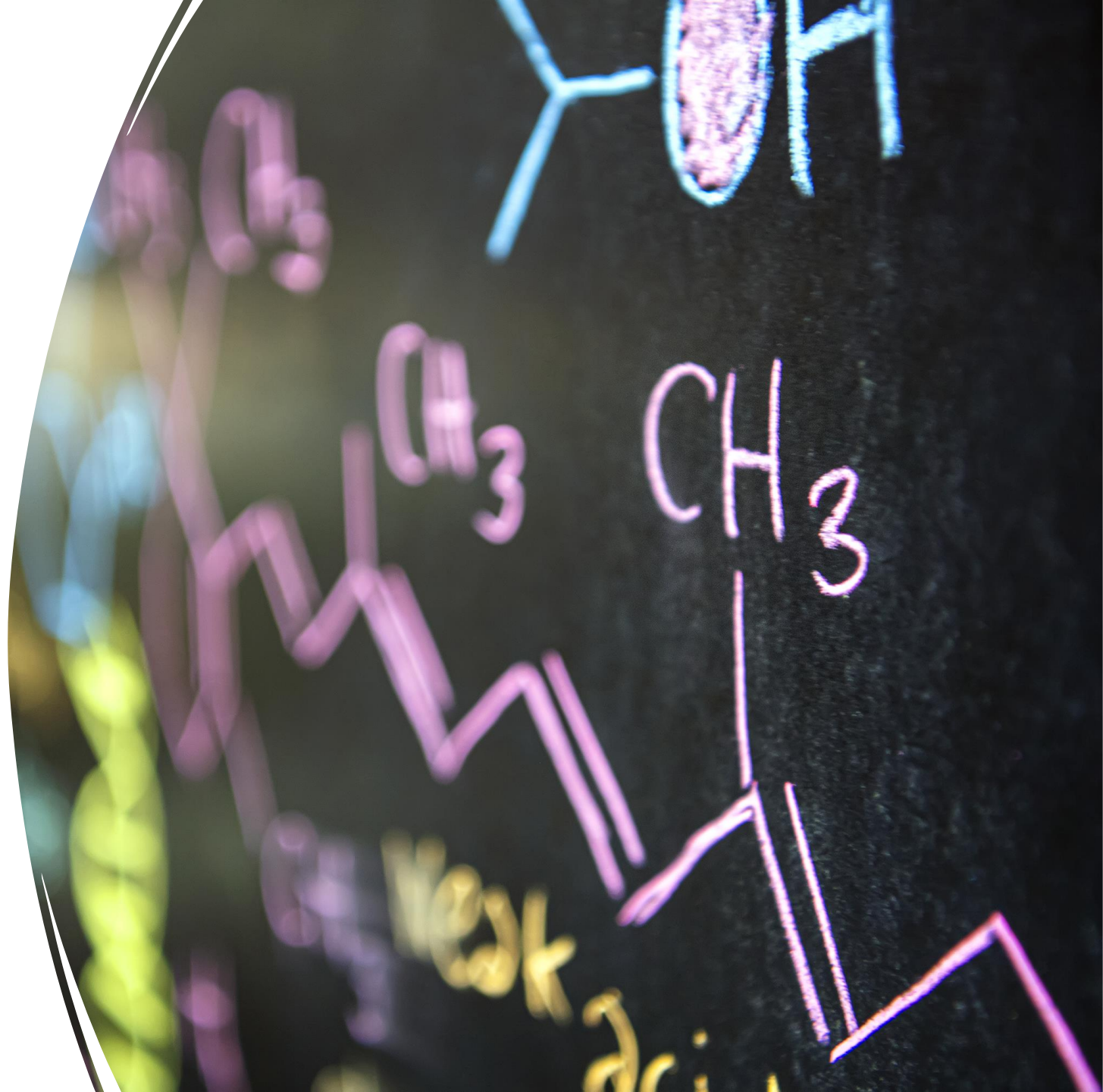
About MPI-CEC

Research Discipline

- Catalysis for chemical energy
- Water splitting, bio-catalysts, chemical production

Magnitude

- 218 Researchers
- 22 Research Groups
- 3 Departments



Diverse requirements

Synthesis Lab

- Focus on:
 - Creating new chemicals
- Workflow:
 - Plan, Synthesize, Characterize, Analyze results, Iterate, Test in some application
- Challenges:
 - Multi-disciplinary
 - Harmonize ELN with diverse data sources



Diverse requirements

Self-service facilities

- Focus on:
 - routine measurements
- Can be performed with minimal training
- Workflow:
 - Schedule, Measure, Retrieve data, Analyze data,
- Challenge:
 - Harmonize users' ELN with instrument
 - Associate data with sample



Diverse requirements

Testing facilities

- Focus on:
 - Behavior in applications
 - Testing parameters
- Workflow:
 - Plan, Schedule, Measure, Retrieve data, Analyze data
- Challenges:
 - Analysis
 - Linking data-sample-conditions



Diverse requirements

Large facilities

- Focus on:
 - Characterizing
- Workflow:
 - Plan, Schedule, Measure, Retrieve data, Analyze data
- Major challenges:
 - Data sizes
 - Integrating with home ELN



Roles and skills



Student (Master/PhD)

Skills

- Conducts experiments
- Documents results

IT interactions

- Measure things
- Interact analysis software



Post-Doc

Skills

- Designing experiments
- Analysis workflows
- Documents results

IT Interactions

- Interact measurement software
- Interact analysis software
- Supervise students
- Review results



Principal investigator

Skills

- Provides research questions
- Supervises research of Students and Post-Docs
- Administrative Tasks

IT Interactions

- Interacts with management software

Roles and skills



Technical Staff

Skills

- Maintenance Laboratory equipment
- Keep services running
- Administer stock
- Supervise experiments

