# mimacom

## Modern Data Engineering 101
## Von Data Warehouse bis Data Mesh

# Agenda

mimacom

# Julia Riedel



## Background
- Data Engineer
- Specialized on Azure / Databricks
- In D&A for 3,5 years

## At Mimacom
- Design and Implement new Data Platforms for Customers

# Thomas Konstantinides



## Background
- In the IT since 2002
- Backend / frontend / tech lead / hands-on architect
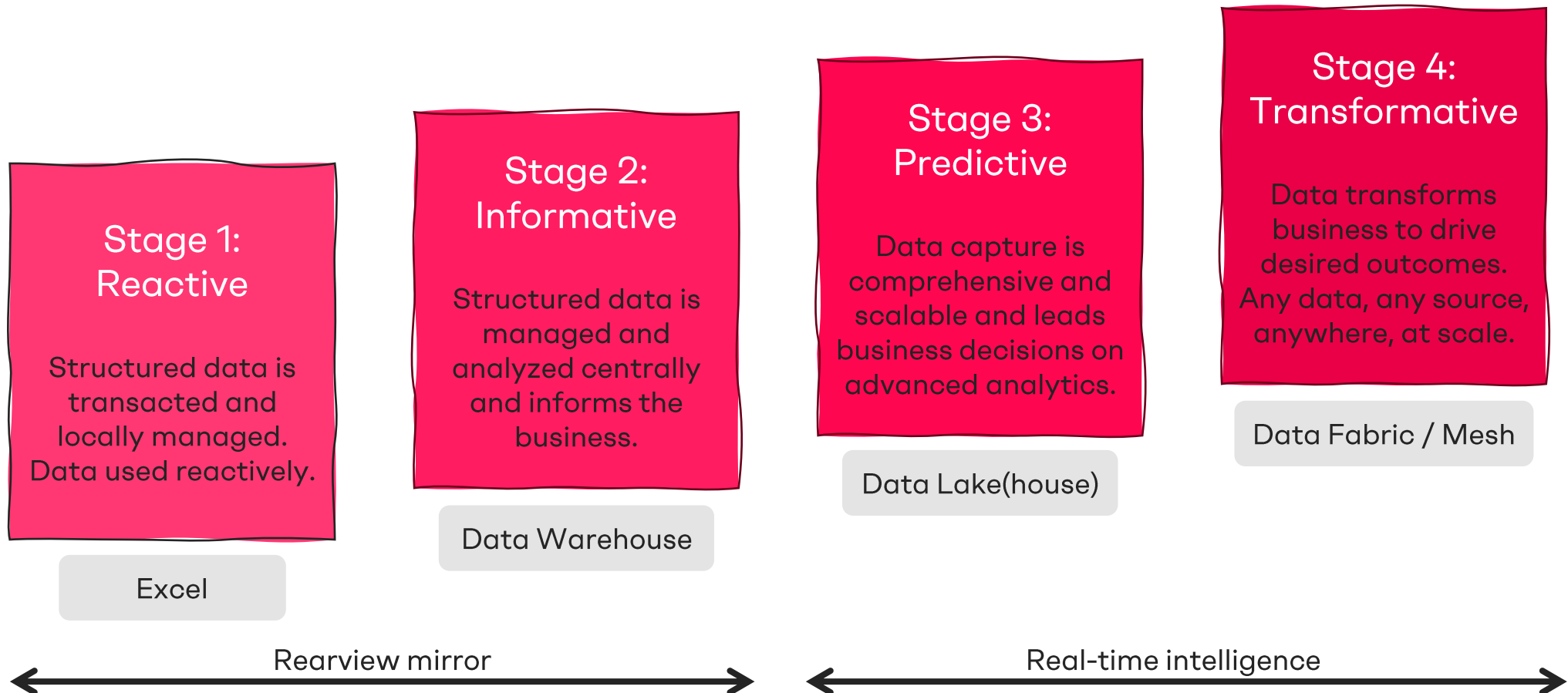- Now focus on data engineering

## At Mimacom
- Part of a customer's data lakehouse team
- Support transition to a data mesh and development of the governance layer

