



Towards Data Democratization

# Evolving AWS datalake to data lakehouse

Data Platform Engineering

Jake Demšar, April 2024

# Our datalake and its problem

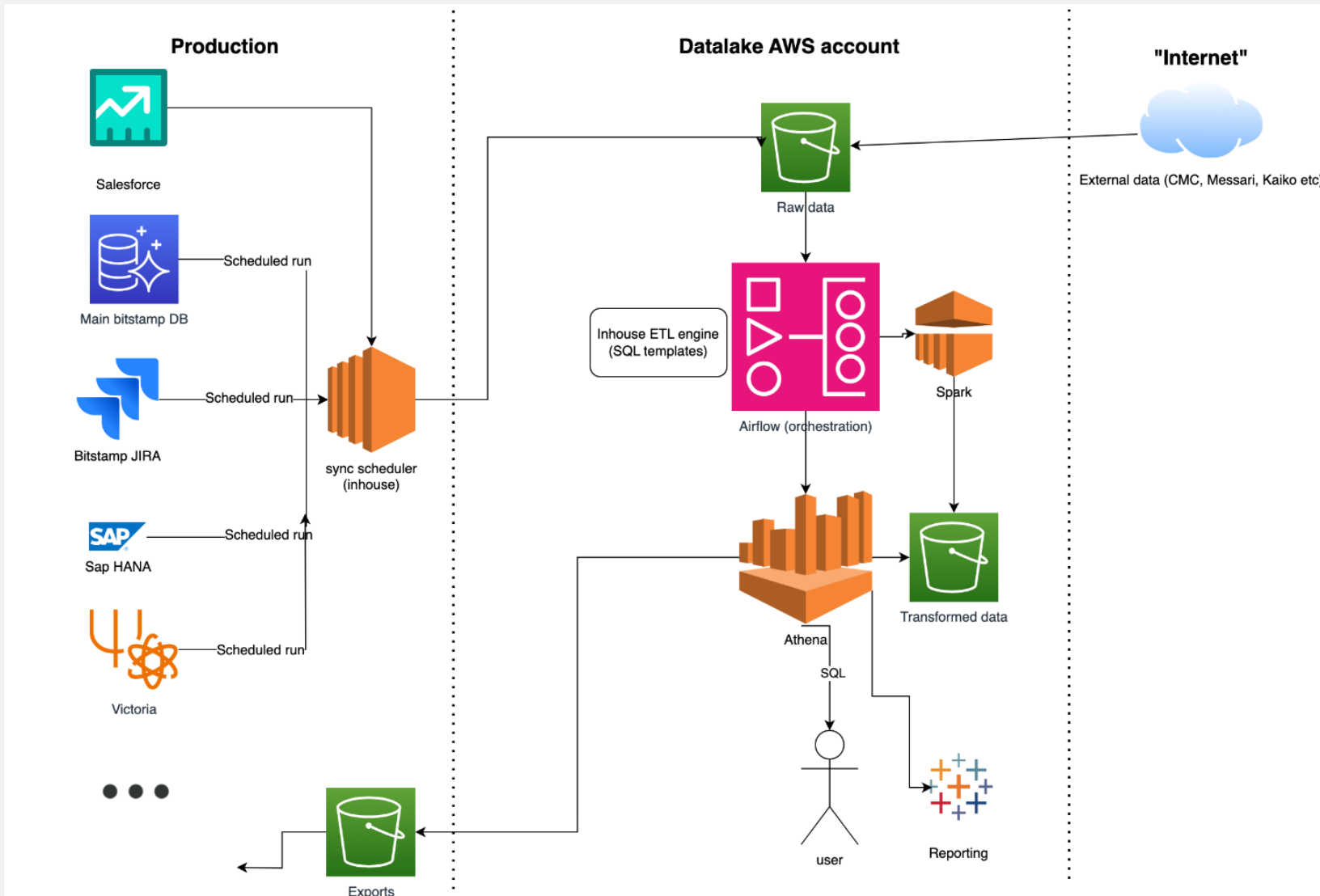


## A short story

1. We conceptualise our **datalake**
2. Its slowly but surely degenerates into a data **swamp**
3. The **pain** is not great enough
4. The datalake **grows**, the userbase **grows**, the demands **grow**
5. The **pain is great enough**
6. We salvage the situation by introducing data **lakehouse**
7. **< we are here >**
8. The proposed golden paths and pattern sustain a healthy data-driven culture

# Athena-based datalake

B



1. Hoard raw data in **S3**
2. Expose raw data in Athena
3. Use **in-house ETL** (SQL templates) + Athena to store transformed data (**.parquet**) in S3
4. Expose parquets via **Glue catalog** ("thin layer")
5. Serve data to users via **Athena SQL interface**

Orchestration: **Airflow**  
Data catalog: **Amundsen**

**IDEMPOTENT**

# A neat solution!



Endless capacity of S3



Serverless nature of Athena



Declarativity and simplicity of SQL



A small team (2 to 5 members) can efficiently manage and further develop the datalake

