# An Evolution of Data Platform Architectures in Azure

## Lambda, Kappa, Delta, Data Mesh
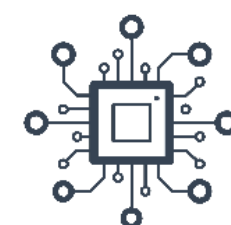
$\lambda$ $\kappa$ $\delta$

# Paul Andrew | Technical Architect in Azure CoE

**Microsoft** MVP Most Valuable Professional

avanade

Mr Paul Andrew Consulting Ltd
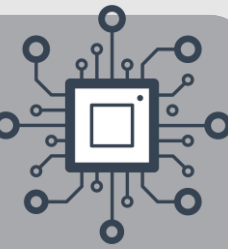
@MrPaulAndrew    In/MrPaulAndrew    MrPaulAndrew.com    c/MrPaulAndrew
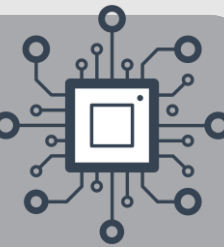
# What is the answer to life, the universe and everything?

**Answer:**
42 ❌

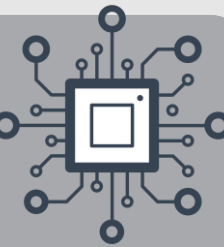

**Answer:**
It depends! ✅

# What is big data?

**Answer:**
It depends! ✓

**Answer:**
"Any data that you cannot process
   in the time that you have/want ✓
     using the technology you have."

*- Buck Woody*

*@BuckWoodyMSFT*

**V**olume
**V**elocity
**V**ariety
**V**eracity
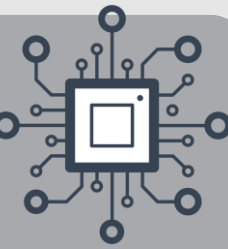**V**alue

# What is the goal of our data solutions?
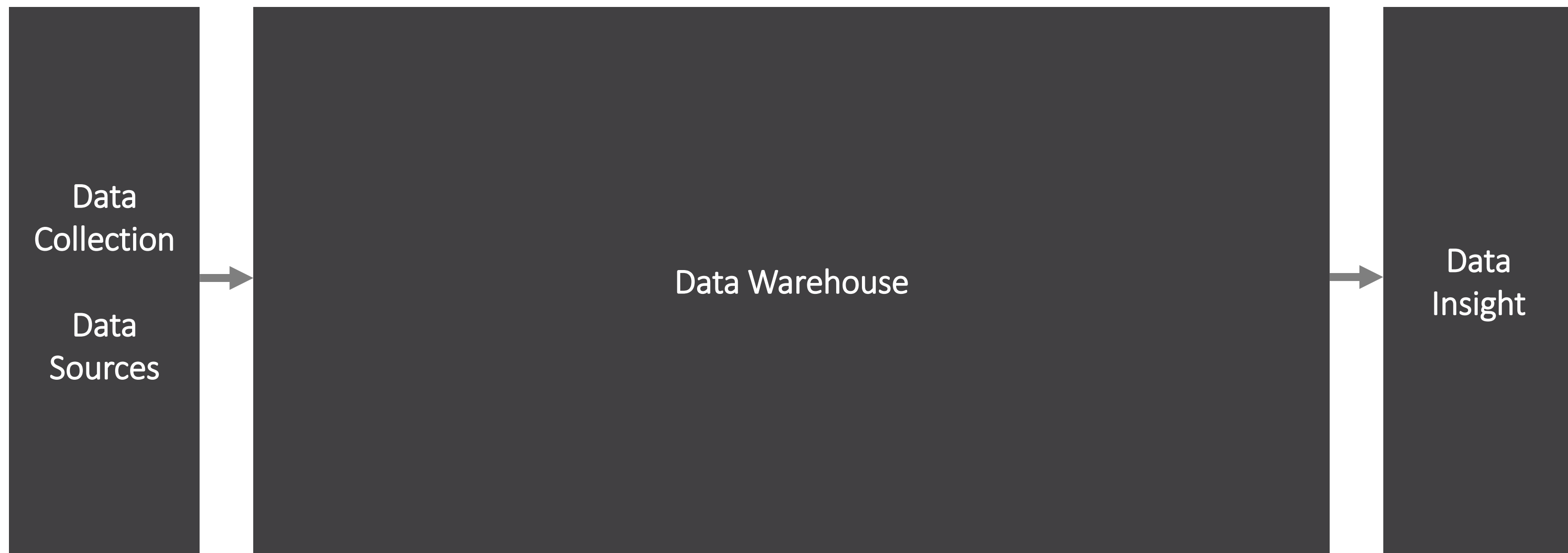
Data
Collection

Data
Sources

*Paul's Magic Box -*
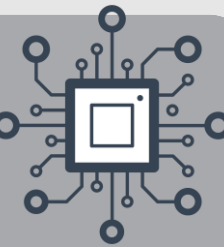*From the Hogwarts School of*
*Witches & Wizardry*