# Why should we care?

- o Software Development is changing

- o Data is becoming more important

- o Newer architectures bring data and software development closer together

- o A "normal" software developer product team will often also take care about data topics
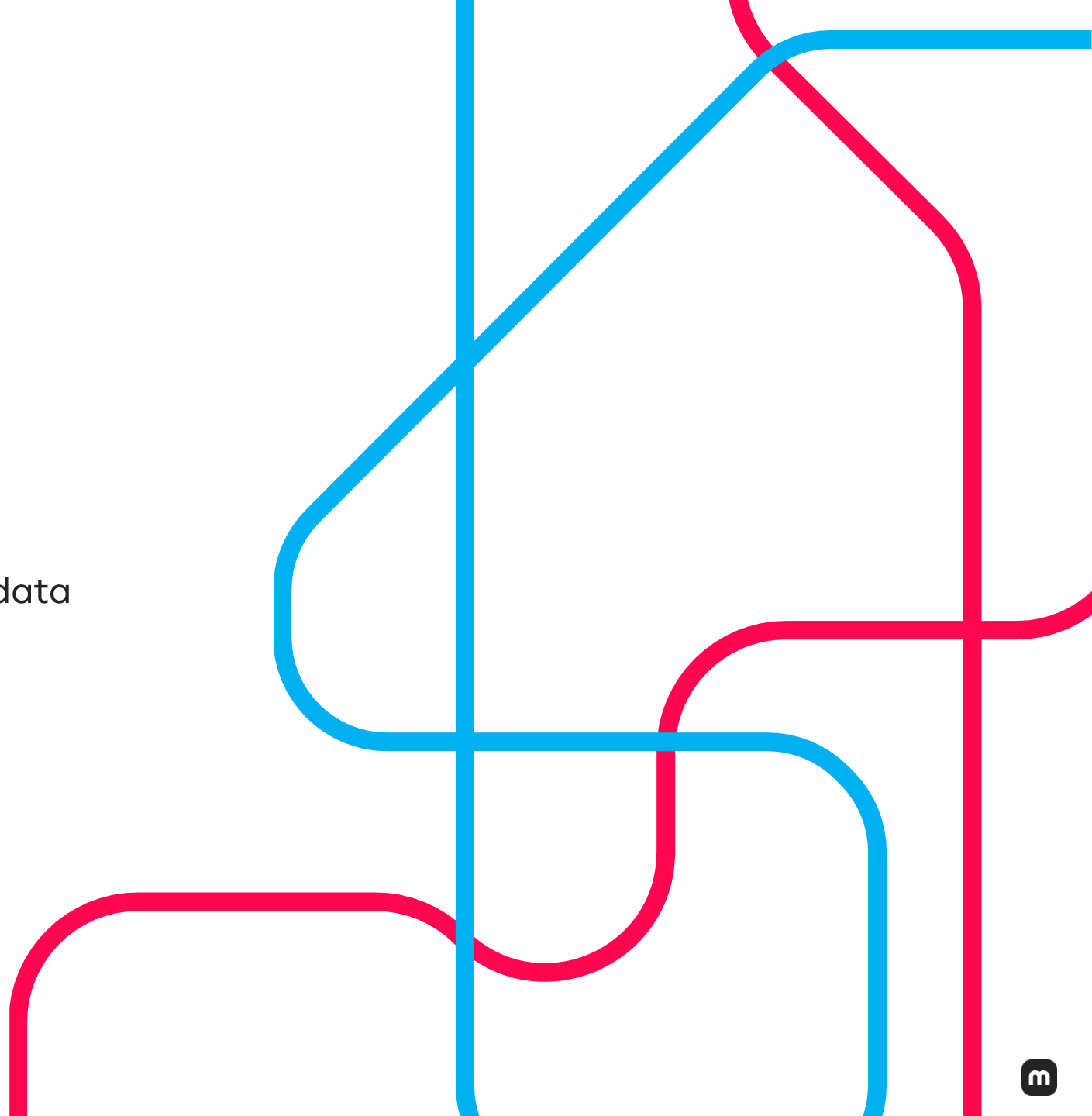
# Becoming a data-driven company

**Stage 1:
Reactive**

Structured data is transacted and locally managed. Data used reactively.

Excel

**Stage 2:
Informative**

Structured data is managed and analyzed centrally and informs the business.

Data Warehouse

**Stage 3:
Predictive**

Data capture is comprehensive and scalable and leads business decisions on advanced analytics.

Data Lake(house)

**Stage 4:
Transformative**

Data transforms business to drive desired outcomes. Any data, any source, anywhere, at scale.

Data Fabric / Mesh

← Rearview mirror →

← Real-time intelligence →

Source: Enterprise data maturity stages from "Deciphering Data Architectures" by James Serra

# OLTP vs. OLAP

- What's the difference?

- Why process the data?

- ETL: From operational data to analytical data

# What's the difference?

## OLTP
**Online Transactional Processing**

- Transactional
- Normalized
- Simple Queries (Read, Insert, Update)
- Current Data

## OLAP
**Online Analytical Processing**

- Analytical
- Denormalized
- Complex Queries (joins and aggregation)
- Historical Data

# What's the difference?

## OLTP

- Transactional
- Normalized
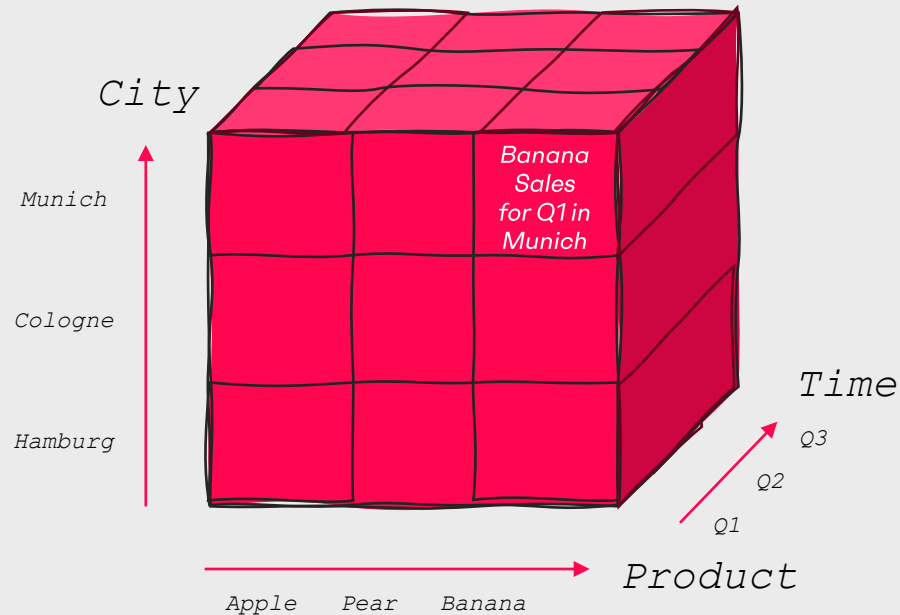- Simple Queries (Read, Insert, Update)
- Current Data