Cca 1PB of data

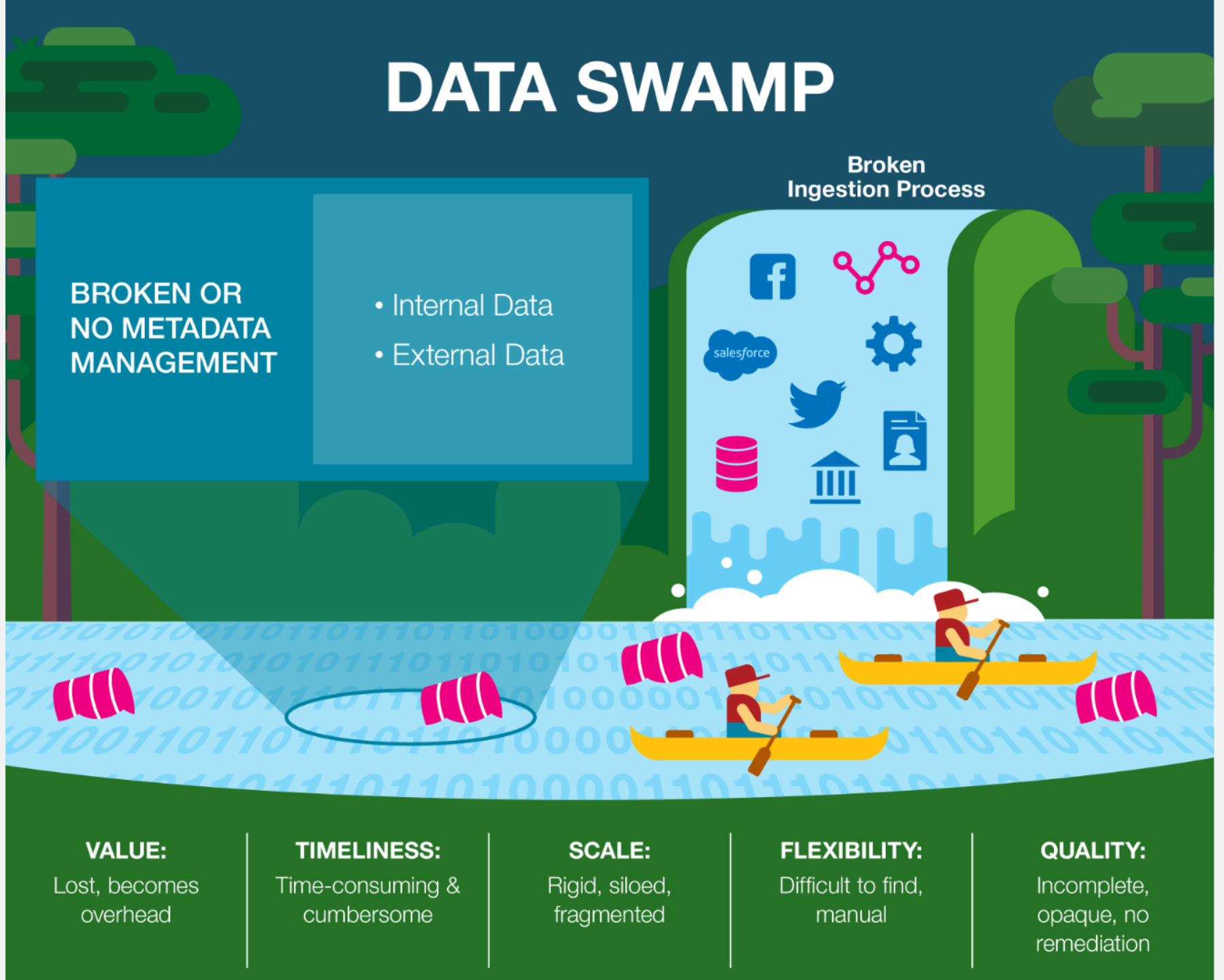
In 1000+ tables synced

500+ interdependent tables calculated

50-100 DAGs run

Serving internal and external data needs

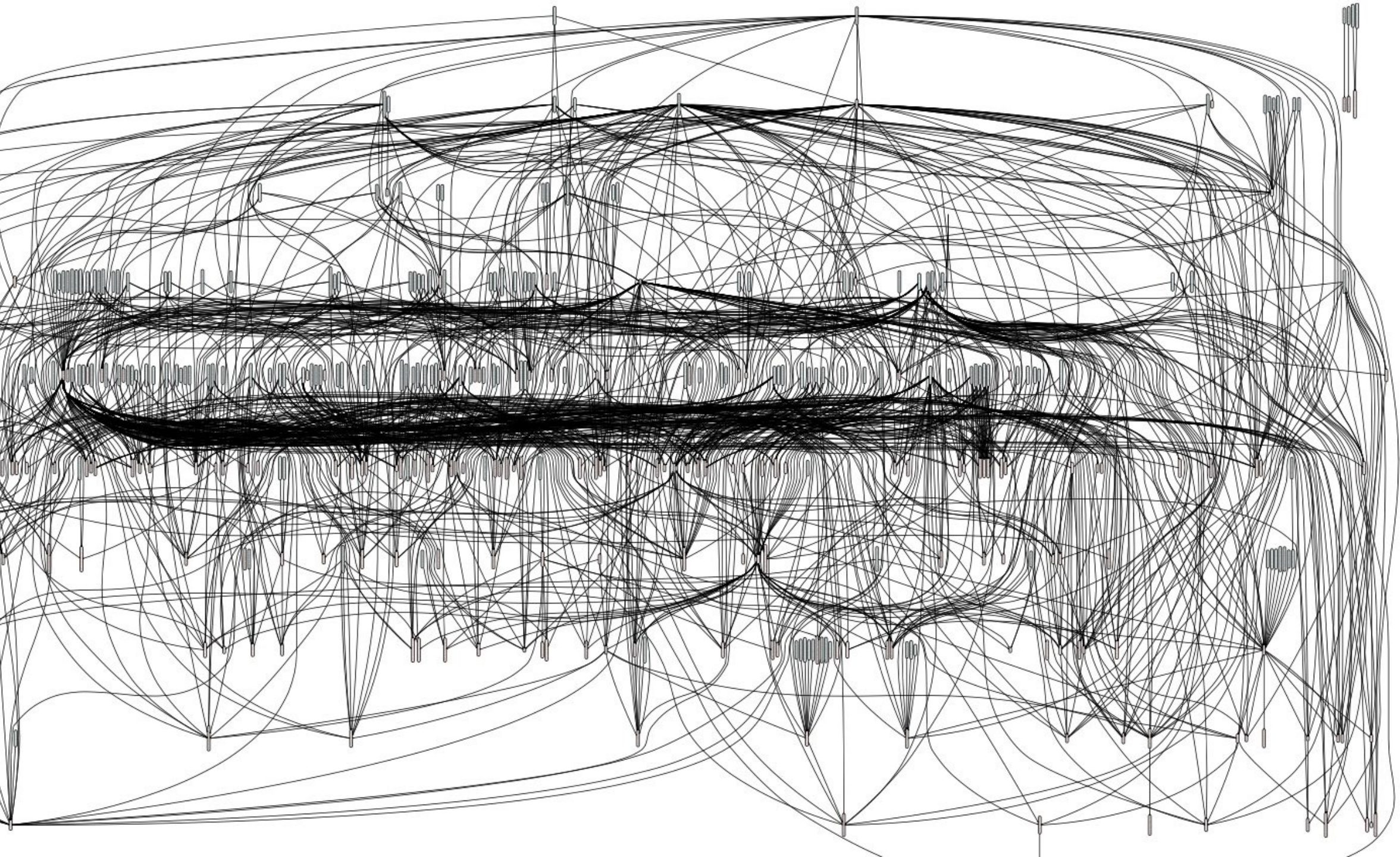
# What is the problem, then?