IT Interaction

- Interact with monitoring software



Group Leaders

Skills

- Project management

IT Interaction

- Interact with management software



Directors

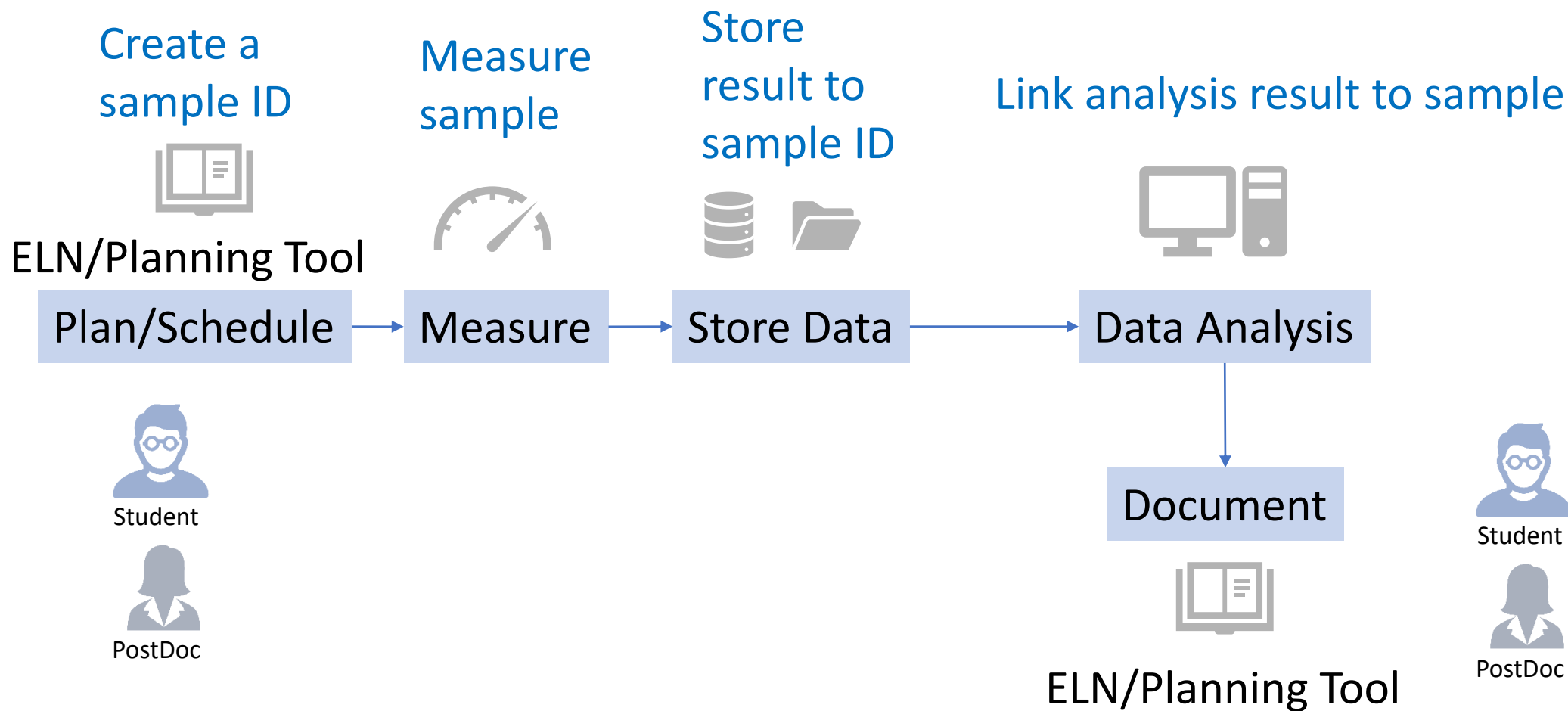
Skills

- Management

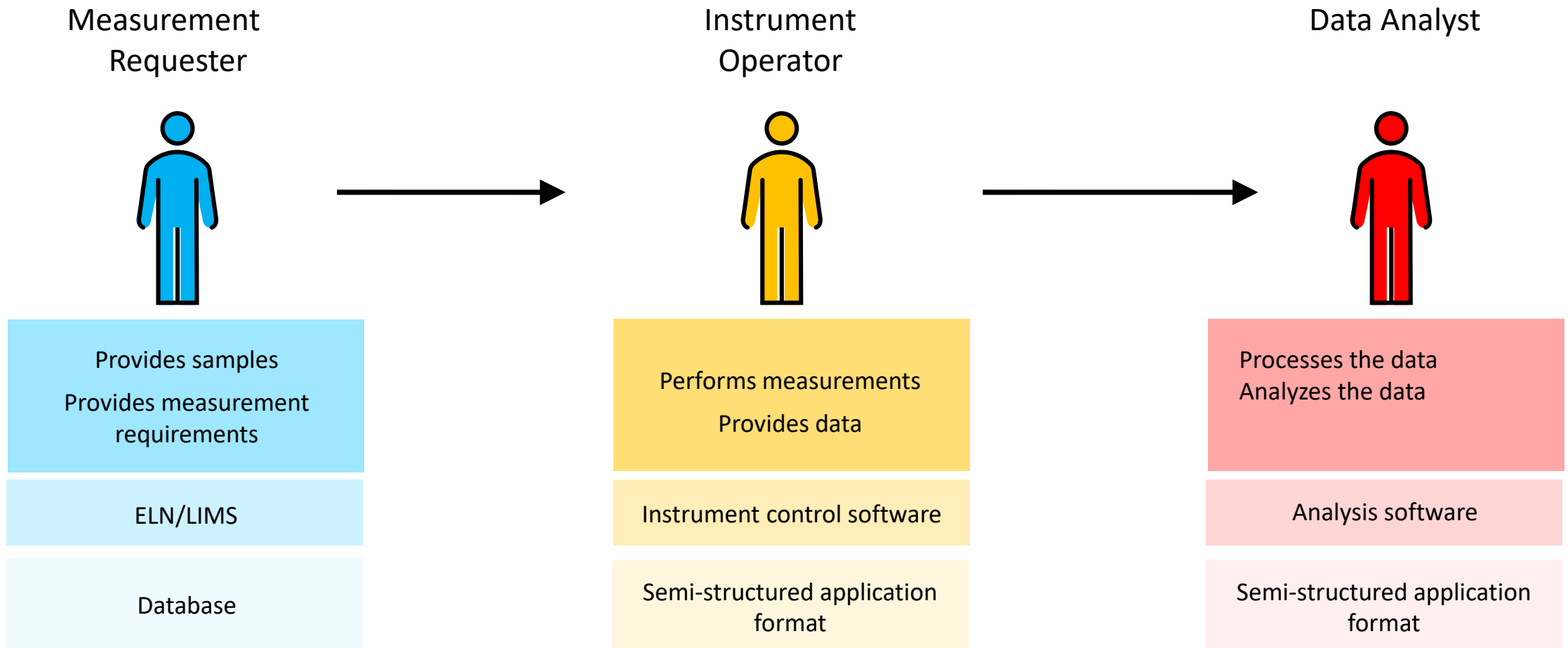
IT Interaction

- View dashboards and reports

Use Case: Self-service testing facility



Different Users; Same Data



2. Problem Statements



Problem Space

Problem Statement

- Data assets are not organized.
- It is distributed across many locations, with no contextual metadata.
- Thus, searching and organizing tools cannot be used to utilize the data.
- Knowledge cannot be automatically extracted from it.



Problem Statements

- Researchers spend too much time on repetitive manual work, related to organizing, searching and processing data.
- Takes away from value-added work, increases errors, leads to re-work.



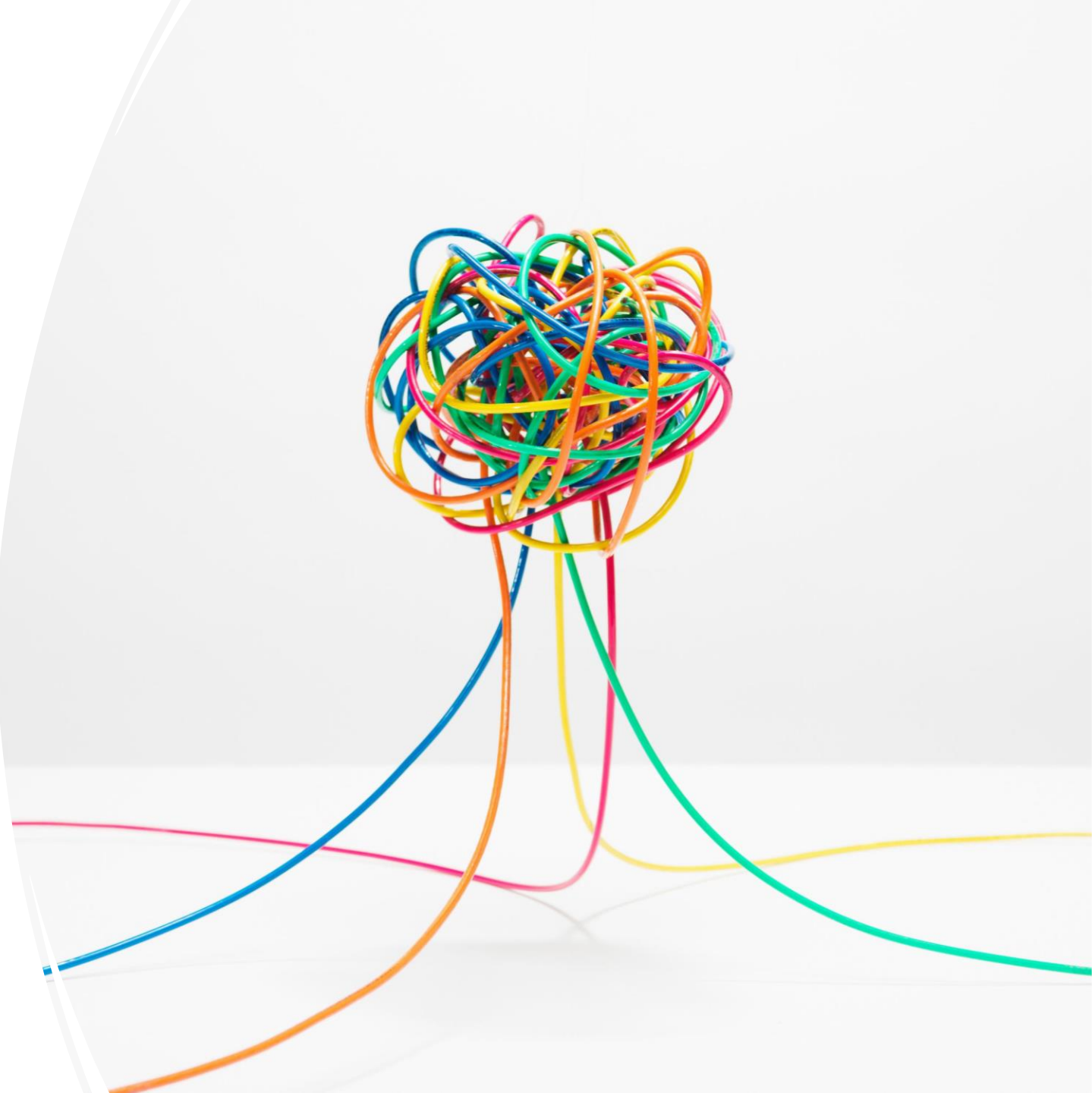
Problem Statements

- Users do not have easy access to all their data assets.



Problem Statement

- It is difficult, sometimes impossible, to trace back the origin of a research result.
- Leads to excessive time spent searching when report revisions are needed.
- Decreases knowledge retention.



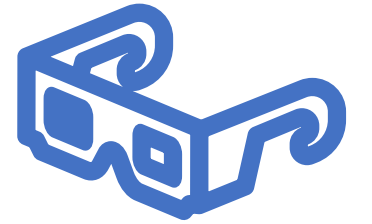
A top-down view of a wooden surface covered with a variety of board game components. In the upper right, there are several chess pieces, including a black king and a light-colored pawn. To their left are two dice, one green and one yellow, and a red die with white pips. A blue domino with white dots is also visible. In the center, there are several cylindrical pieces in green, blue, and yellow, along with a black circular piece featuring concentric rings. A white die with black pips is in the lower left. A black domino with white dots is at the bottom. A red circular piece with a decorative pattern is on the right. A black and white card with a spade symbol is partially visible on the left. The background is a light-colored wooden surface with a dark, checkered pattern.

- The structure of the organization's data assets are not suitable for large-scale analysis algorithms.
- Unable to utilize modern algorithms for meta-analysis.