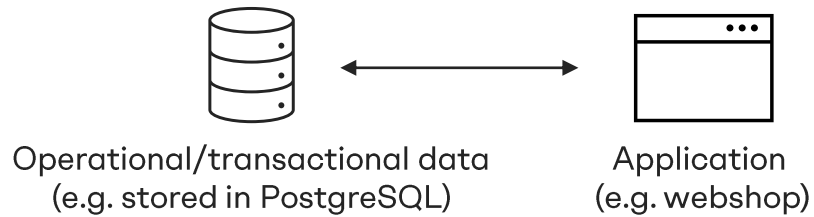City Table
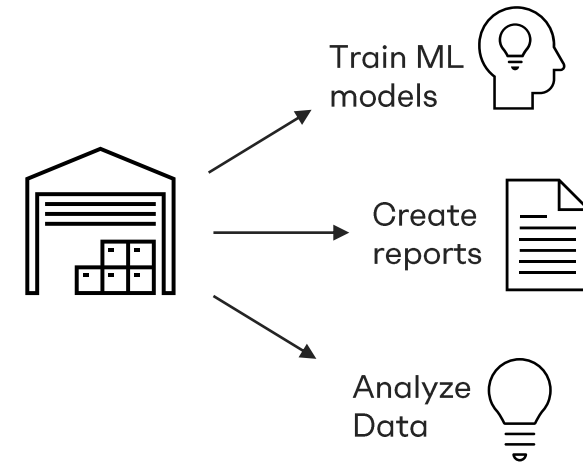
Product Table

Time Table

# What's the difference?



*City*

Munich

Cologne

Hamburg

*Banana Sales for Q1 in Munich*

*Time*

Q3

Q2

Q1

*Product*

Apple    Pear    Banana

## OLAP

- Analytical
- Denormalized
- Complex Queries (joins and aggregation)
- Historical Data

# Why process the data?

OLTP

OLAP



Operational/transactional data
(e.g. stored in PostgreSQL)

Application
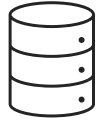(e.g. webshop)

Train ML models

Create reports

Analyze Data
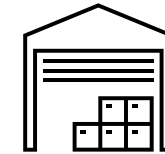
# ETL: From operational data to analytical data

ETL (= extract, transform, load)

Extract data from source system (Database, S3 Bucket, Kafka, API)

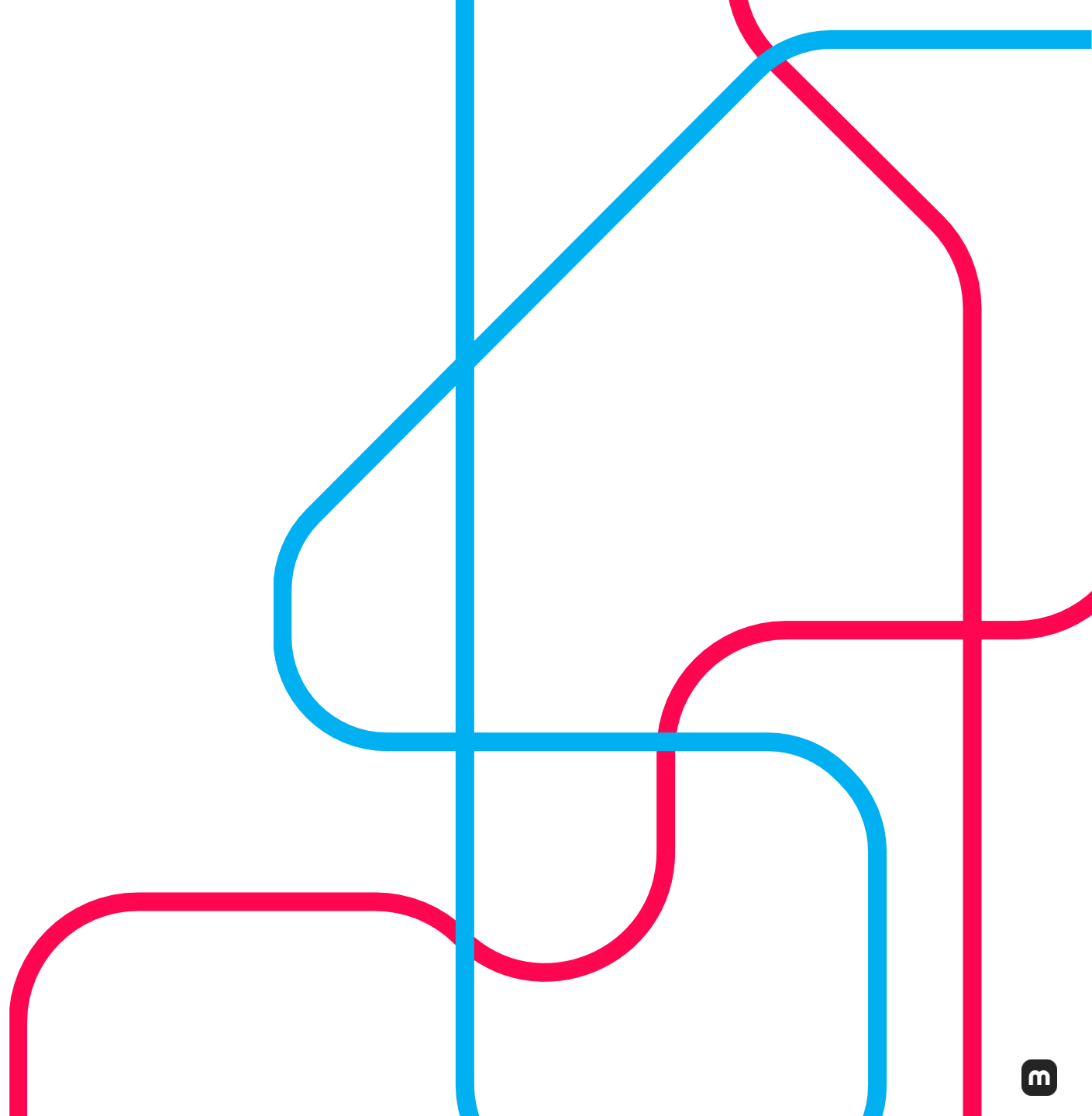Transform data to format suitable for usage

Load data into desired data storage

Analytics / ML /...