Data
Insight

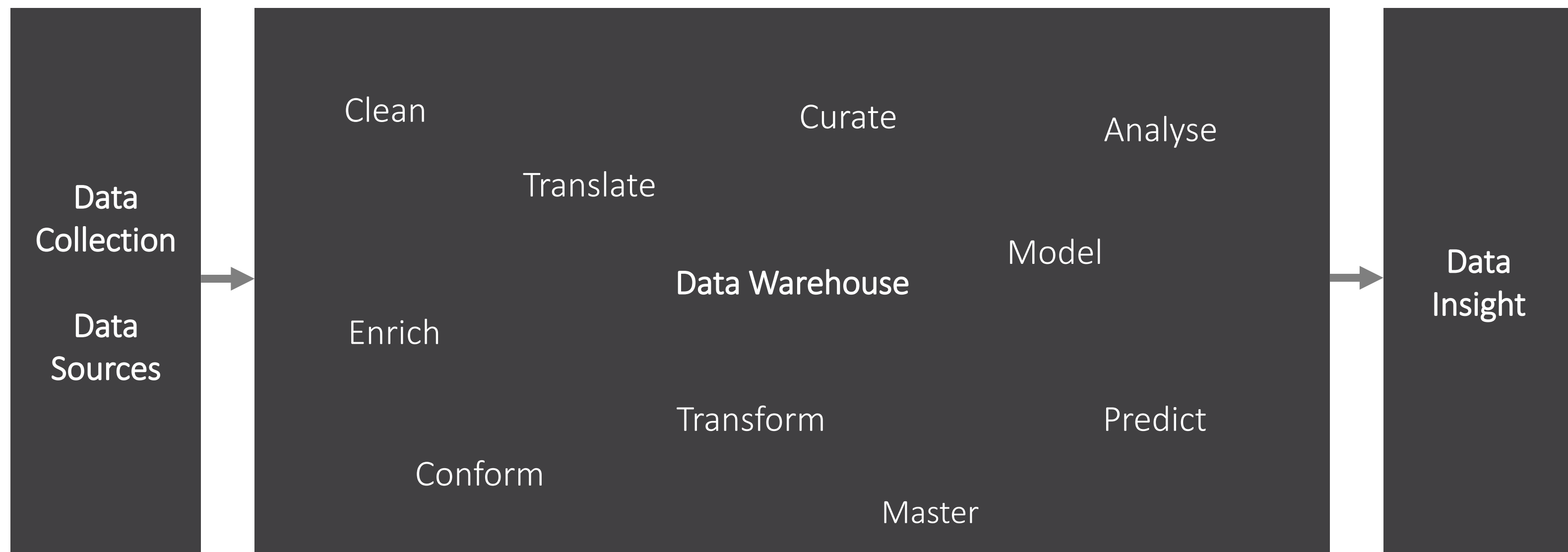*Data = Information = Knowledge = Power*

# What is the goal of our data solutions?

| Data Collection<br><br>Data Sources | → | Data Warehouse | → | Data Insight |

*Data = Information = Knowledge = Power/Insights*

# What is the goal of our data solutions?

**Data Collection**

**Data Sources**

→

Clean

Translate

Curate

Analyse

Model

**Data Warehouse**

Enrich

Transform

Predict

Conform

Master

→

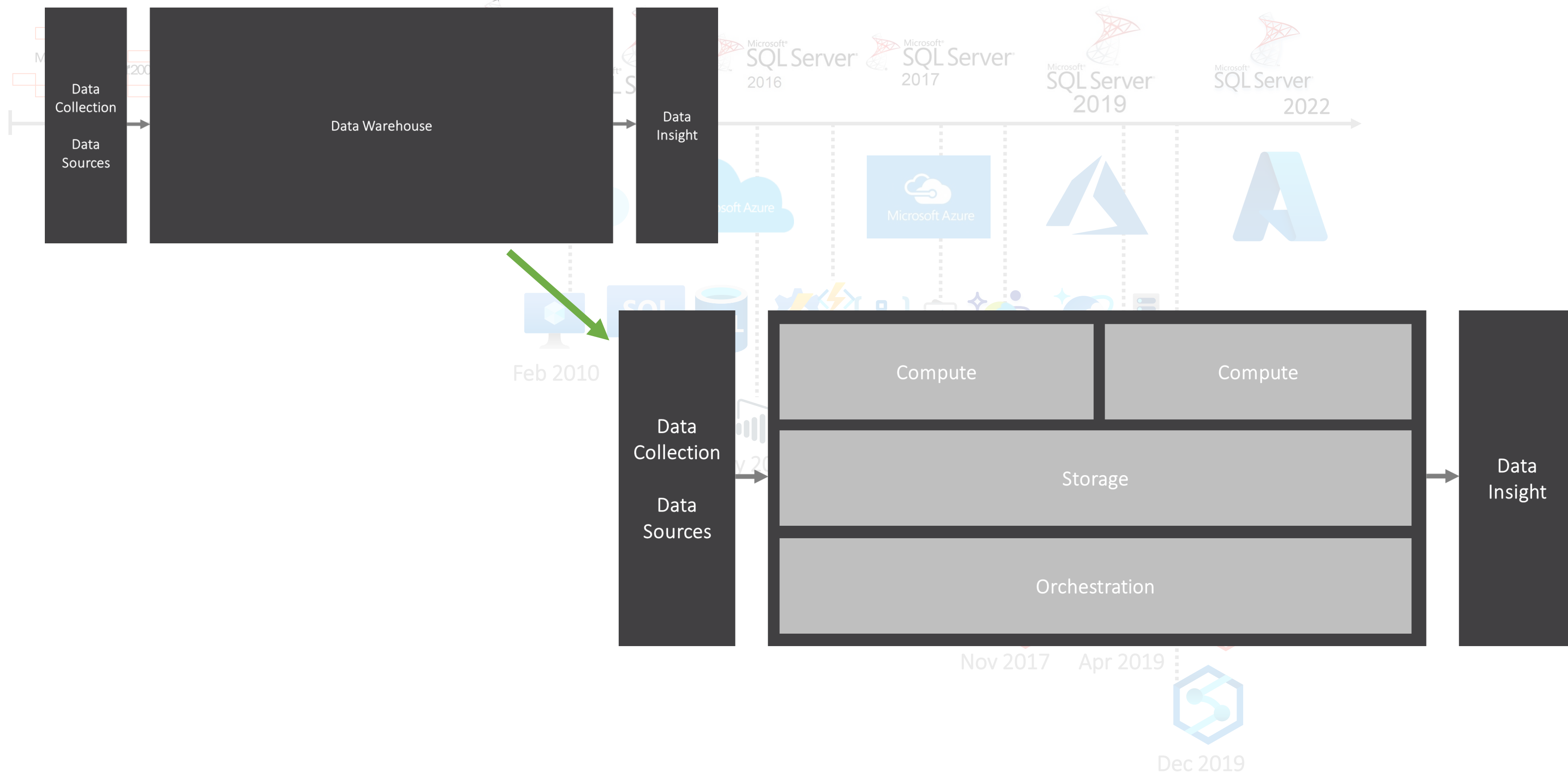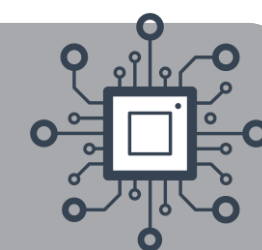**Data Insight**