# Data mesh-mess





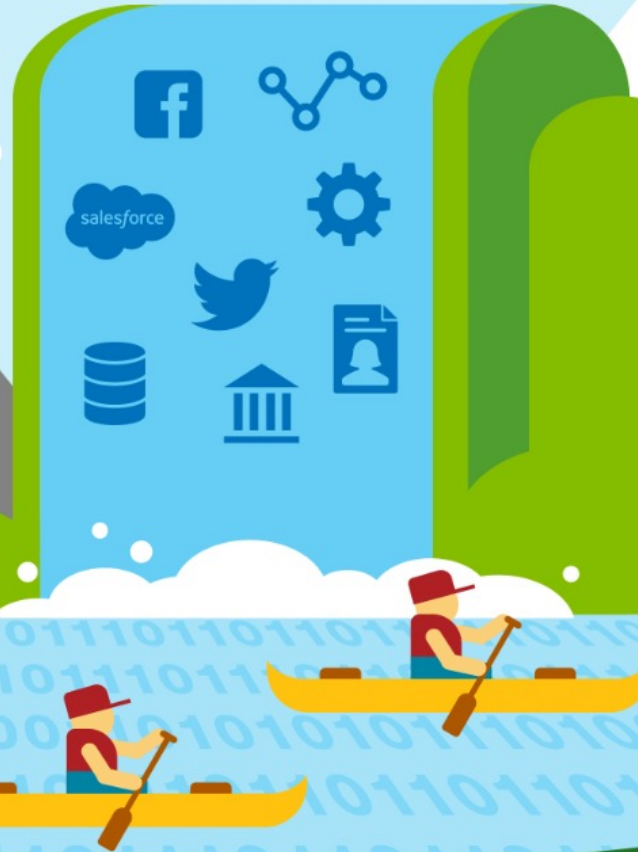
# DATA LAKE

## METADATA MANAGEMENT

- Processes
- Properties
- Relationships
- Tags

- Web Server Logs
- Databases
- Social Media
- Third Party Data
- CRM Data

## Streamlined Ingestion Process



### VALUE:

Added,  
self-service, truly  
data-driven

### TIMELINESS:

Always ready,  
easy to find

### SCALE:

Robust  
infrastructure  
supports growth

### FLEXIBILITY:

Easily modified,  
automated &  
streamlined

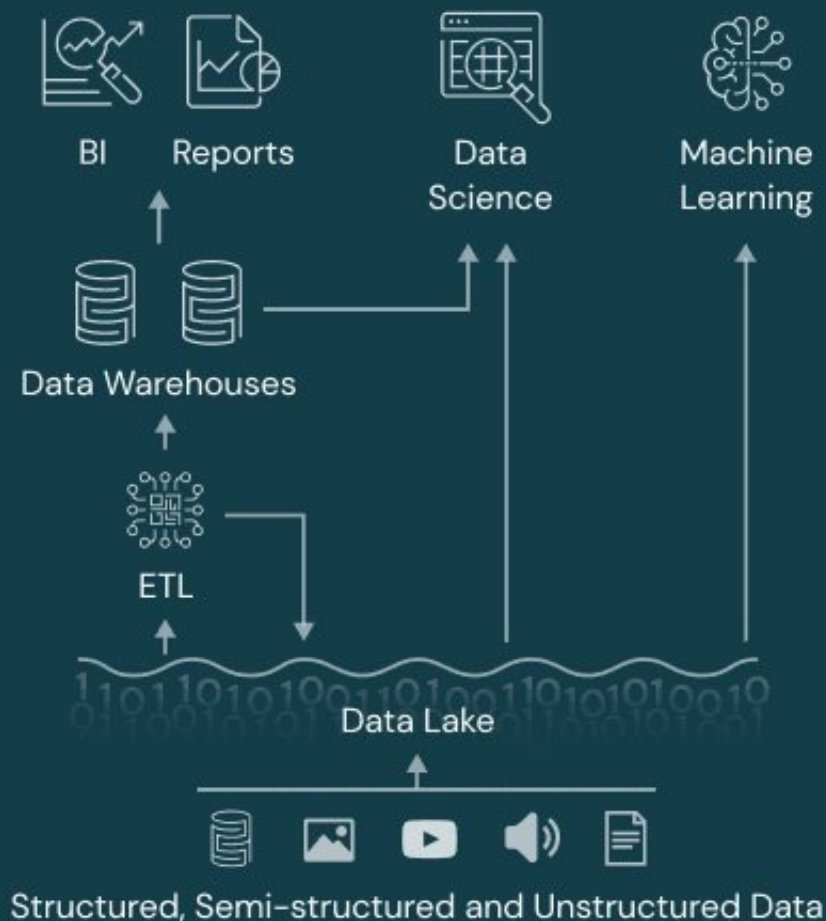
### QUALITY:

Explicit visibility,  
easily understood  
& trustworthy

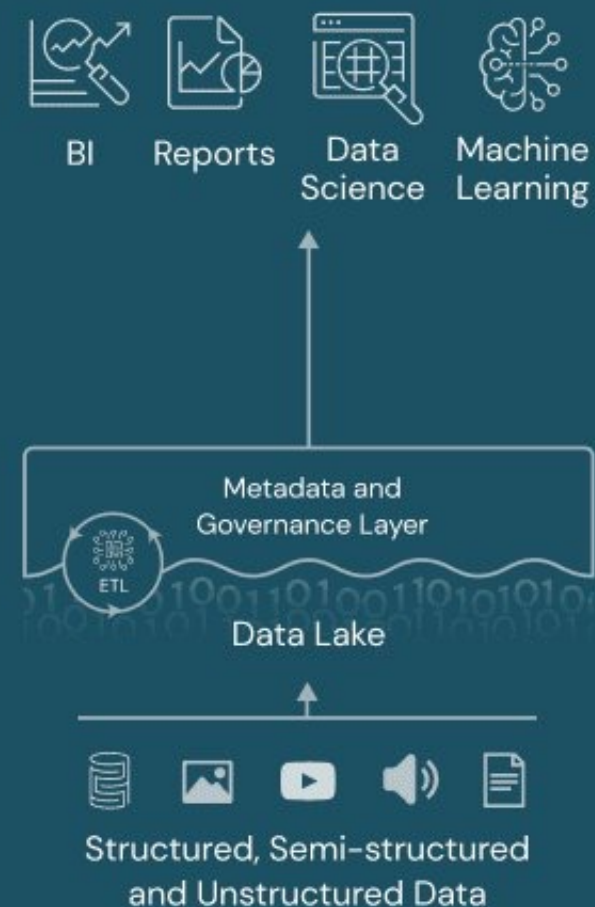
## Data Warehouse



## Data Lake



## Data Lakehouse





# Why lakehouse?



|                         | Data lake   | Data lakehouse  |
|-------------------------|---|---|
| Data format             | Open format   | Open format   |
| Types of data           | All types: Structured data, semi-structured data, textual data, unstructured (raw) data | All types: Structured data, semi-structured data, textual data, unstructured (raw) data |
| Data access             | Open APIs for direct access to files with SQL, R, Python and other languages            | Open APIs for direct access to files with SQL, R, Python and other languages            |
| Reliability             | Low quality, data swamp   | <b>High quality, reliable data with ACID transactions</b>                               |
| Governance and security | Poor governance as security needs to be applied to files                                | <b>Fine-grained security and governance for row/columnar level for tables</b>           |
| Performance             | Low   | <b>High</b>   |
| Scalability             | Scales to hold any amount of data at low cost, regardless of type                       | <b>Scales to hold any amount of data at low cost, regardless of type</b>                |
| Use case support        | Limited to machine learning   | <b>One data architecture</b> for BI, SQL and machine learning                           |
| Total cost of ownership | Low   | Low   |

# How?

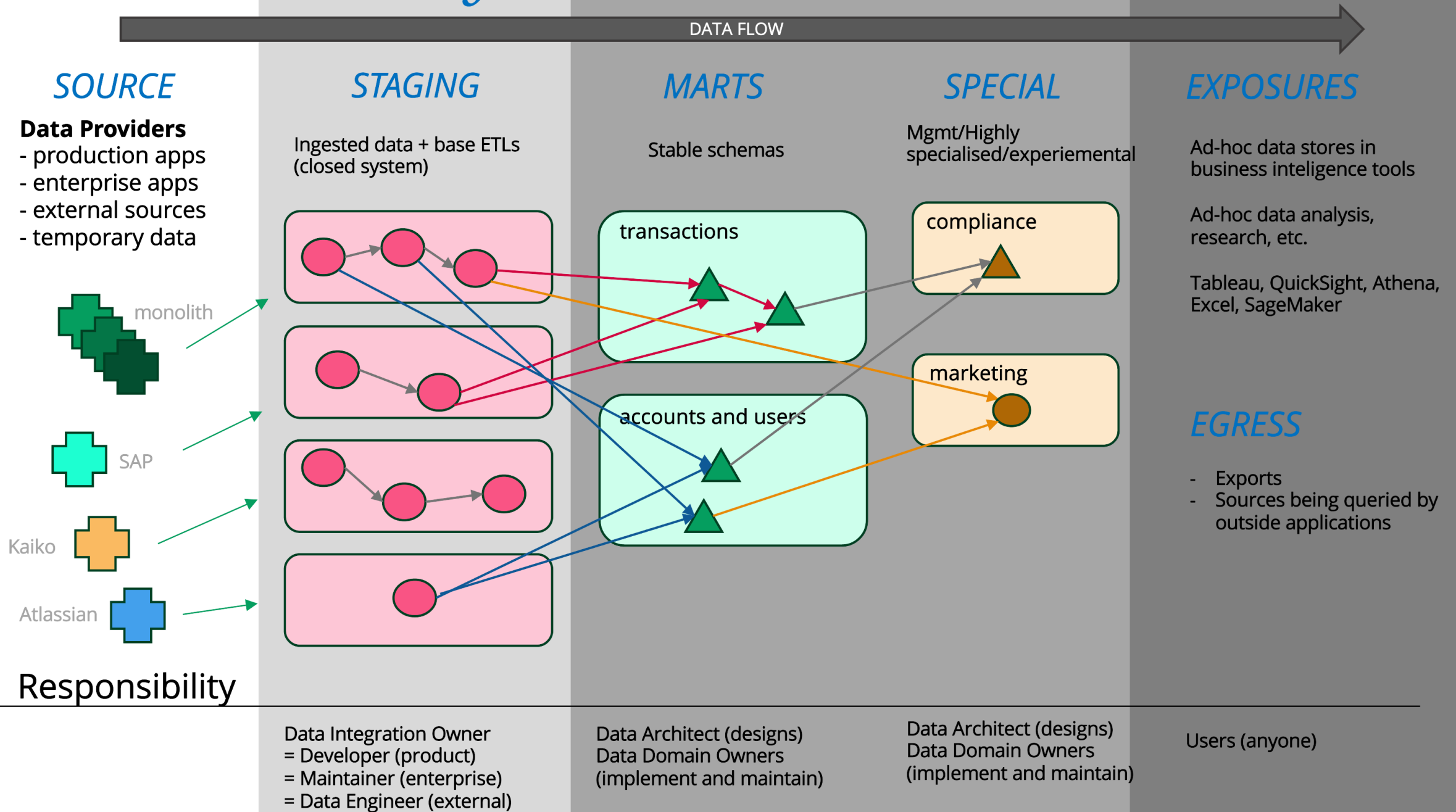


## Organisational changes

- **Shift left**
- Use **config-based** approach
- Introduce **end-to-end testing**
- Provide toolkits, **golden paths** and conventions.
- **Devs shouldn't think too much.**
- **Data governance** and ownership are maintained on all layers.