ELT ( = extract, load, transform)

# Data Storage Paradigms

- Data Warehouse

- Data Lake

- Data Lakehouse

- Technologies

# Data Warehouse

**Goal**: Support decision making process and reporting & visualizations

- Schema-on-write (ETL)
- Relational schema
- Structured data

## Advantages / Disadvantages

⭘ Optimized for downstream BI consumption
⭘ Pay for the peak of user load
⭘ No support for unstructured data
⭘ Limited use-case Support

# Data Lake

**Goal**: provide a cheap storage for data

- Schema-on-Read (ELT)
- Structured and unstructured data
- Data in generic and open file formats
- Often combined with a data warehouse

## Advantages / Disadvantages
- Low-cost storage systems with file API
- Lack of basic management features
- When used with an additional DWH double costs for storage
- It's hard to use the data in the lake

# Data Lakehouse

**Goal**: combine advantages of Data Warehouse and Data Lake

- Open direct-access data formats
- Open table formats like Delta Lake
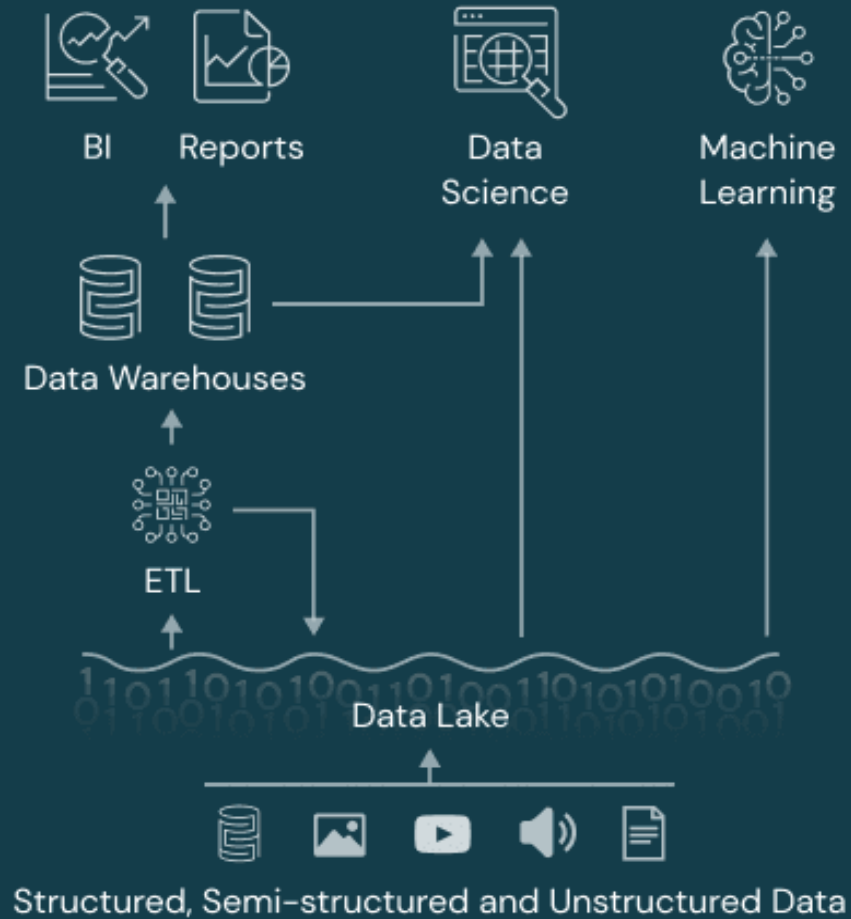- First-class support for machine learning and data science workloads

**Advantages / Disadvantages**
- Performance and management features of data warehouses
- Fast, direct I/O for advanced analytics workloads
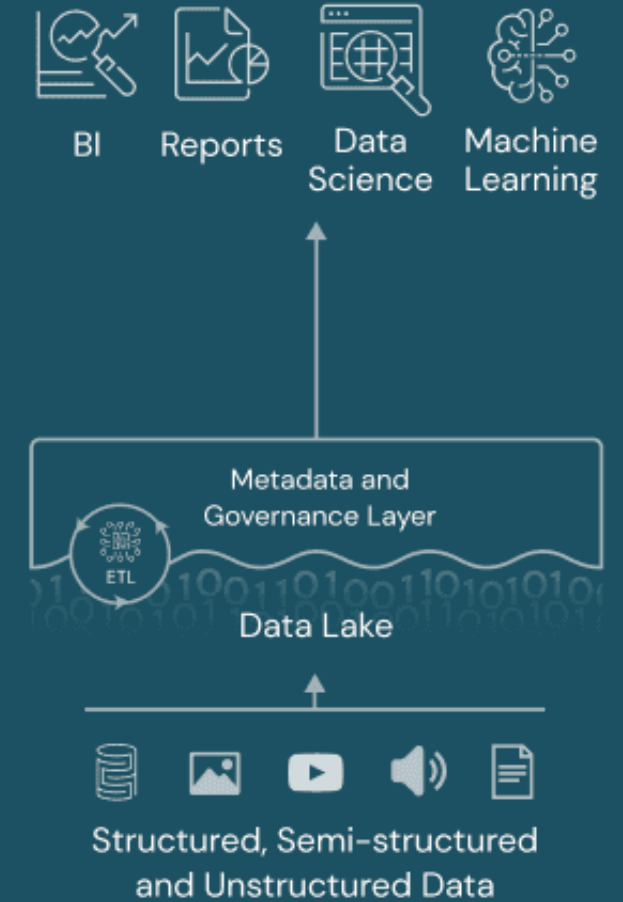- Less performance than DWH for some use cases