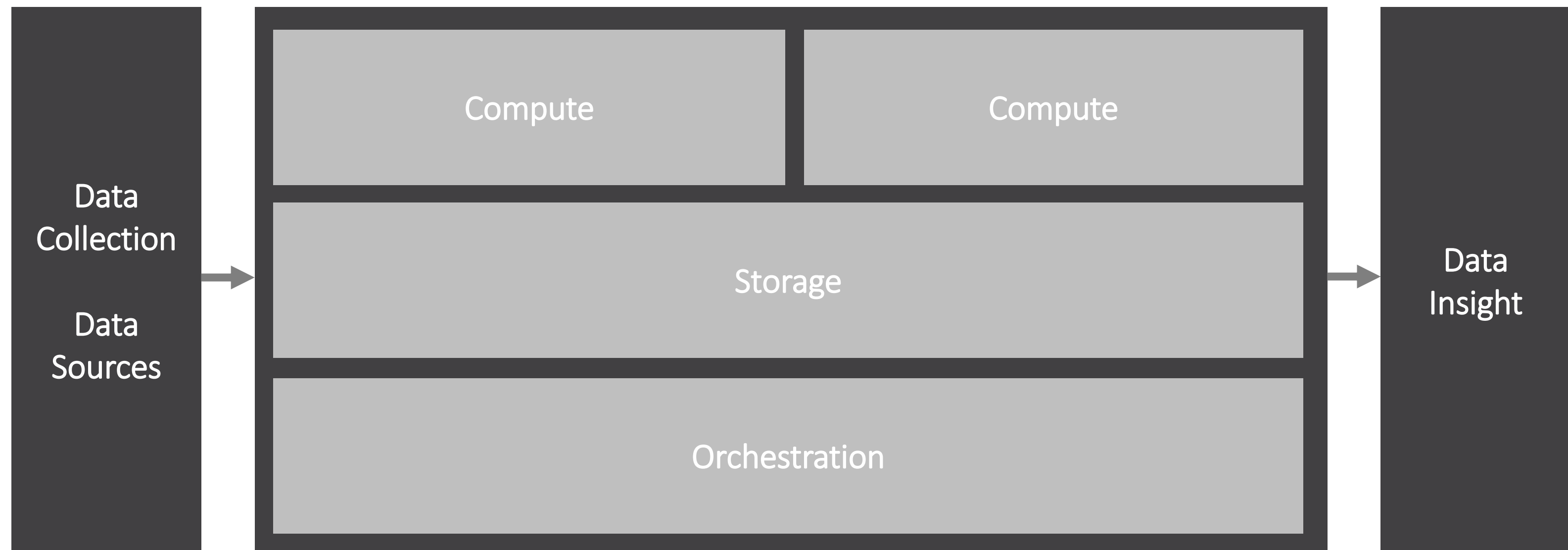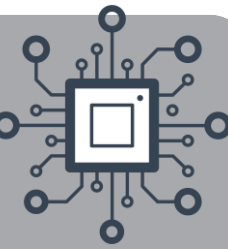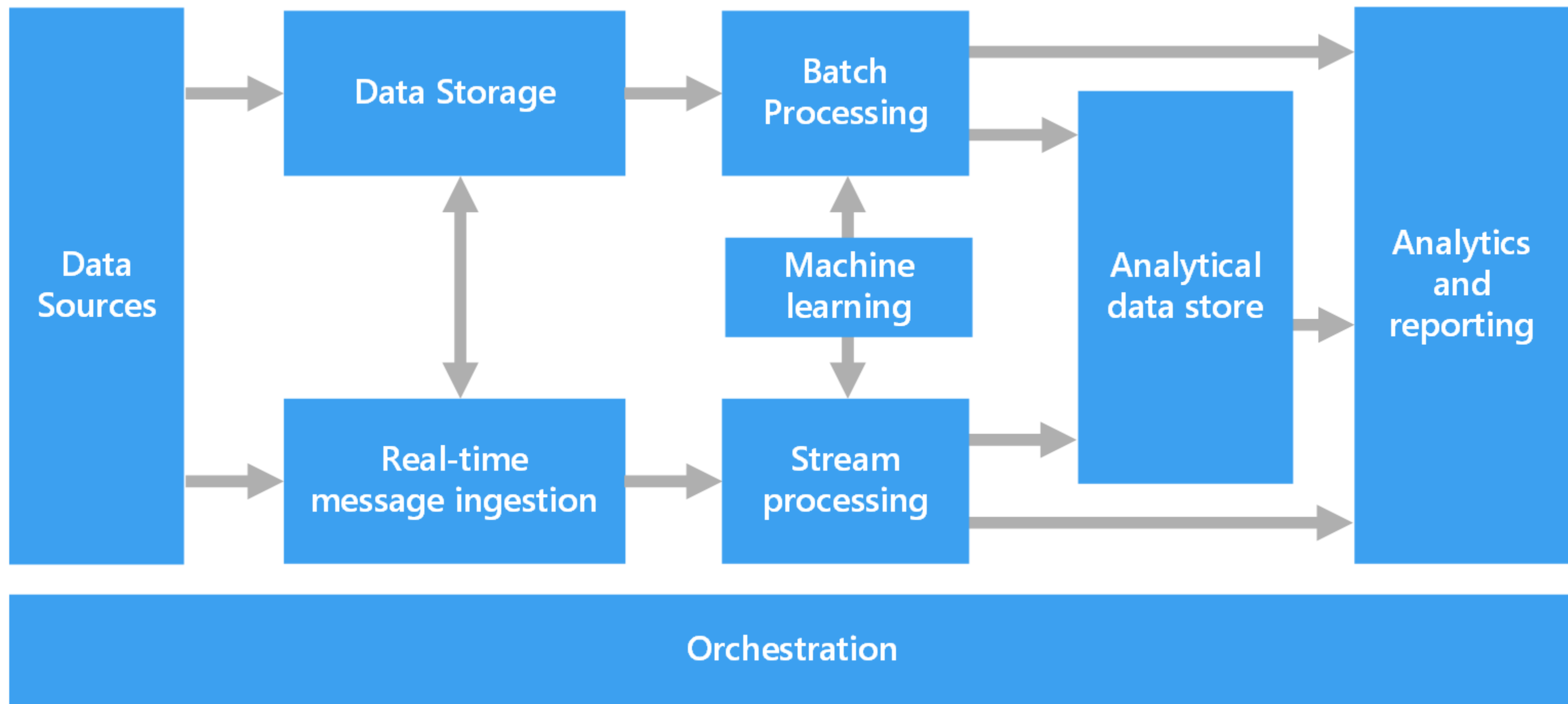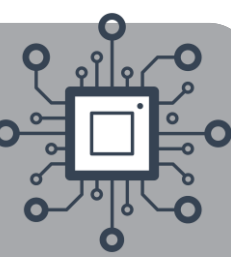*Data = Information = Knowledge = Power/Insights*