Provide **platform**

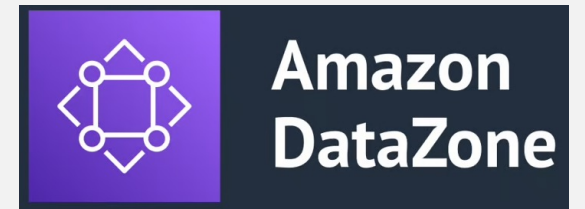
# Domains in Layers



# How?

## Technological changes

- Change:
  - AWS Athena → Iceberg
  - Inhouse ETL tooling → DBT (data build tool)
  - Add governance layer → (AWS Data Zone)



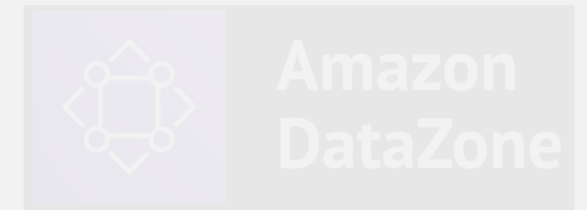
Provide **platform**



# How?

## Technological changes

- Change:
  - AWS Athena → **Iceberg**
  - Inhouse ETL tooling → DBT (data build tool)
  - Add governance layer → (AWS Data Zone)

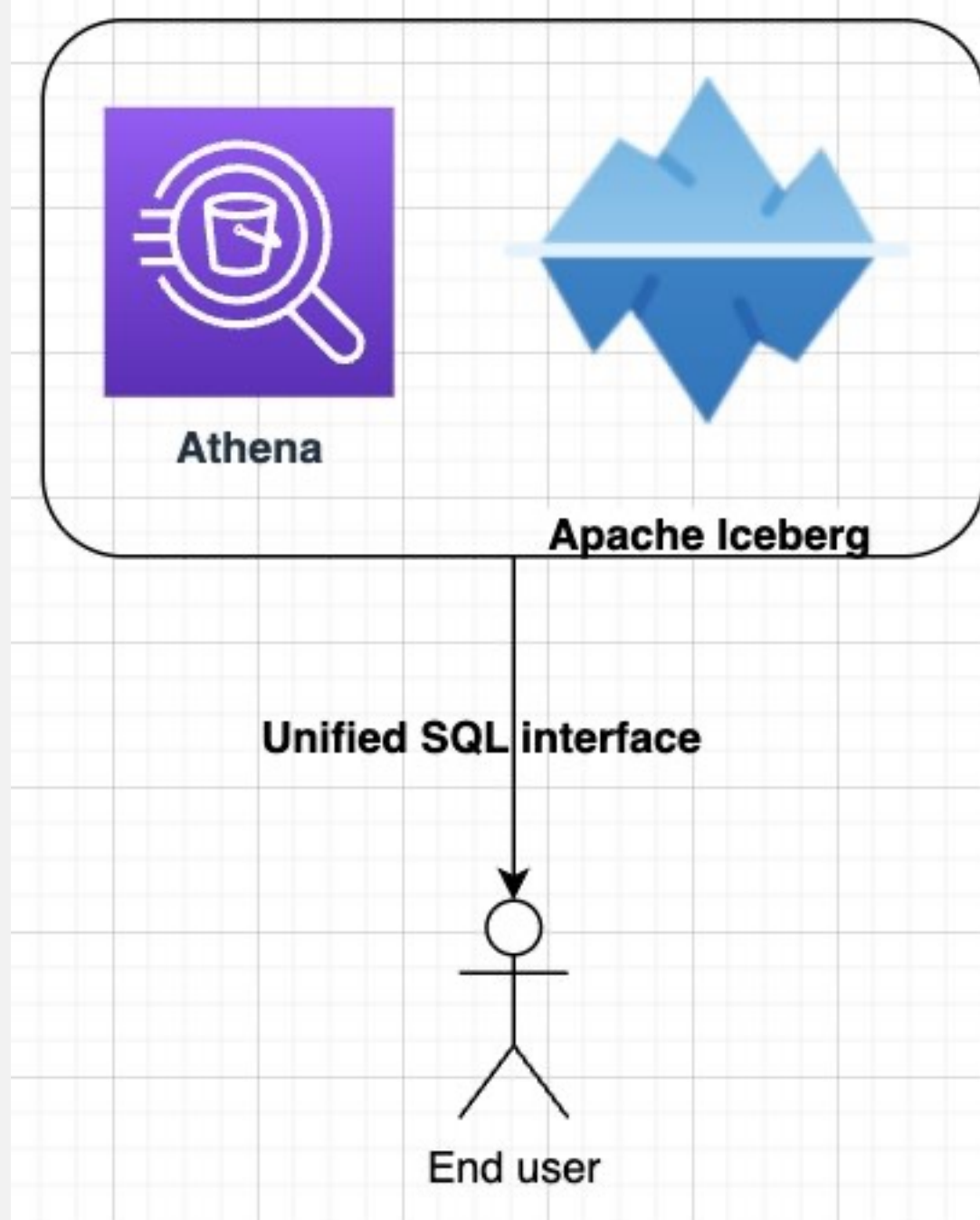


Provide **platform**

# From Athena to Iceberg

- ACID transactions
- CRUD operations
- Schema evolution
- Time travel
- Rollback
- Data Compaction
- Compatible with Athena
- Common SQL