# Data Warehouse

BI  Reports

Data Warehouses

ETL

Structured Data

# Data Lake

BI  Reports  Data Science  Machine Learning

Data Warehouses

ETL

Data Lake

Structured, Semi-structured and Unstructured Data

# Data Lakehouse

BI  Reports  Data Science  Machine Learning

Metadata and Governance Layer

ETL

Data Lake

Structured, Semi-structured and Unstructured Data
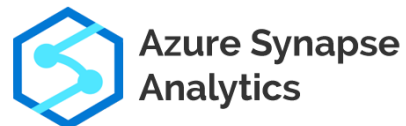
# Technologies

Multi language engine for
data engineering

- Supports Scala, Java,
  Python and R
- Can make use of SQL
- Integrates with most
  relevant frameworks and
  formats

# Technologies

### Azure Synapse Analytics

**Azure ELT / Lakehouse Solution**

- ELT Pipelines
- Serverless / Dedicated SQL Pools
- Read/Write to Data Lake

APACHE **Spark**™

**databricks**

**amazon** S3

❄ **snowflake**®

Microsoft® **SQL Server**®

# Technologies



**databricks**

**Lakehouse Solution**

- Data transformations via Spark
- Read/Write to Lakehouse
- Read/Write to Data Lake
- Combination with Azure Synapse possible

# Technologies

**amazon S3**

### Data Lake Solution

- Combination with Azure Synapse, AWS Glue possible
- Storage for Data Lake or Lakehouse

**APACHE Spark™**

**Azure Synapse Analytics**

**databricks**

**snowflake®**

**Microsoft® SQL Server®**

# Technologies

### snowflake®

**Lakehouse Solution**

- Data transformations via SQL
- All on one solution or only DWH
- Combination with Azure Synapse possible

APACHE Spark™

Azure Synapse Analytics

databricks

amazon S3

Microsoft® SQL Server®

# Technologies



### DWH Solution

- On-prem
- SQL Server Integration Services (SSIS) as ETL Solution
- traditional solution

# Data Architecture Patterns
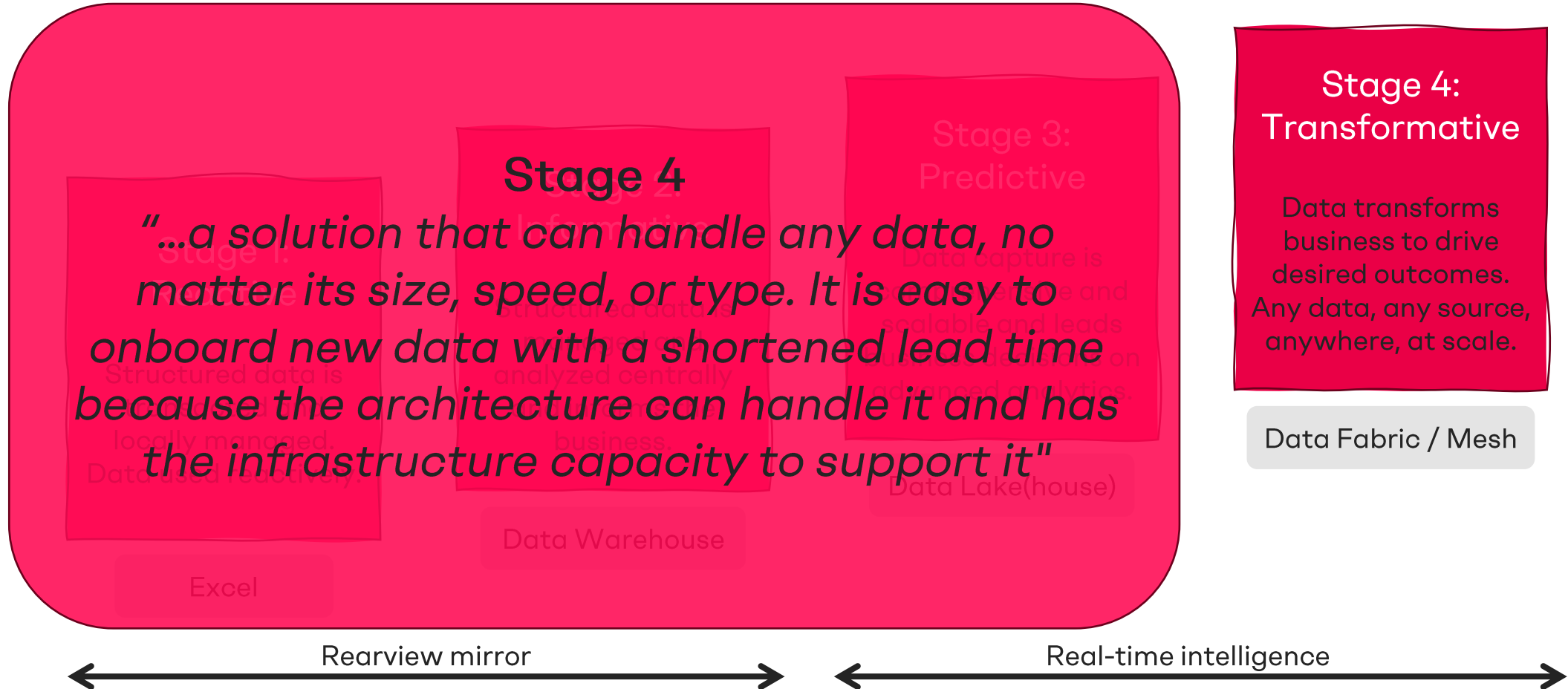
- Data Fabric
- Data Mesh

# Why more?

**Stage 4**

*"...a solution that can handle any data, no matter its size, speed, or type. It is easy to onboard new data with a shortened lead time because the architecture can handle it and has the infrastructure capacity to support it"*

Stage 1:

Stage 2:
Information

Stage 3:
Predictive

Excel

Data Warehouse

Data Lake(house)

### Stage 4: Transformative

Data transforms business to drive desired outcomes. Any data, any source, anywhere, at scale.