Example of governance problems

	A	B	C	D	E	F	G	H	I
1	ID	building_name	building_identificati	building_identificatic	room_number	floor_num	storage_name	storage_characteristic	
2	1	Physikgebäude	D	Physikgebäude	110	TKG	Kuhlschrank 2 +4Â°C		
3	4	Physikgebäude	D	Physikgebäude	102	TKG	102-1-Hazards		
4	5	Physikgebäude	D	Physikgebäude	102	TKG	102-2-Solvents		
5	6	Physikgebäude	D	Physikgebäude	103	TKG	Chem	einfacher Schrank ohne Absaugung	
6	7	Physikgebäude	D	Physikgebäude	104	TKG	104-1-Hazards		
7	8	Physikgebäude	D	Physikgebäude	104	TKG	104-3-Bases		
8	9	Physikgebäude	D	Physikgebäude	104	TKG	104-4-Acids		
9	10	Physikgebäude	D	Physikgebäude	104	TKG	Betriebsmitteltonne	Spannringfass 30L oder 60L zur Entsorgung	
10	11	Physikgebäude	D	Physikgebäude	107	TKG	107-3-Hazards		
11	12	Physikgebäude	D	Physikgebäude	107	TKG	107-4-Hazards		
12	13	Physikgebäude	D	Physikgebäude	107	TKG	107-1-Solvents		
13	14	Physikgebäude	D	Physikgebäude	107	TKG	Kuhlschrank		
14	15	Physikgebäude	D	Physikgebäude	107	TKG	107-5-Acids		
15	16	Physikgebäude	D	Physikgebäude	107	TKG	Betriebsmitteltonne	Spannringfass 30L oder 60L zur Entsorgung	
16	17	Physikgebäude	D	Physikgebäude	110	TKG	110-2-Hazards		
17	18	Physikgebäude	D	Physikgebäude	110	TKG	110-1-Inorg Acids		
18	19	Physikgebäude	D	Physikgebäude	110	TKG	Betriebsmitteltonne	Spannringfass 30L oder 60L zur Entsorgung	
19	20	Physikgebäude	D	Physikgebäude	111	TKG	Glovebox A	nur für ungefährliche Substanzen geeignet	
20	21	Physikgebäude	D	Physikgebäude	111	TKG	Glovebox B	nur für ungefährliche Substanzen geeignet	
21	22	Physikgebäude	D	Physikgebäude	101	TKG	Glovebox C	nur für ungefährliche Substanzen geeignet	
22	24	Physikgebäude	D	Physikgebäude	113	TKG	F90 Gase	F90	
23	25	Physikgebäude	D	Physikgebäude	113	TKG	F90 Gase Formiergas	F90	
24	26	Physikgebäude	D	Physikgebäude	113	TKG	Betriebsmitteltonne	Spannringfass 30L oder 60L zur Entsorgung	
25	27	Physikgebäude	D	Physikgebäude	142	TKG	Hausdienste	ganzer Raum	

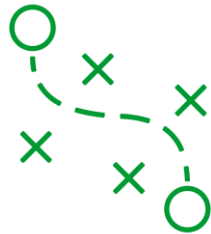
3. Data Governance



Solution Space

What is a data strategy?

- A data strategy can be considered as an approach that enables to derive new knowledge from data
- A comprehensive data strategy is the basis for the successful implementation of data related projects



A **Data Strategy** describes the ...

- ... **organizational structure** for the successful use of data with relevant **processes** for dealing with data
- ... required **skills & roles**
- ... **technology** and tools

Goals of a Data Strategy

Remove
Information Silos



Keep knowledge of
research



Make research
more accessible



Enable meta
research

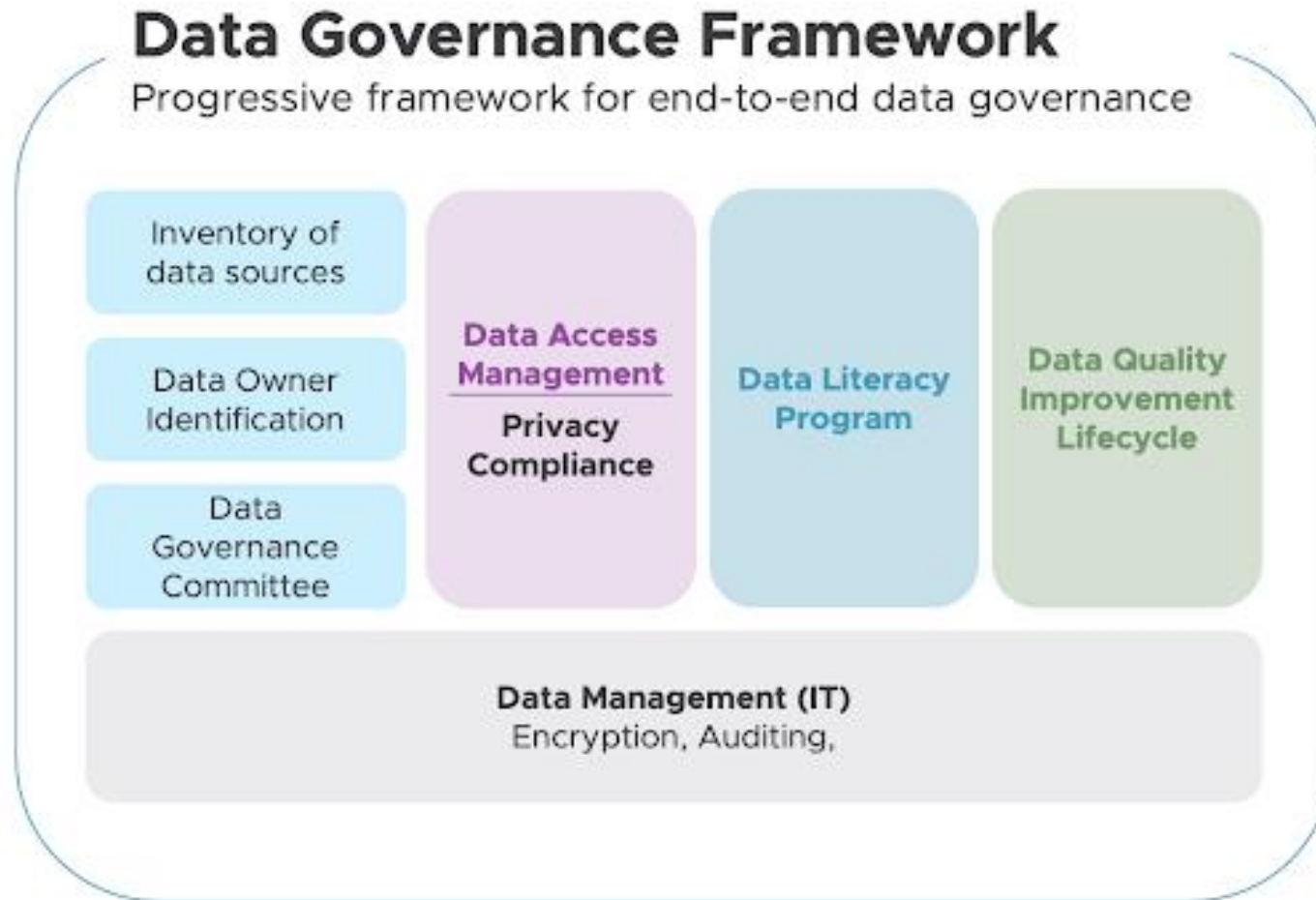


Improve
sustainability



<https://datagovernance.com/goals-and-principles-for-data-governance/>

Data Governance Concepts



Roles and skills in a business context

Data Engineer

Transform and harmonise data

Data Architect

Provide data processing concepts

Data Scientist

Analyse and model data

Data Artist

Visualize data

Data Custodian

Data storage and security

Data Steward

Steering and household data

Data Security Admin.

Data security concepts

Domain Expert

Domain knowledge

Data Evangelist

Explores data potential

Data Governance Roles



Student (Master/PhD)

DG Roles

- Data engineer
- Data scientist



Post-Doc

DG Roles

- Data scientist
- Data architect
- Data engineer
- Domain expert



Principal investigator

DG Roles

- Data steward
- Data evangelist
- Domain expert

Data Governance Roles



Technical Staff

DG Roles

- Data steward



Group Leaders

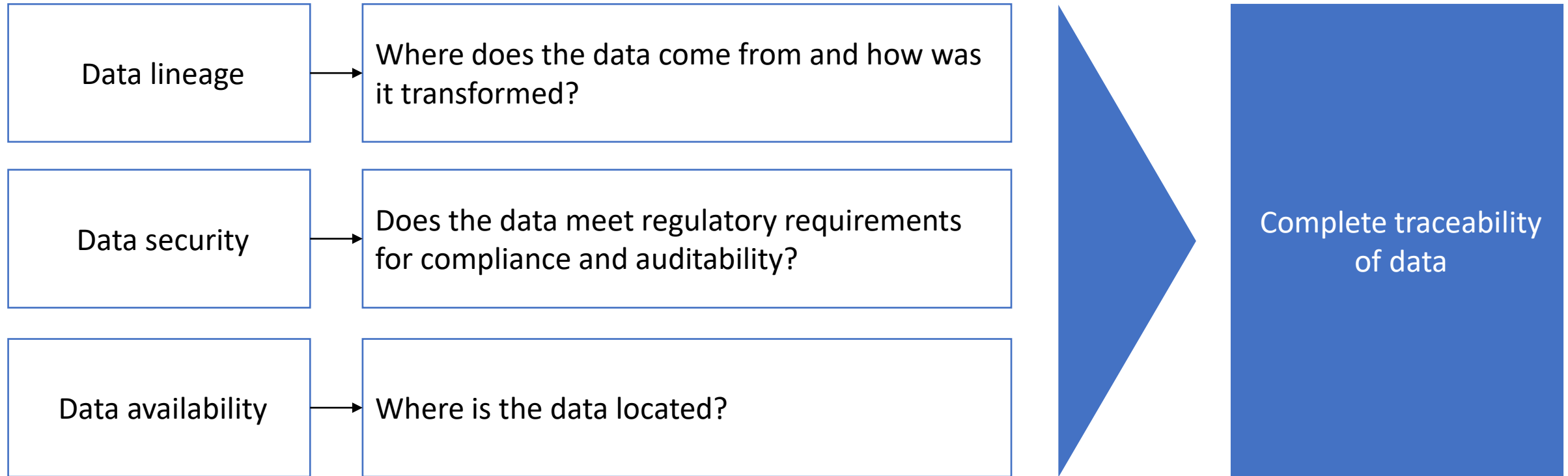
DG Roles

- Data steward
- Data evangelist
- Domain expert

Take-away message

- Data governance has specialized roles
- Academia often fulfills these roles by loading responsibility onto operational staff
- It may be time to have dedicated data-governance roles in academia
- However, due to scale and diversity at academia organizations, it may need a different structure than is used in industry

Data Governance Concepts



What is metadata?