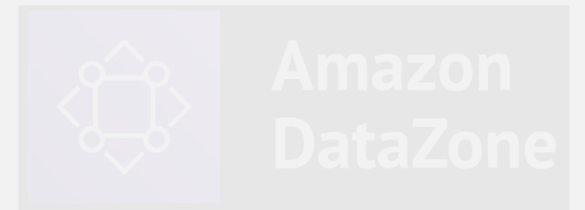
TL;DR – **Athena++ for Data Lakehouse**



# How?

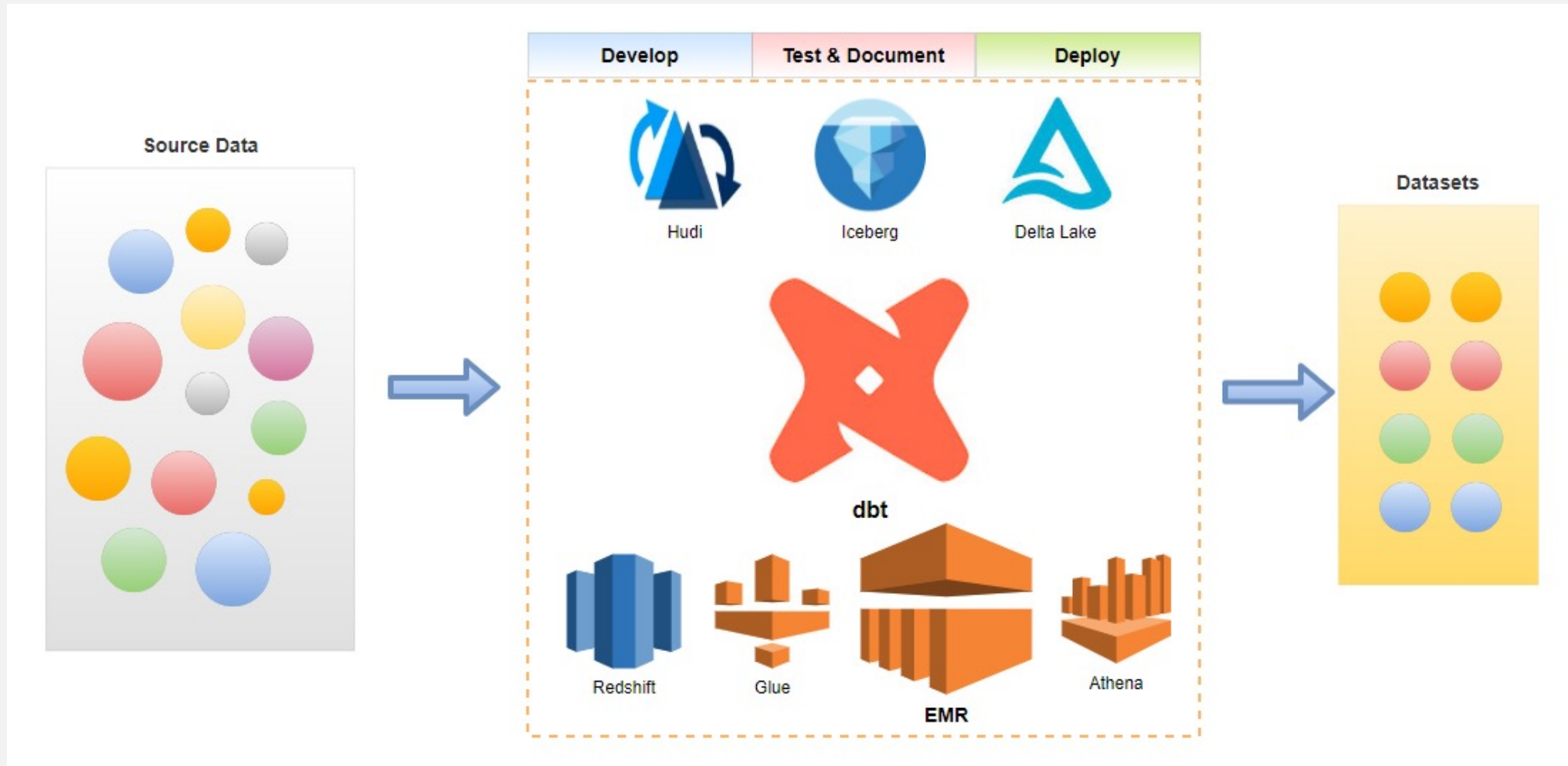
## Technological changes

- Change:
  - AWS Athena → Iceberg
  - Inhouse ETL tooling → **DBT (data build tool)**
  - Add governance layer → (AWS Data Zone)



Provide **platform**

# DBT - Data Build Tool



Implementing robust ETLs reduces to “writing SQL”



# DBT

= SE best practices + data

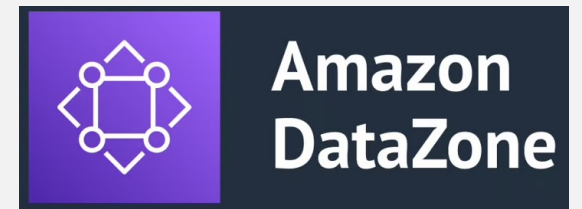
- **Declarative:** transform user-written code into raw SQL and run it against an engine
- **Git** workflow
- **Data lineage**
- **Testing**
- Support various **materialization** strategies
- In practice: “**SQL and Jinja**”

```
jaffle_shop
├── README.md
├── analyses
├── seeds
│   └── employees.csv
├── dbt_project.yml
├── macros
│   └── cents_to_dollars.sql
├── models
│   ├── intermediate
│   │   ├── finance
│   │   │   ├── _int_finance__models.yml
│   │   │   └── int_payments_pivoted_to_orders.sql
│   ├── marts
│   │   ├── finance
│   │   │   ├── _finance__models.yml
│   │   │   ├── orders.sql
│   │   │   └── payments.sql
│   │   └── marketing
│   │       ├── _marketing__models.yml
│   │       └── customers.sql
│   ├── staging
│   │   ├── jaffle_shop
│   │   │   ├── _jaffle_shop__docs.md
│   │   │   ├── _jaffle_shop__models.yml
│   │   │   ├── _jaffle_shop__sources.yml
│   │   │   ├── base
│   │   │   │   ├── base_jaffle_shop__customers.sql
│   │   │   │   └── base_jaffle_shop__deleted_customers.sql
│   │   │   ├── stg_jaffle_shop__customers.sql
│   │   │   └── stg_jaffle_shop__orders.sql
│   │   └── stripe
│   │       ├── _stripe__models.yml
│   │       ├── _stripe__sources.yml
│   │       └── stg_stripe__payments.sql
│   └── utilities
│       └── all_dates.sql
├── packages.yml
├── snapshots
├── tests
│   └── assert_positive_value_for_total_amount.sql
```