Data Fabric / Mesh

Rearview mirror ⟷

Real-time intelligence ⟷

Enterprise data maturity stages from "Deciphering Data Architectures" by James Serra

# Scalability issues

- Source data stored in different clouds and on-premise
- Different ways to access data sources
- Integrating new data sources takes long
- Hard to find needed data and know who can grant access
- Governance processes can create bottlenecks
- Data Warehouses cannot serve all use cases
- Data Lake can become "Data Swamp"

**Data Integration**

**Data Governance**

**Data Democratization**

**Scalability**

*Image from https://medium.com/@armandovazquez/navigating-the-waters-designing-a-data-lake-to-avoid-the-murky-depths-of-a-data-swamp-d67f5600c27c*

# Scalability issues
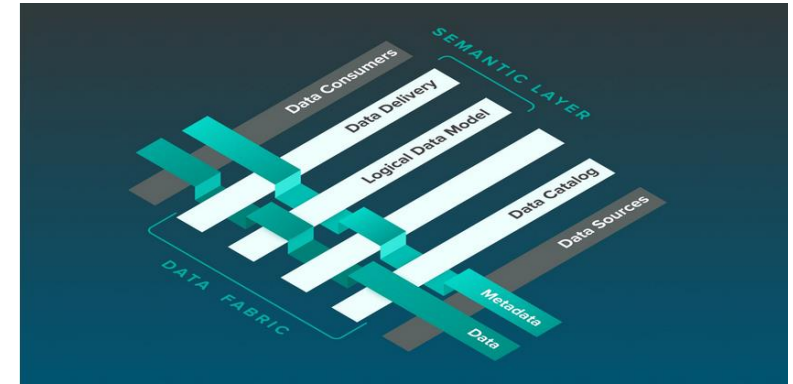
- Source data stored in different clouds and on-premise
- Different ways to access data sources
- Integrating new data sources takes long
- Hard to find needed data and kn~~~~~~~
- Governance pr~~~~~~~~~~~~~~checks
- D~~~~~~~~ ~~~~ serve all use cases
- ~~~~ Lake can become "Data Swamp"

Decreases business agility



**Data Integration**

**Data Governance**

**Data Democratization**

**Scalability**

*Image from https://medium.com/@armandovazquez/navigating-the-waters-designing-a-data-lake-to-avoid-the-murky-depths-of-a-data-swamp-d67f5600c27c*

# Data Fabric

# What is data fabric?

Data fabric is an <u>architecture</u> that facilitates the <u>end-to-end integration</u> of various data pipelines and cloud environments through the use of <u>intelligent and automated systems</u>.