- Data about data 
- Metadata can describe a single piece of data, a dataset or collection.
- Standard types of metadata:
 - **Descriptive**: information about **who** created a resource, **what** it is **about** and **what** it **includes** (e.g. title, author, subjects, keywords etc.)
 - **Structural**: information about the **way** data elements are **organized**, their **relationship** and the **structure** they exist in (e.g. ER-model)
 - **Administrative**: information about the origin of resources, their type and access rights (e.g. file type, date of creation etc.)

Table with 4 books, created by Joe Dow

Descriptive METADATA

ID	Title	Author	Year	Cover	Edition	Price
1	Emma	Austen	1815	Paperback	20th	€ 4,30
2	Dracula	Stoker	1897	Hardback	15th	€ 10,00
3	Ivanhoe	Scott	1820	Hardback	8th	€ 20,00
4	Kidnapped	Stevenson	1886	Paperback	11th	€ 3,50

Origin of resources: Book store

Access rights: read only - everyone; write and read – Joe Dow

Created on 5. January 2019

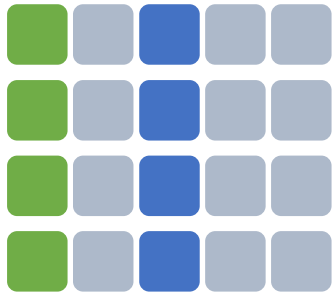
Administrative METADATA

Unsorted table;
related to sales data;
key is 'ID'

Structural METADATA

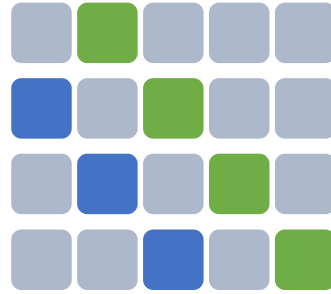
What is data structure?

▶ Data structure is the particular way of
ORGANIZING & **STORING** digital information.



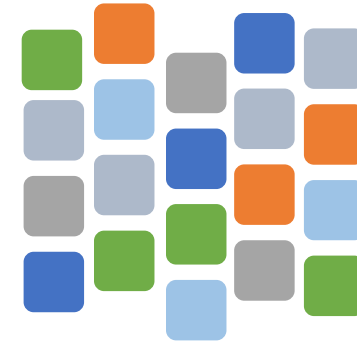
Structured data:

Information with a **specific**
and **high degree** of
organization (tabular form)



Semi-structured data:

Information with **some**
degree of organization






Unstructured data:

Information with **no pre-**
defined organizational
structure



How do we store data?

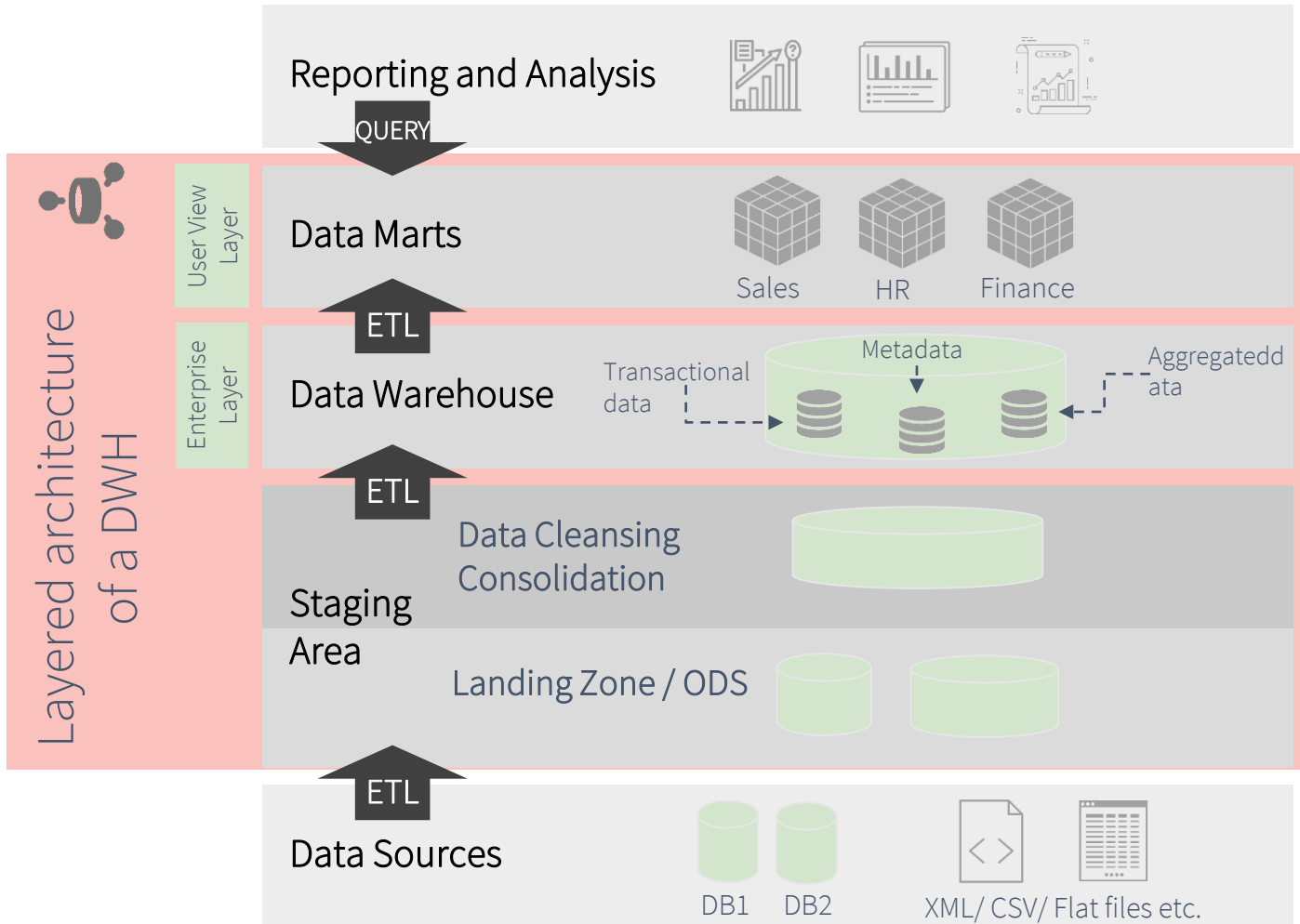
	Databases		File System	Object Storage
Definition	<p>Collection of stored data organized in and determined by the data model underlying the database (e.g. ER-diagram)</p> 		<p>Store data in a space with a pre-defined scheme (e.g. a file hierarchy)</p> 	<p>Store data as objects in a space with no pre-defined scheme</p> <p>Every file has a unique identifier, so it can be found (e.g. URL)</p> 
	<p>Relational DB</p> <ul style="list-style-type: none"> Based on relational data model (data is organized in tables) Structured Query Language (SQL) for querying the database 	<p>Non-relational DB</p> <ul style="list-style-type: none"> Based on any data model other than the relational model Examples include key-value stores, document stores and graph databases 		
Examples	<ul style="list-style-type: none"> MS SQL Server, MySQL, Oracle 	<ul style="list-style-type: none"> MongoDB, Cassandra, Neo4J 	<ul style="list-style-type: none"> NTFS (Windows), Hadoop Distributed File System, ext3 	<ul style="list-style-type: none"> AWS S3 buckets Azure Blob Storage

4. Data Architectures



Solution Space

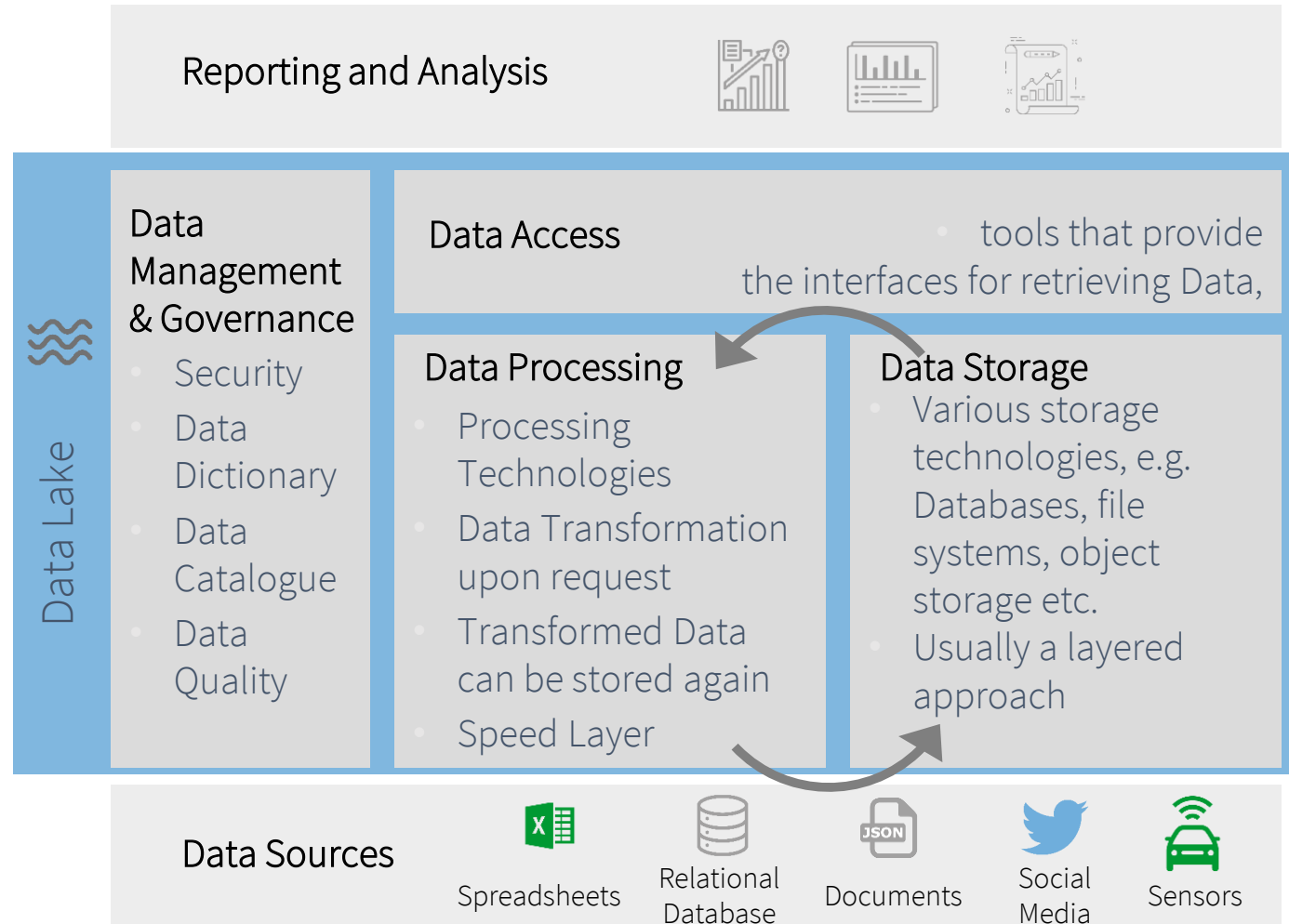
What is a Data Warehouse?