# An Evolution of Data Platforms

# An Evolution of Data Platforms

# A Reference Architecture?

| Data Collection / Data Sources | | Data Insight |
|---|---|---|
| → | Compute / Compute / Storage / Orchestration | → |

# Microsoft's Components of a ~~Big~~ Data Architecture



Diagram showing: Data Sources → Data Storage → Batch Processing; Data Sources → Real-time message ingestion → Stream processing; Machine learning connects to Batch Processing and Stream processing; Batch Processing and Stream processing → Analytical data store → Analytics and reporting; Orchestration spans across the bottom.
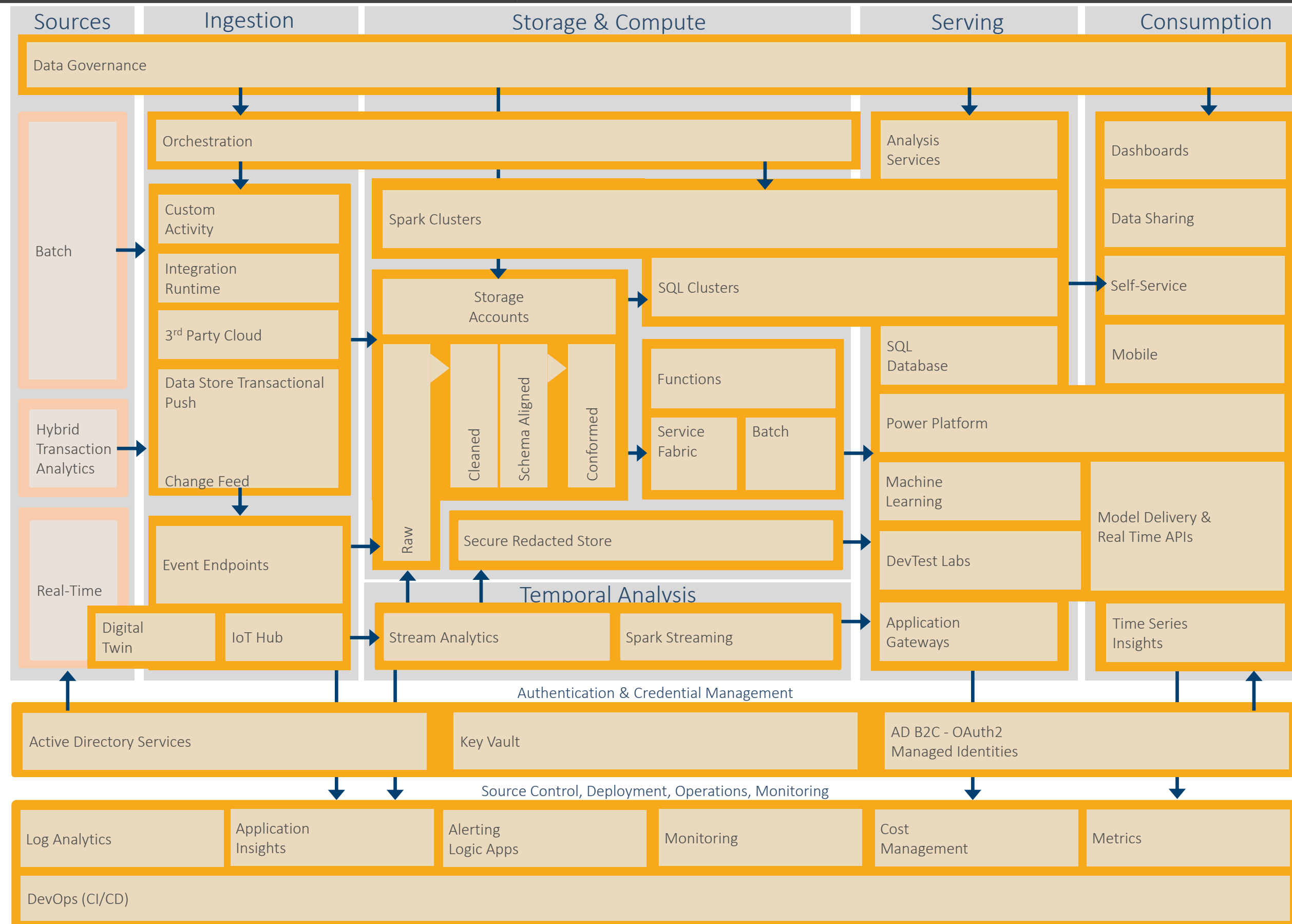
# A Logical Data Architecture

# A Logical Data Architecture



Left reference diagrams:

- Data Collection / Data Sources → Data Warehouse → Data Insight
- Data Collection / Data Sources → Compute / Compute / Storage / Orchestration → Data Insight
- Data Sources → Data Storage → Batch Processing / Machine learning → Analytical data store → Analytics and reporting; Real-time message ingestion → Stream processing; Orchestration

Main architecture — columns: Sources | Ingestion | Storage & Compute | Serving | Consumption

**Data Governance**

**Sources**
- Batch
- Hybrid Transaction Analytics
- Real-Time

**Ingestion**
- Orchestration
- Custom Activity
- Integration Runtime
- 3rd Party Cloud
- Data Store Transactional Push — REST
- Change Feed
- Event Endpoints
- Digital Twin
- IoT Hub

**Storage & Compute**
- Spark Clusters — Apache Spark
- Storage Accounts
- Raw / Cleaned / Schema Aligned / Conformed
- Secure Redacted Store — SQL
- SQL Clusters
- Functions
- Service Fabric / Batch
- Temporal Analysis
  - Stream Analytics
  - Spark Streaming — Apache Spark

**Serving**
- Analysis Services
- SQL Database — SQL
- Power Platform
- Machine Learning
- DevTest Labs
- Application Gateways

**Consumption**
- Dashboards
- Data Sharing
- Self-Service
- Mobile
- Model Delivery & Real Time APIs
- Time Series Insights

**Authentication & Credential Management**
- Active Directory Services
- Key Vault
- AD B2C - OAuth2 Managed Identities

**Source Control, Deployment, Operations, Monitoring**
- Log Analytics
- Application Insights
- Alerting Logic Apps
- Monitoring
- Cost Management
- Metrics

**DevOps (CI/CD)**

Legend:
- Concepts
- Resources
- Layers

avanade

# Delta* Lake

Delta.io

# Databases

**D**ata**B**ase
**M**anagement
**S**ystem

**A**tomicity
**C**onsistency
**I**solation
**D**urability



- Metadata
- Schema Updates
- Self-describing
- Integrity constraints
- Multiple Abstracted Views
- Multi Process Handling
- Compute and Storage Collaboration
- Indexing
- Row Level 'Upserts'
- Data Structure
- Transactional Processing
- Redundancy

SQL
My

**O**nline

**L**ine

**T**ransactional

**P**rocessing

**O**ffline

**A**nalytical

**T**ransactional

**P**rocessing

Application Data

**E**xtract
**T**ransform
**L**oad

Data Collection

Data Sources

Data Warehouse

Data Insight

# Data Lakes



Self-describing

Metadata

Schema Updates

Integrity constraints

Multiple Abstracted Views

Multi Process Handling

Compute and Storage Collaboration

Indexing

Row Level 'Upserts'