*From https://www.ibm.com/topics/data-fabric*

...a <u>design concept</u> that serves as an <u>integrated</u> layer (fabric) of <u>data and connecting processes</u>. A data fabric utilizes <u>continuous analytics</u> over existing, discoverable and inferenced <u>metadata assets</u> to support the design, deployment and utilization of <u>integrated and reusable data across all environments</u>, including hybrid and multi-cloud platforms.

*From https://www.gartner.com/smarterwithgartner/data-fabric-architecture-is-key-to-modernizing-data-management-and-integration*



*From https://www.atscale.com/blog/what-is-a-data-fabric/*

# Takeaways

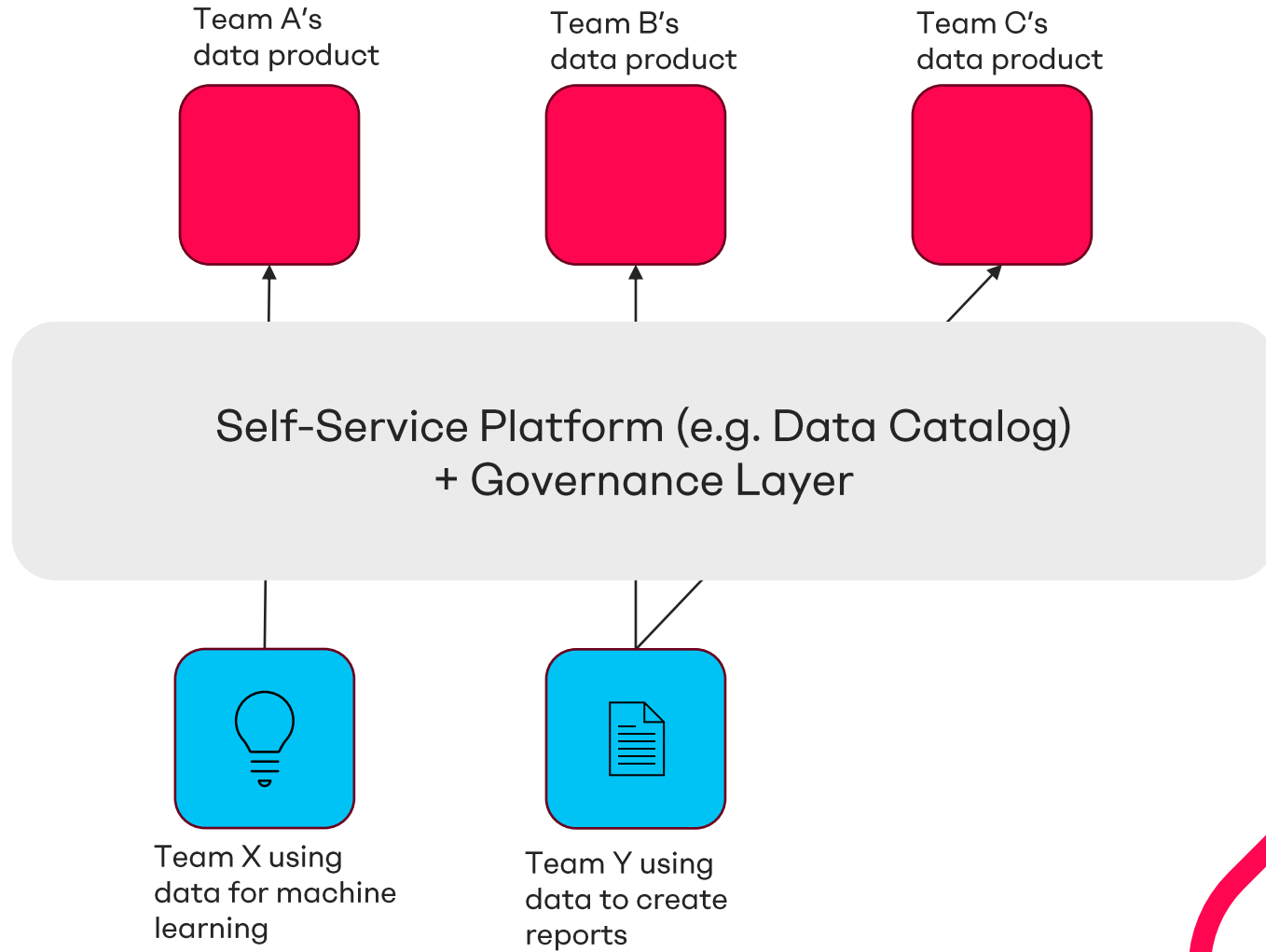## It needs more than just storing data and processing it

Seamless integration and governance across the entire data landscape

- o Integration over different data sources
- o Data virtualization
- o Unifying data access over APIs
- o Enabling data discovery
- o Meta data catalog
- o Centralized governance and security
- o Real-time support

Data Mesh

# Data Mesh idea

Team A's
data product

Team B's
data product

Team C's
data product

Self-Service Platform (e.g. Data Catalog)
+ Governance Layer

Team X using
data for machine
learning

Team Y using
data to create
reports

# Addressed issues

**Ownership: The source team provides the data product**

o Domain-driven, decentralized ownership

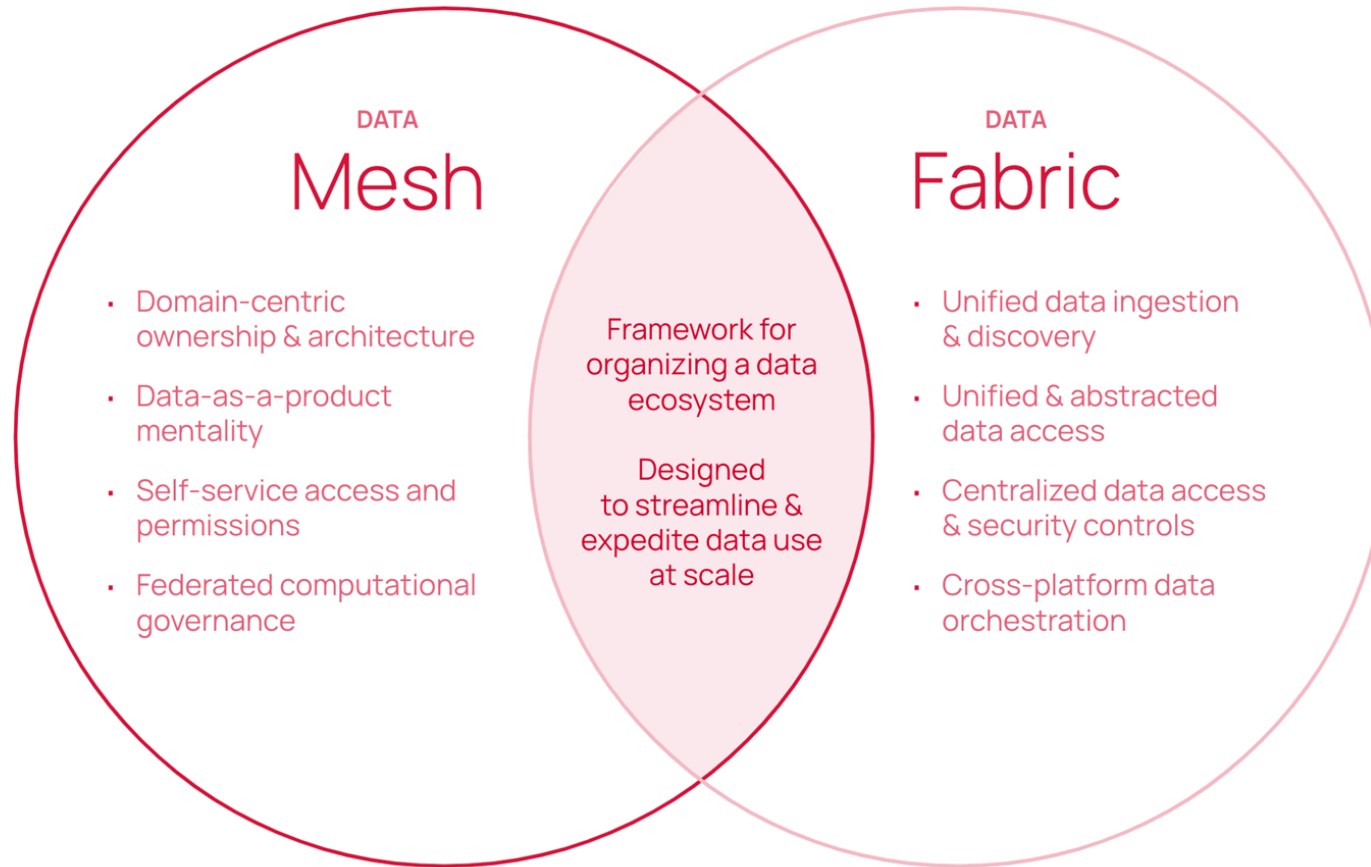o Management of data products within domain-specific teams

**Quality: The source team knows the data best**

o Product thinking