“A Data Warehouse is a subject oriented, integrated, non-volatile and time-variant collection of data in support of management decisions”*

*Source: Immon, WH. Building the Data Warehouse

What is a Data Lake?

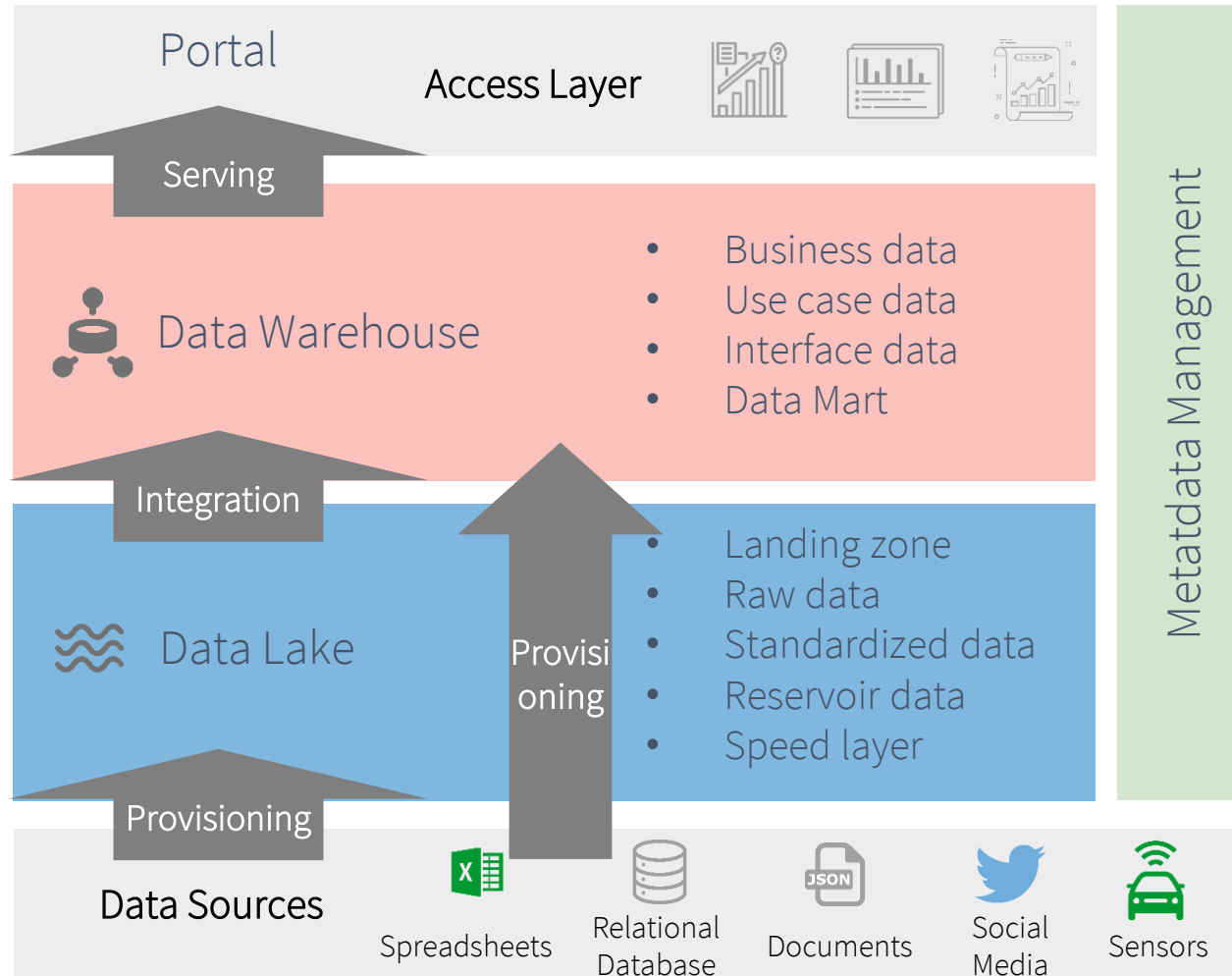


▶ A Data Lake is a modular system of data storage and processing technologies. Like a DWH it is a logical concept rather than a tangible entity.

Benefits of a Data Lake

- stores all kinds of data, e.g.:
 - Structured Tables
 - Text Documents
 - Pictures
- Central entry point for data access

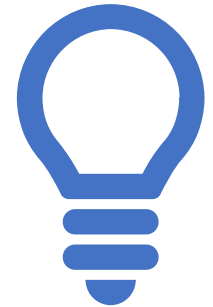
How can both architectures be combined?