Data Structure

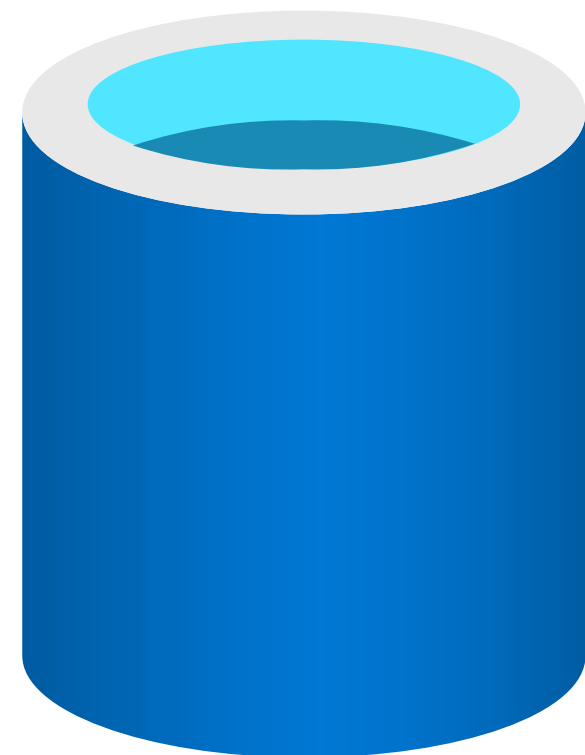Transactional Processing

Redundancy

**V**olume
**V**elocity
**V**ariety
**V**eracity
**V**alue

## Data Lakes are good, but they still lack some of the basic <u>ACID</u> functionality needed for data processing.

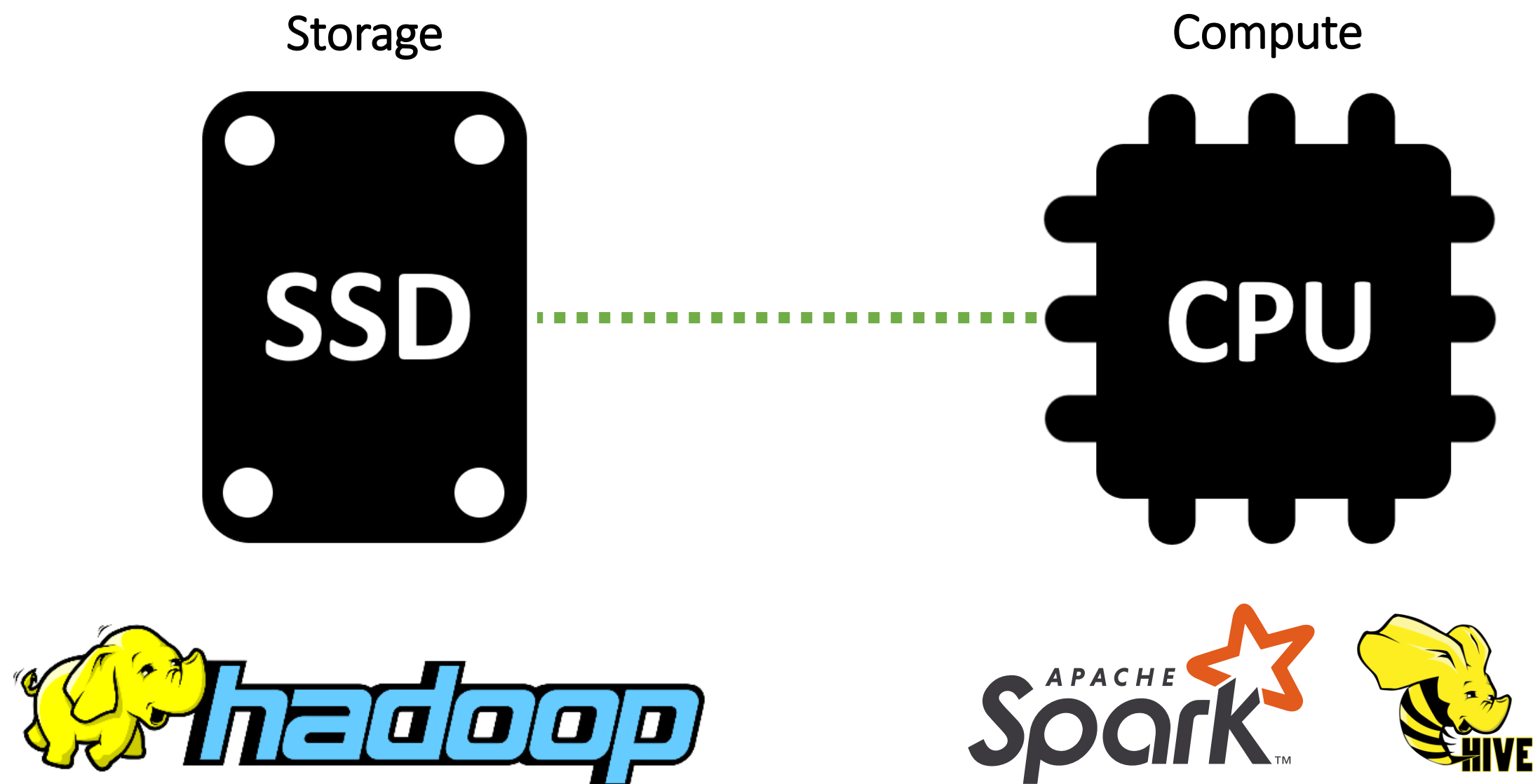We are/were trying to use Data Lakes for everything (to replace Databases).