**Scaling: The central team can become a bottleneck**

o Self-service platform

o Federated Governance

# Data Mesh vs. Data Fabric

## DATA
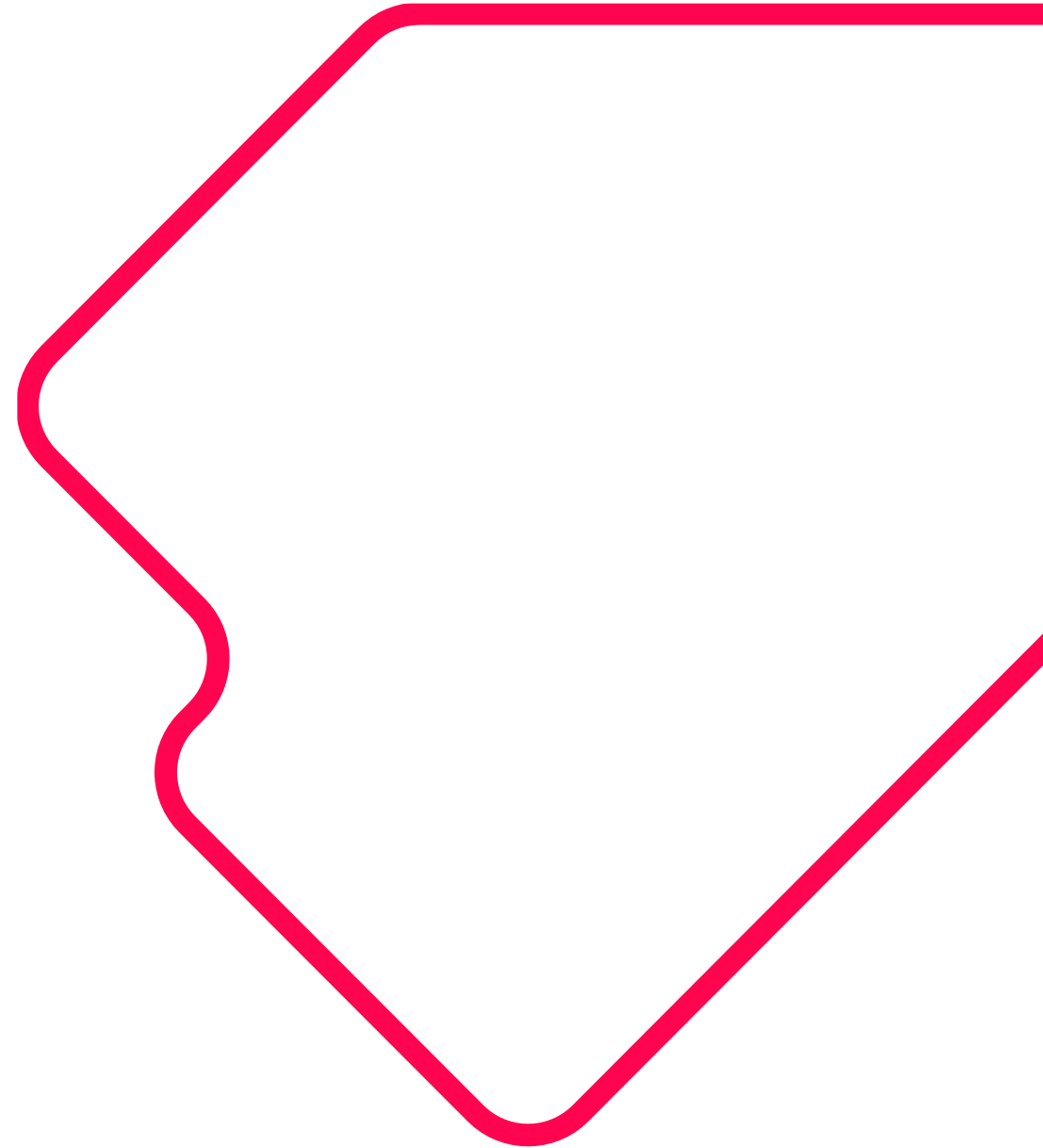### Mesh

- Domain-centric ownership & architecture

- Data-as-a-product mentality

- Self-service access and permissions

- Federated computational governance

Framework for organizing a data ecosystem

Designed to streamline & expedite data use at scale

## DATA
### Fabric

- Unified data ingestion & discovery

- Unified & abstracted data access

- Centralized data access & security controls

- Cross-platform data orchestration

*From  https://www.immuta.com/blog/data-mesh-vs-data-fabric/*

# What should we use?

It depends (of course) 😁

# Topics to consider

- Current organizational data maturity level

- The 6 vs of (big) data processing: volume, velocity, variety, variability, veracity, value

- Structured/semistructured/unstructured/binary data

- Number of data sources

- Experience and amount of data engineers/analysts

- Short-term goals vs. strategic vision

- You are (most likely) not Facebook, Google, Microsoft or NASA

# Apache Spark

# Tools for ETL and Data Analysis

Data
Frame
based

**Apache Spark**

Distributed processing
for large data sets

**Polars, Pandas**

Single node processing for
smaller amounts of data

SQL
based

**DuckDB, AWS
Athena**

Abstraction layer
allowing to use plain SQL

Integrated

**Azure Data Flows,
AWS Glue**

Abstractions on top of Spark,
which make it easier to use

# Demo

# Kommt uns an unserem Stand besuchen!

mimacom

Hegel Foyer – Stand 27