Advantage

- One platform for different requirements
- Best of both architectures

5. Technologies

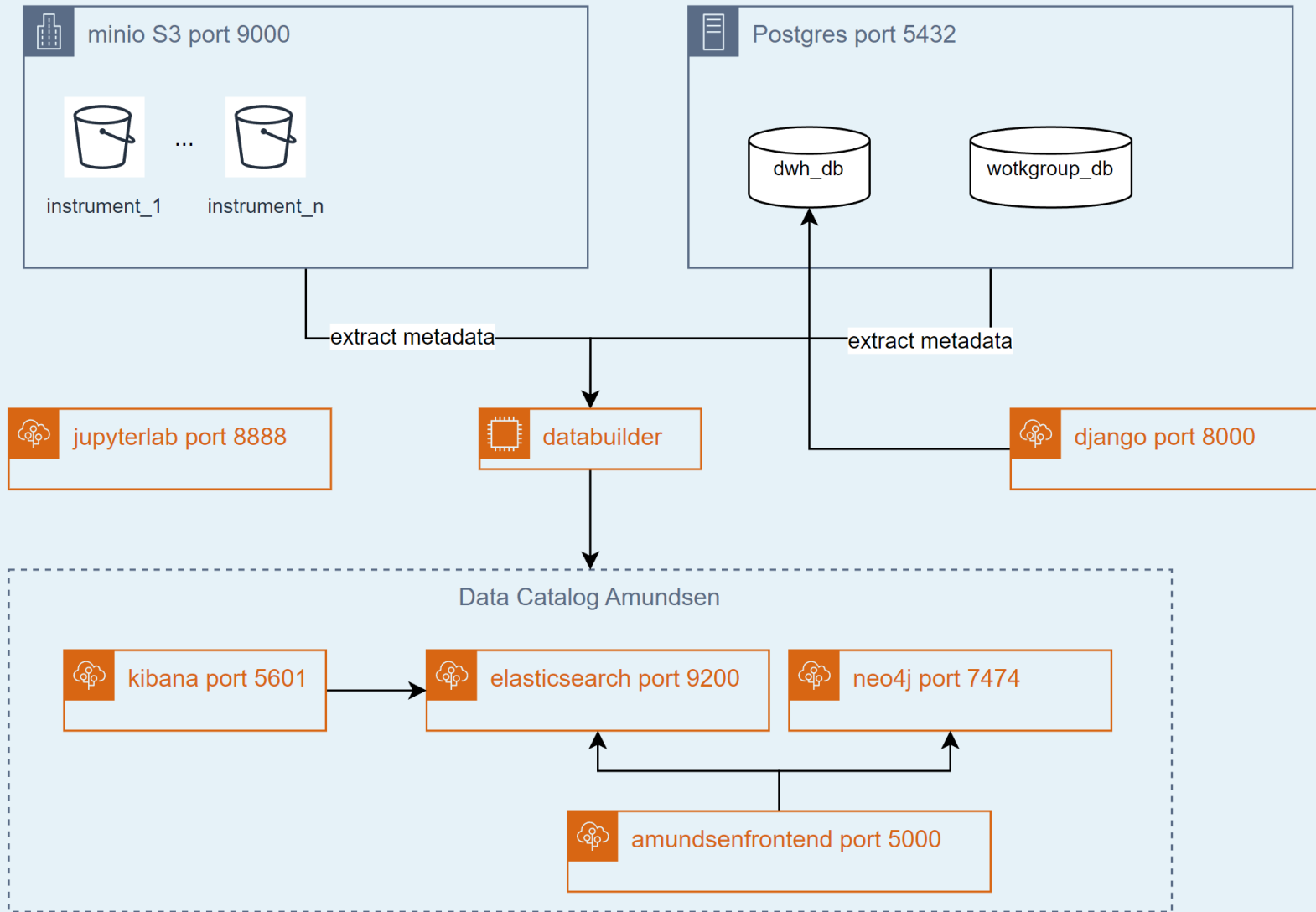


Solution Space



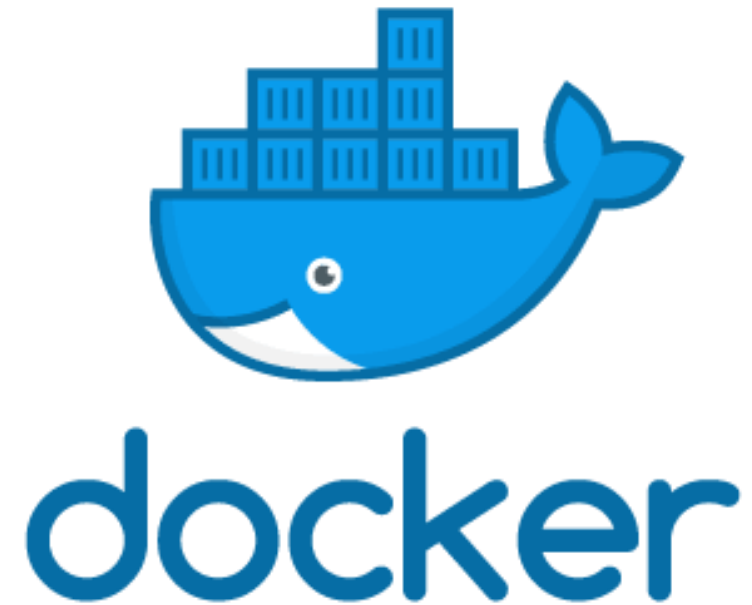
VM: mpicecnet as docker environment

Work in progress



Docker

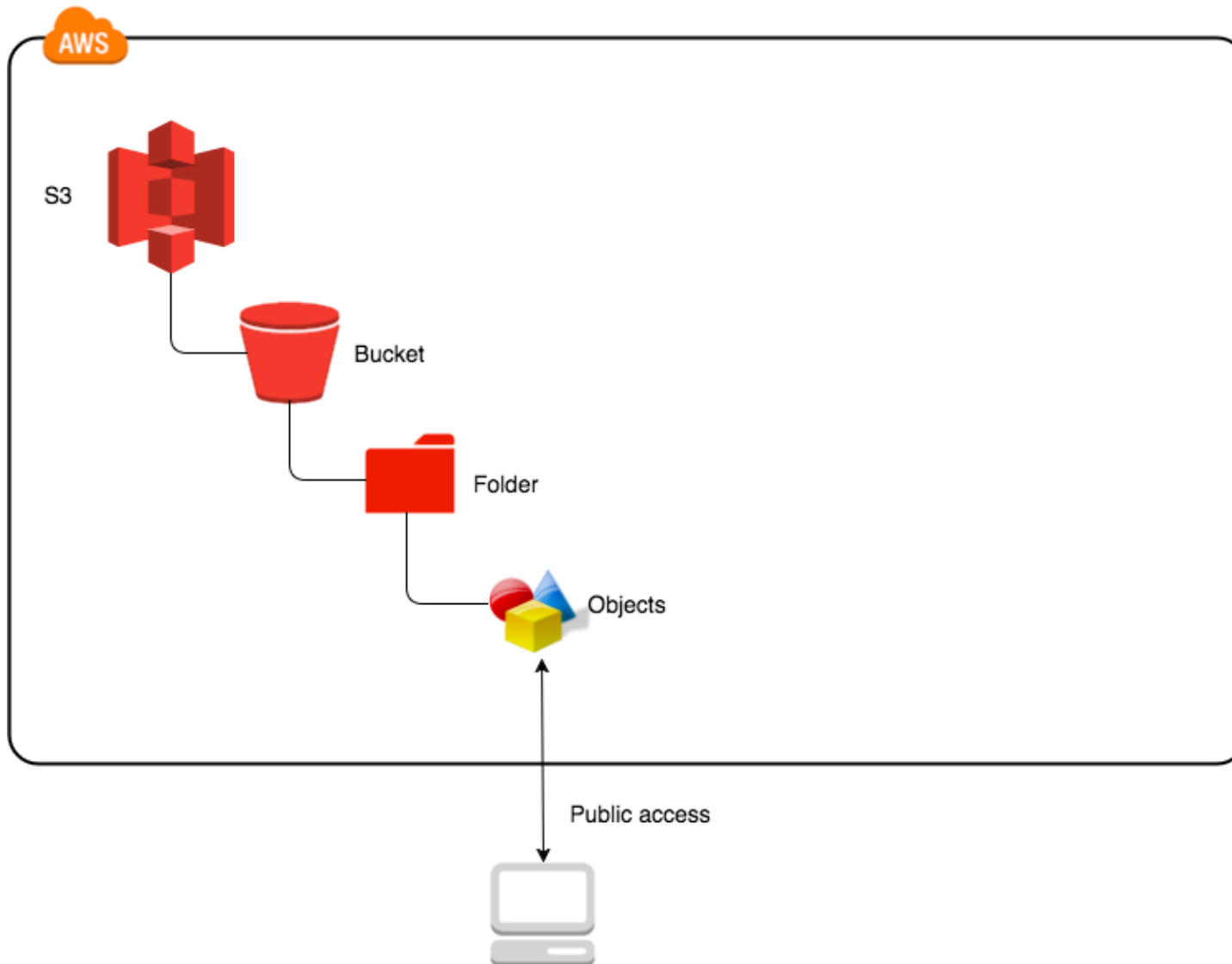
- Containers for running microservices
- Isolated environments
- Needs less resources than a VM



Minio

- Object storage
- Same as Amazon S3
- Can store any kind of data

The Minio logo is displayed in a bold, red, sans-serif font. The letters are thick and blocky, with a slight shadow effect. The 'M' and 'N' have a distinctive shape with a small gap at the bottom. The 'I' is a simple vertical bar. The 'O' is a thick, rounded circle. The entire logo is centered horizontally within a white rectangular area.



MINIO

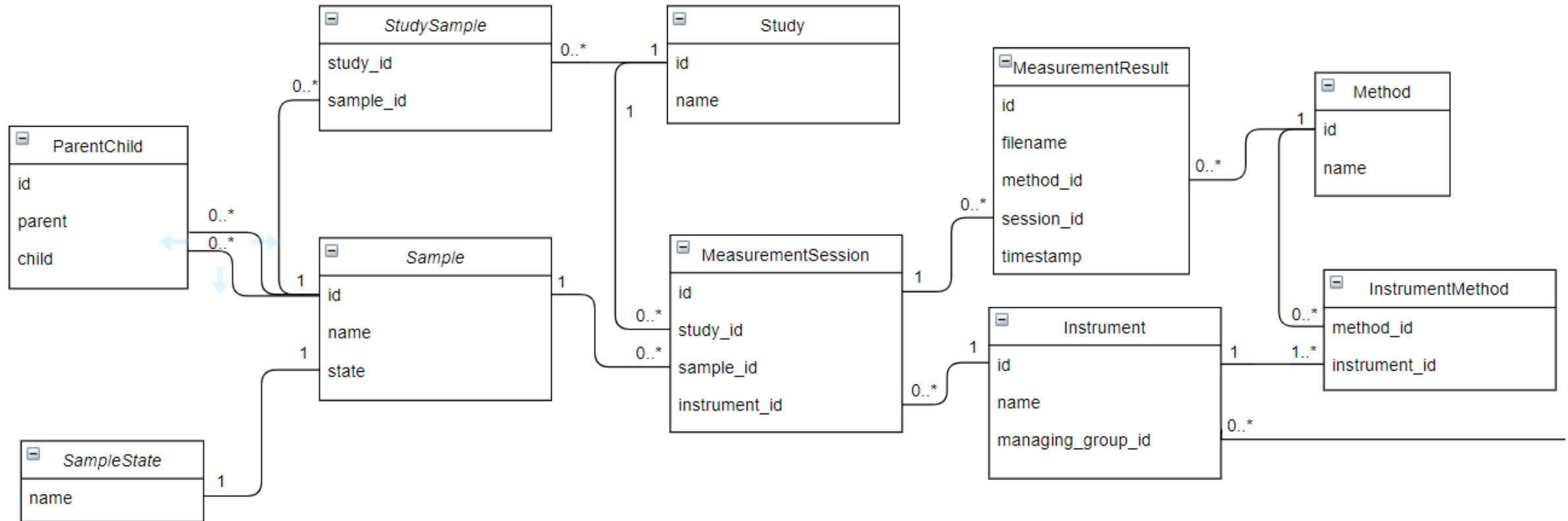
- Buckets hold objects
 - Define access rights
- Objects have immutable metadata
- Client can access data over API

PostgreSQL

- Relational database
- Used for storing structured data
- Transactional data and data warehouse
- Strict and robust data models