VS

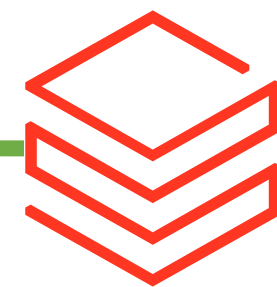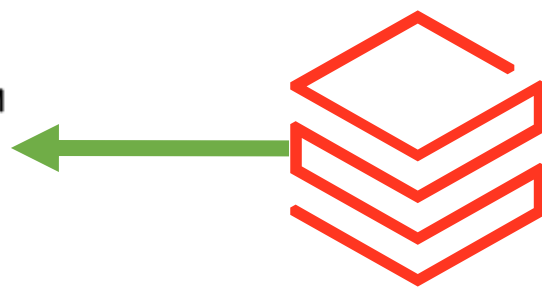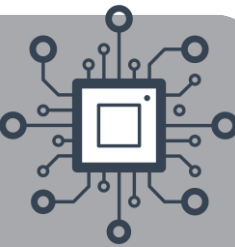| | |
|---|---|
| Scales Up | Scales Out |
| Natural Home for Structured Data | Any Data Structure |
| Storage Limits | No Storage Limits |
| Transactional Resilience | No Transactional Handling |
| Storage & Compute Coupled | Storage & Compute Decoupled |

Data Lakes are good, but they still lack some of the basic <u>ACID</u> functionality needed for data processing.

We are/were trying to use Data Lakes for everything (to replace Databases).

 VS 

| | |
|---|---|
| Scales Up | Scales Out |
| Natural Home for Structured Data | Any Data Structure |
| Storage Limits | No Storage Limits |
| Transactional Resilience | No Transactional Handling |
| Storage & Compute Coupled | Storage & Compute Decoupled |

# Just enable ACID transactional support for Data Lakes...

Storage

Compute



Storage & Compute ~~Decoupled~~ Working Together Again As Friends!