# How?

## Technological changes

- Change:
  - AWS Athena → Iceberg
  - Inhouse ETL tooling → DBT (data build tool)
  - Add governance layer → (**AWS Data Zone**)

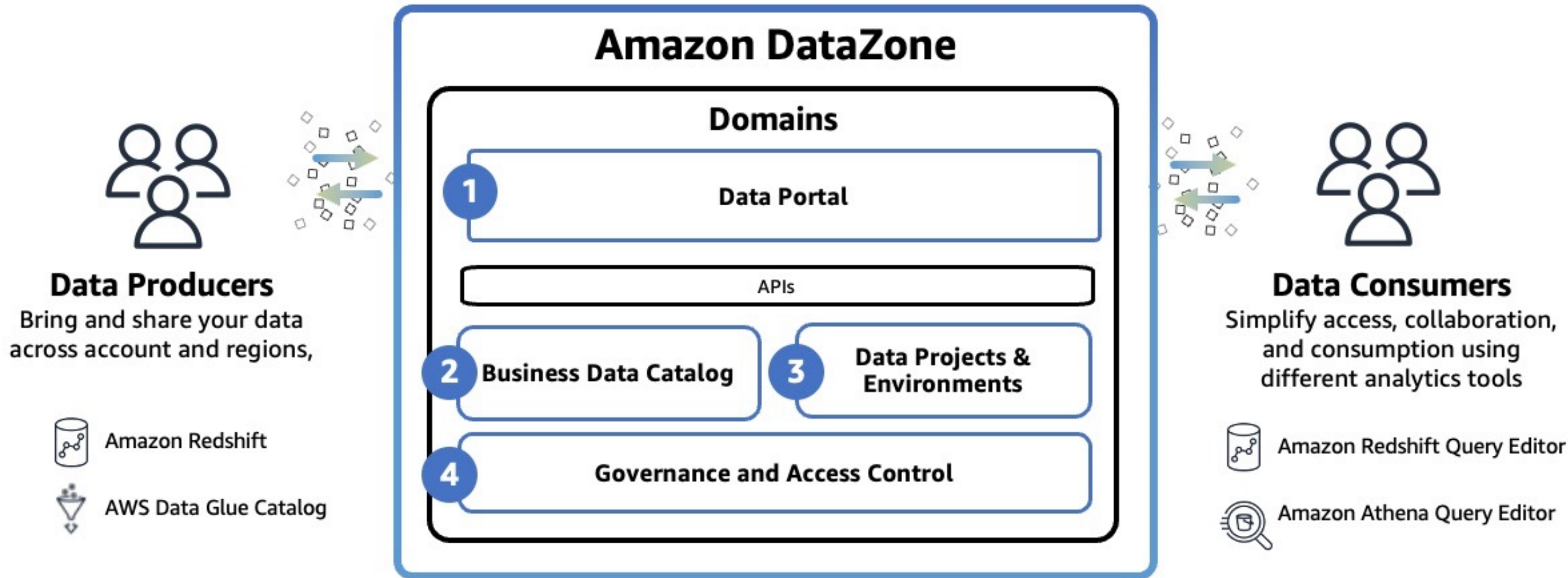


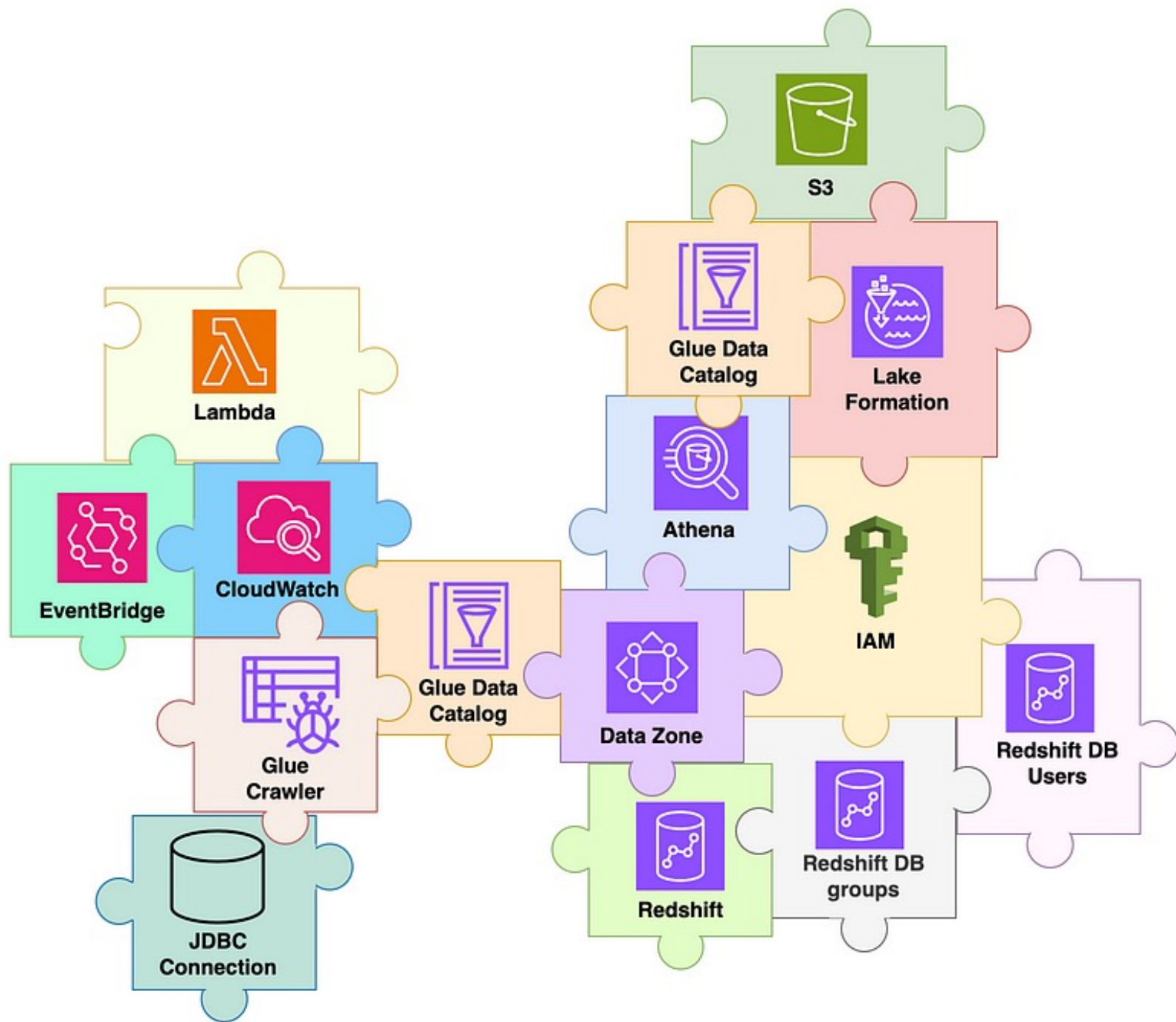
Provide **platform**

# AWS Data Zone to the rescue

Amazon's new offering for data governance

B

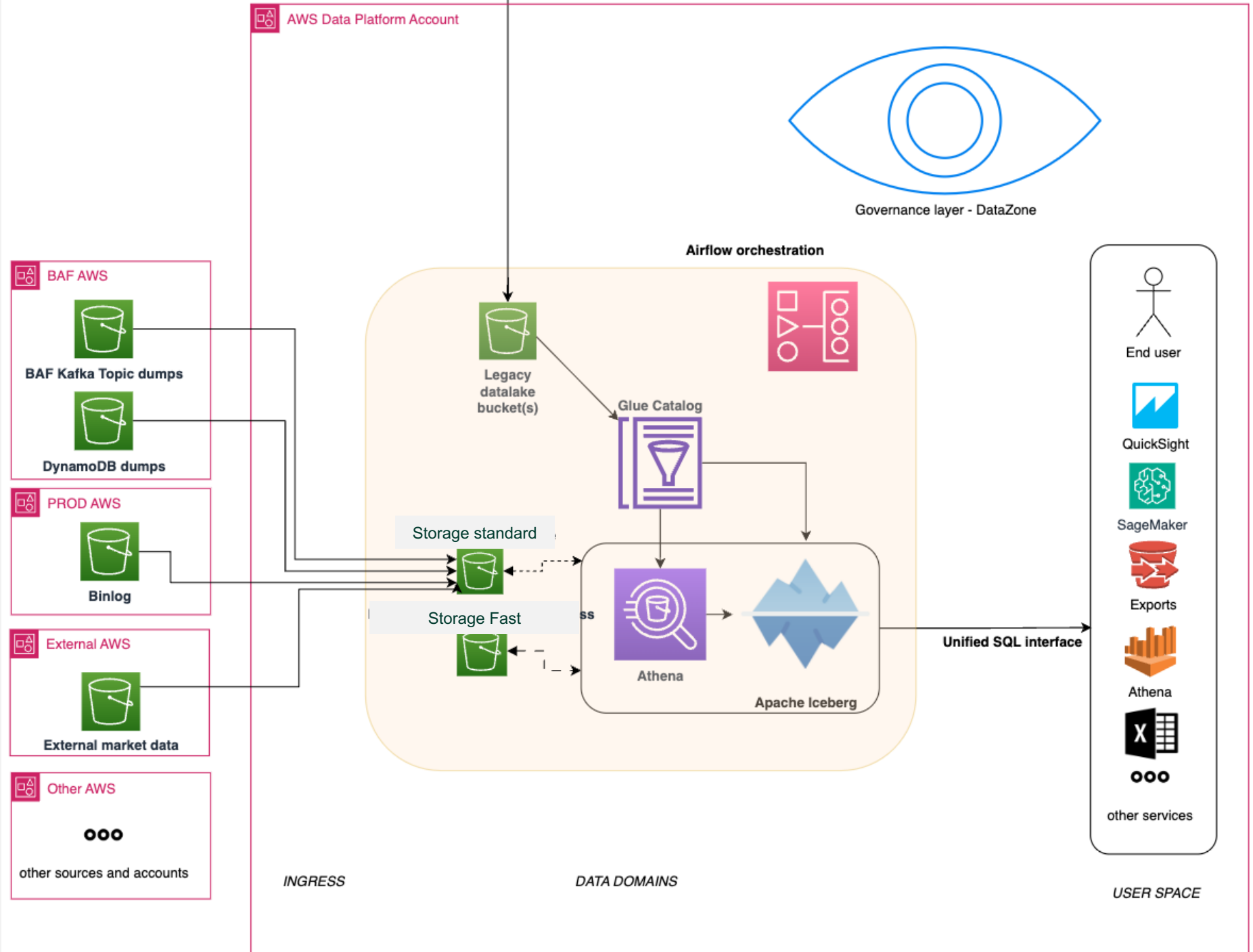






# Lakehouse Infrastructure

- Shift ingress left
- Apache Iceberg-
- All-seeing governance
- Liberated user space



# Highlights



- **Data swamp is salvageable**
- We **embraced legacy** data, we do not fight it (too much)
- We were **lucky** to have a sort of lakehouse already set up
- **Iceberg** is neat, **DBT** is neat
- **Golden paths**, conventions and governance is the name of the game

**Unchecked, at some point, every datalake will become “too big”  
for a single team**

# Challenges



- **Shifting left** and scaling
  - Technologically
  - Culturally
- Onboarding to **DBT**
  - Introducing software practices into traditionally non-software teams
- Adapting AWS **DataZone**
  - New product, pilot customers



Q & A