- Used for storing:
 - Transactional Data
 - Data Warehouse tables



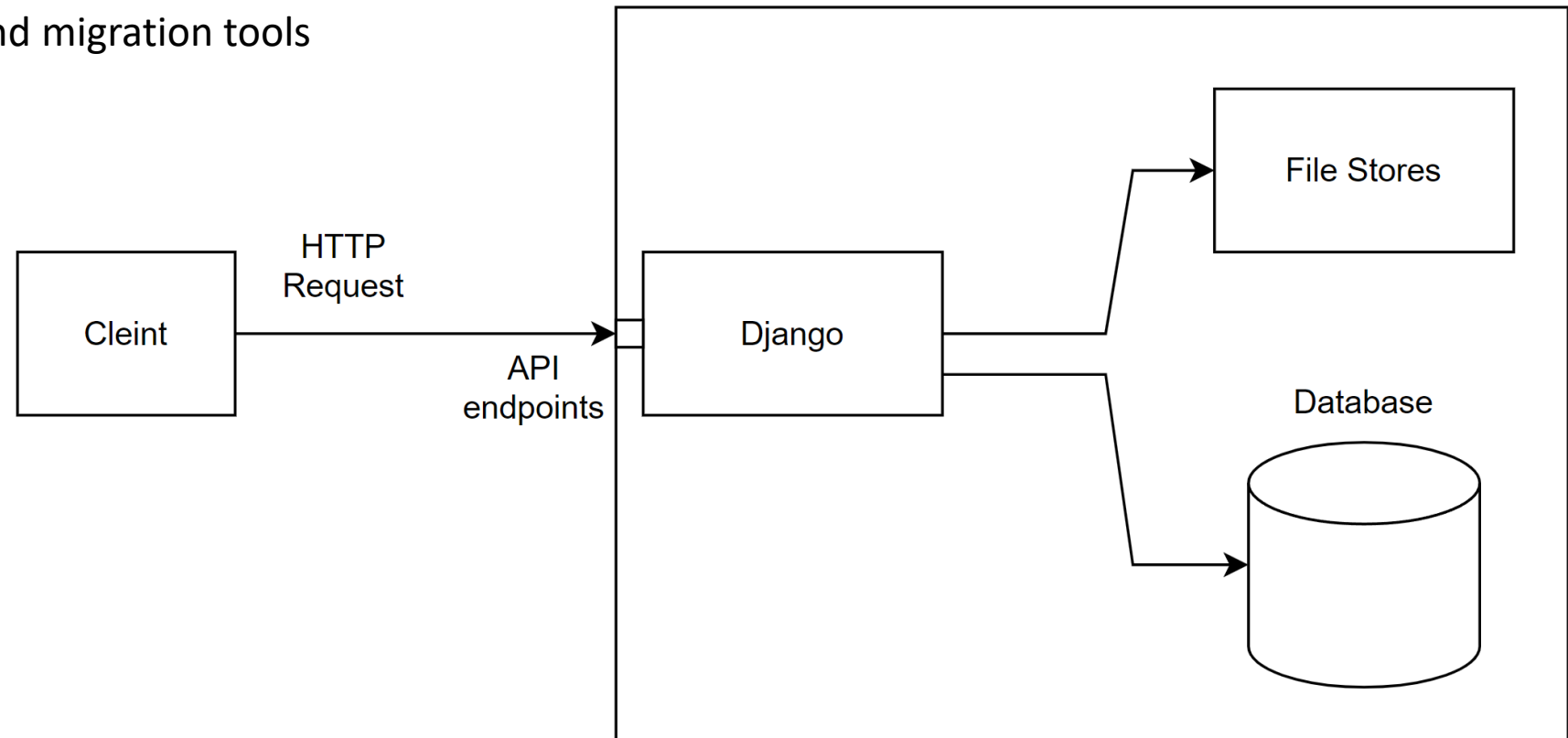
Django

- Back-end Web application framework written in python
- Robust and scalable
- Used for providing REST API for accessing Relational Database

The Django logo, featuring the word "django" in a white, lowercase, sans-serif font, centered on a dark green rectangular background.



- API endpoints provide client with ability to CRUD database entries
- Will be used to connect front-end apps
- Access to File stores
- Django provides Authentication and Authorization to access resources
- Django provides convenient database schema history and migration tools



Amundsen

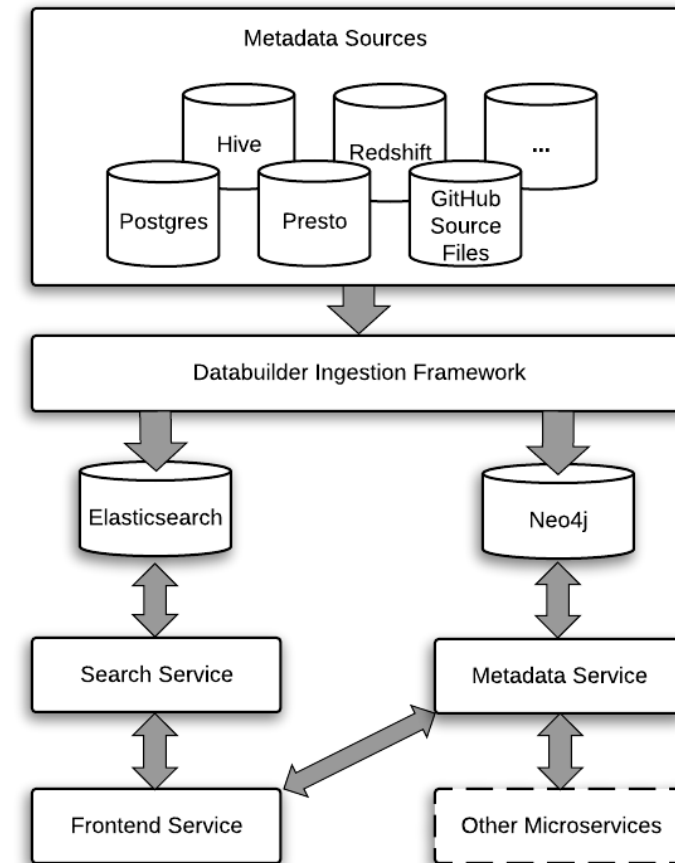
- Data catalog
- PageRank-inspired search algorithm
- Provides REST API for search engine



How does the data catalog work?

- Amundsen collects Metadata via the data builder ingestion framework
- Metadata and lineage data are stored in Neo4j and elasticsearch
- Metadata is made accessible via search interface
- Metadata can be made searchable for all users, whereas the content remains visible only with sufficient access permissions

Amundsen data catalog architecture



Elasticsearch

- Search engine
- Service behind Amundsen's search library
- NoSQL data store
- Stores data as documents (like JSON)



elasticsearch



elasticsearch

Inverted Index

Documents 1 & 2

The bright
blue
butterfly
hangs on
the breeze

Under blue
sky, in bright
sunlight, one
need no
search around



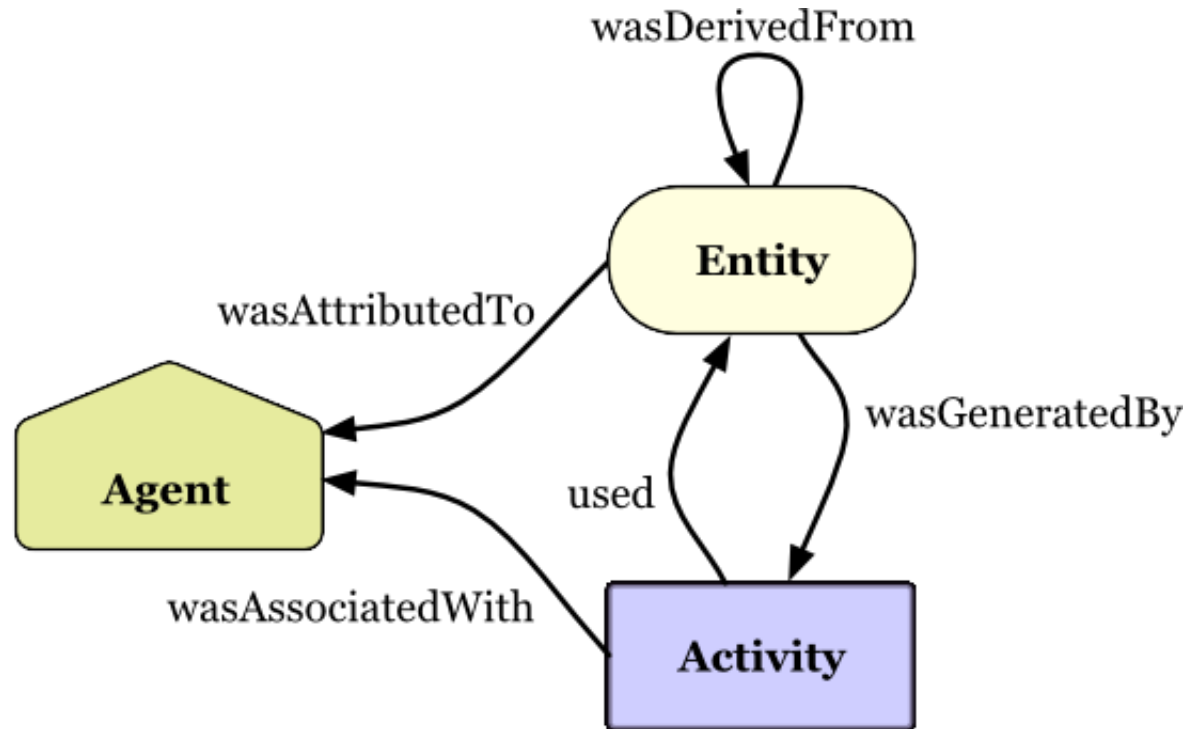
ID	Term	Document
1	butterfly	1
2	blue	1,2
3	bright	1,2
4	retire	2
5	wind	2

Neo4j

- Graph database behind Amunden's metadata service
- Stores 'Triples' (subject, predicate, object)
- Great for traversing relationships



Data Provenance



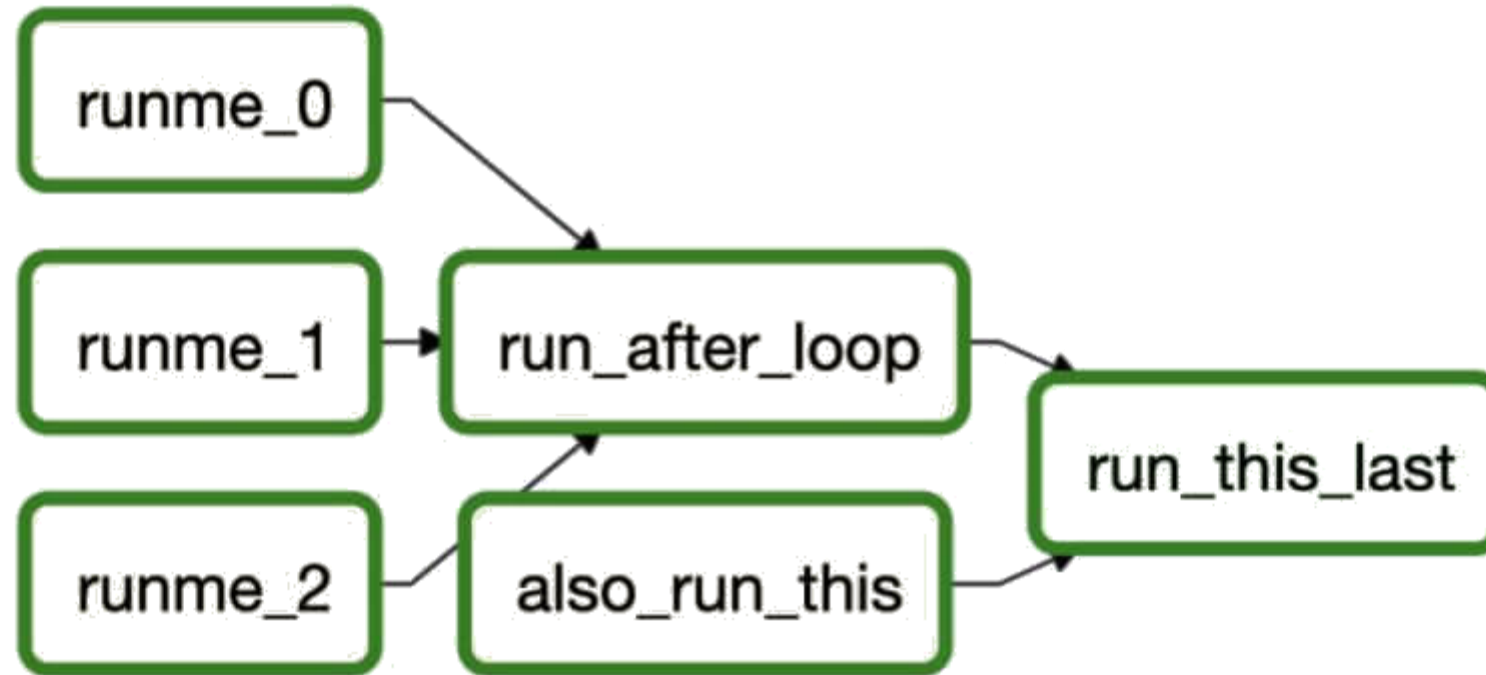
Apache Airflow

- Orchestration tool
- Build for running complex cron jobs
- Will be used for ETL processes



Apache
Airflow

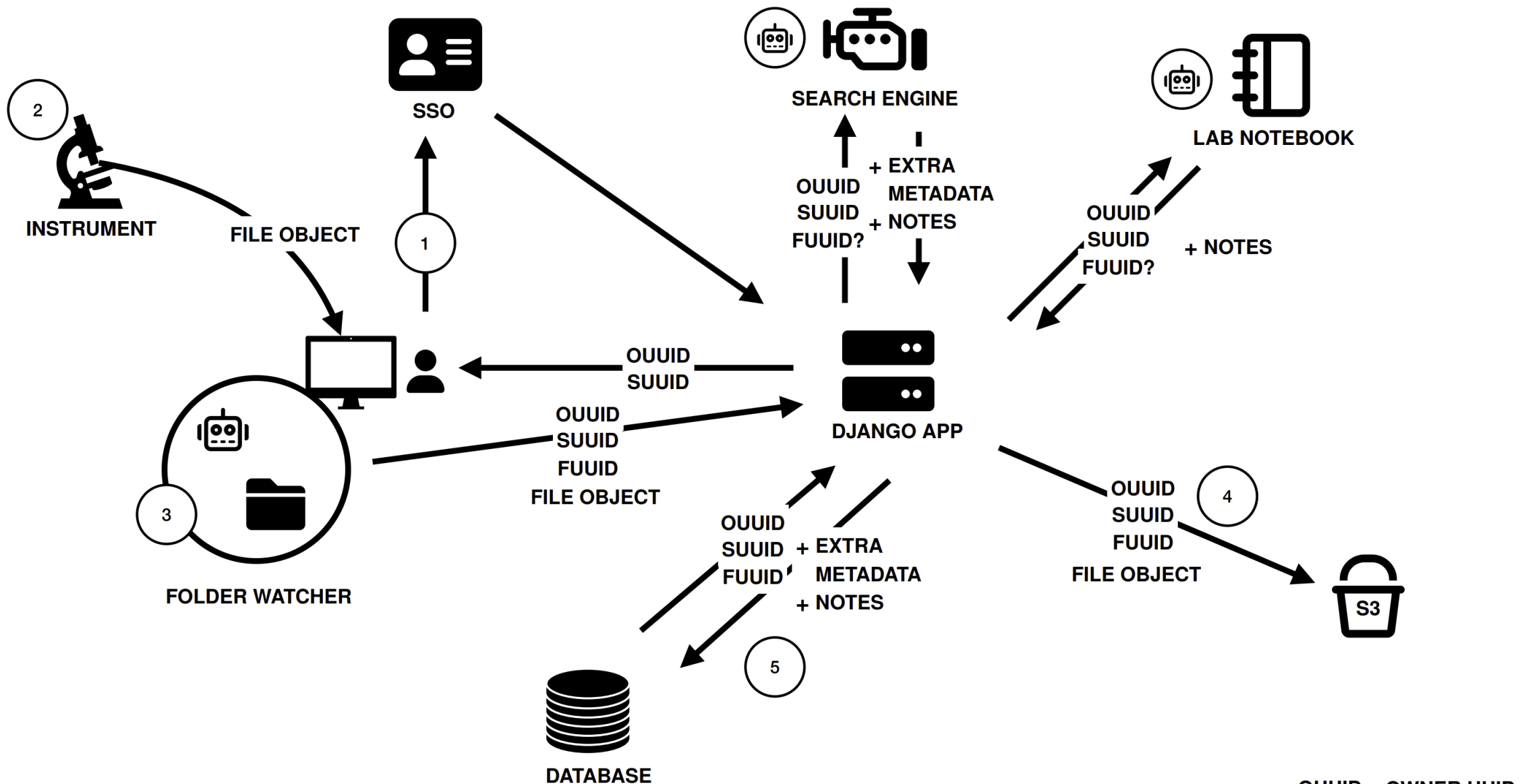
Orchestration



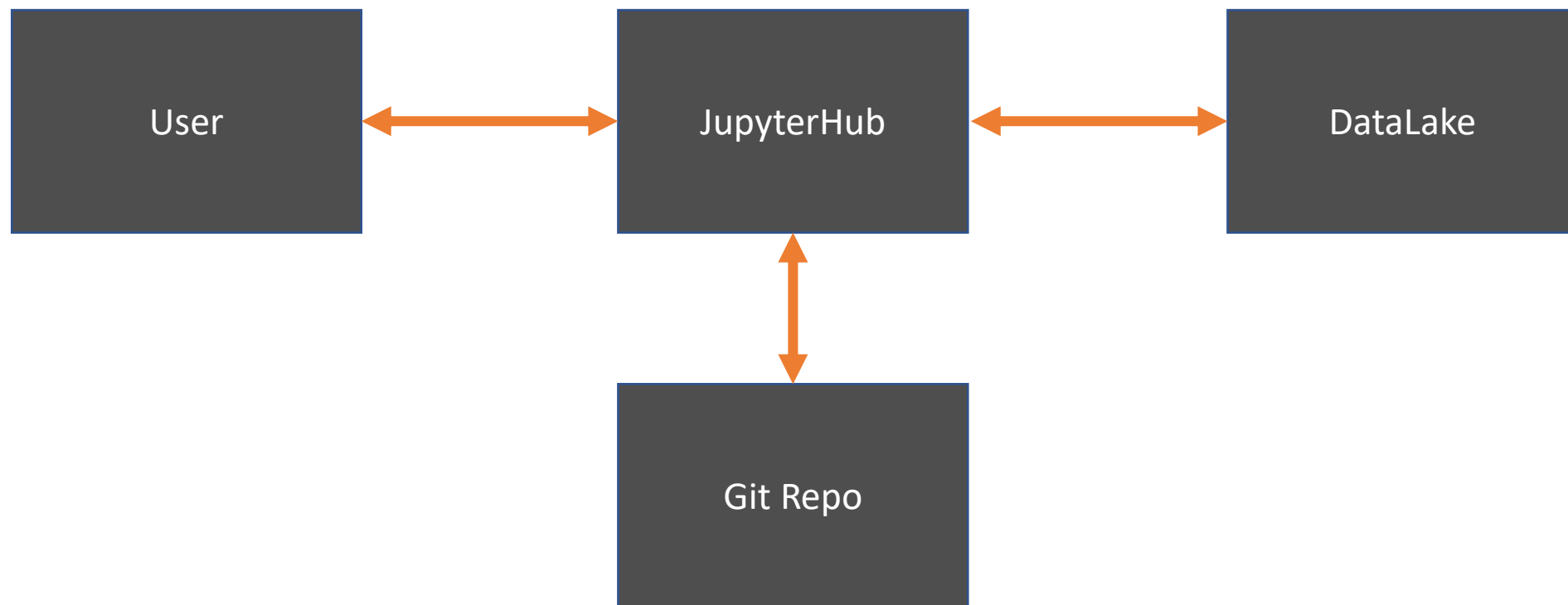
6. Future Directions



Solution Space



OUUID OWNER UUID
SUID SESSION UUID
FUUID FILE UUID



Continued Improvement

- Improved data management architecture
- Get buy-in for having more dedicated data governance roles
- Continually evolve data governance policies



7. Summary

Problem Space

- Diverse requirements
- Non-conventional use case
- Overlapping roles

Solution Space

- Common Data governance Philosophies
- Common Data Governance Architectures
- Lots of tech needed
 - Docker, Object storage (S3), RDBMS (PostgreSQL), Data catalogue (Amundsen), Orchestration (Airflow)
- Change
 - New roles, new infrastructure, new concepts for data management in academia