*February 2019*

https://delta.io
https://github.com/delta-io/delta

*April 2019*

Open Source

Synapse Analytics

Data Factory Data Flows

HD Insight Spark

Microsoft Implementation

Databricks Implementation

*"Delta Lake is an open-source storage layer that brings ACID transactions to Apache Spark™ and big data workloads."*
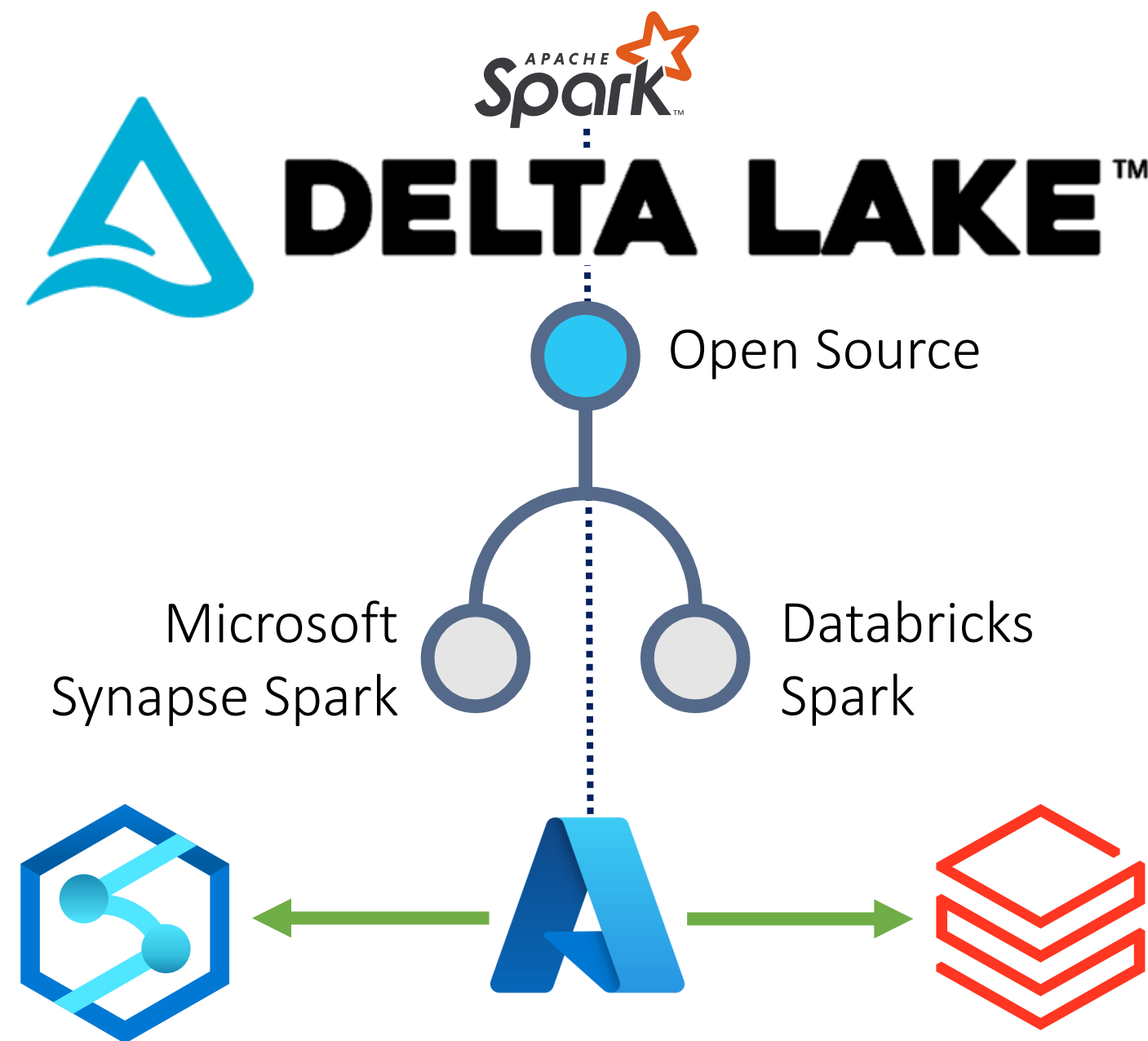
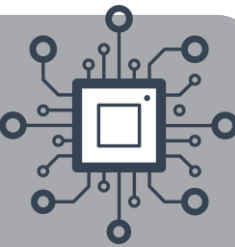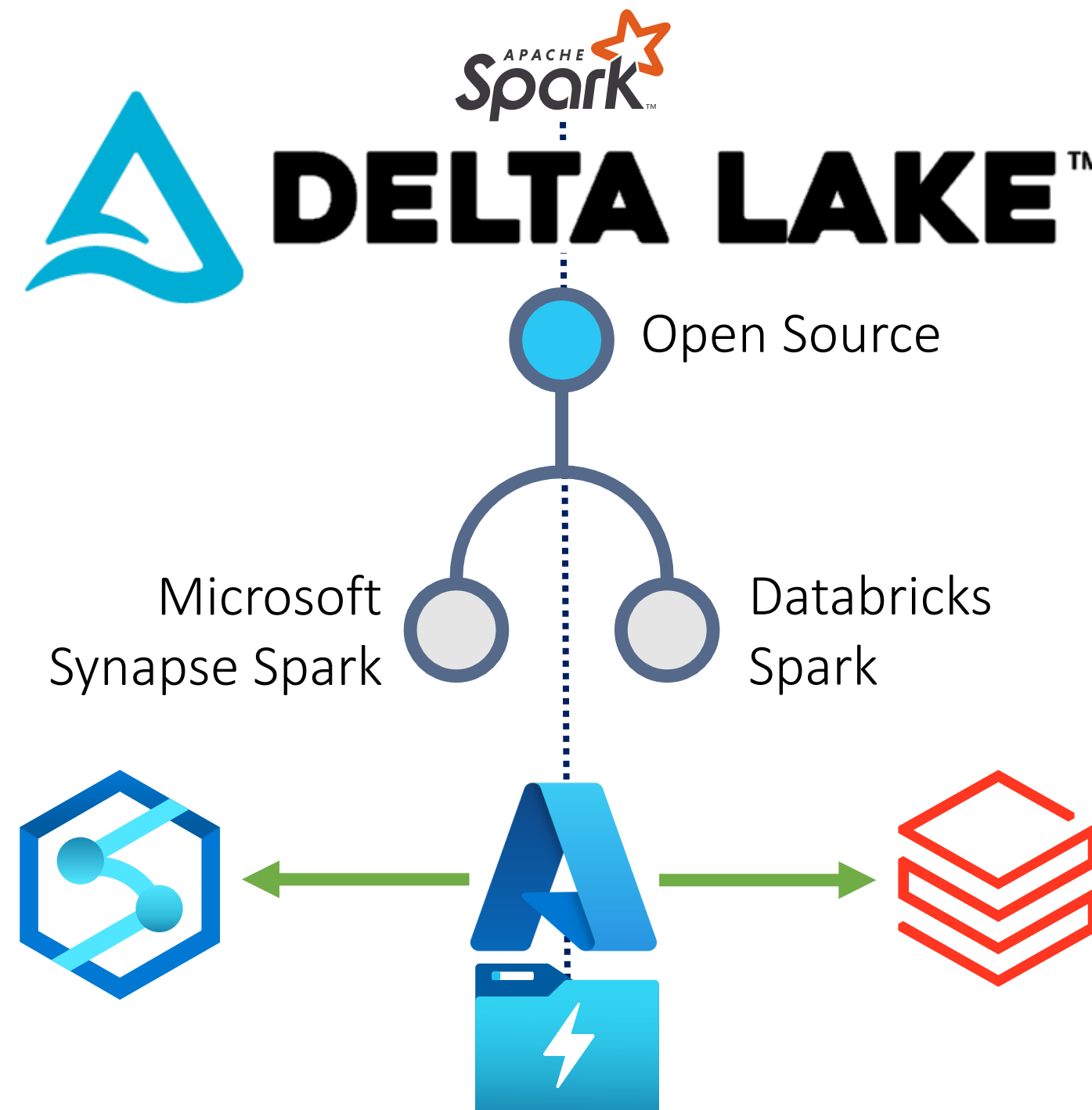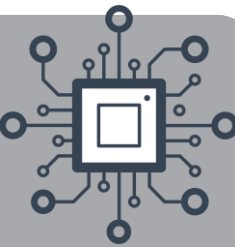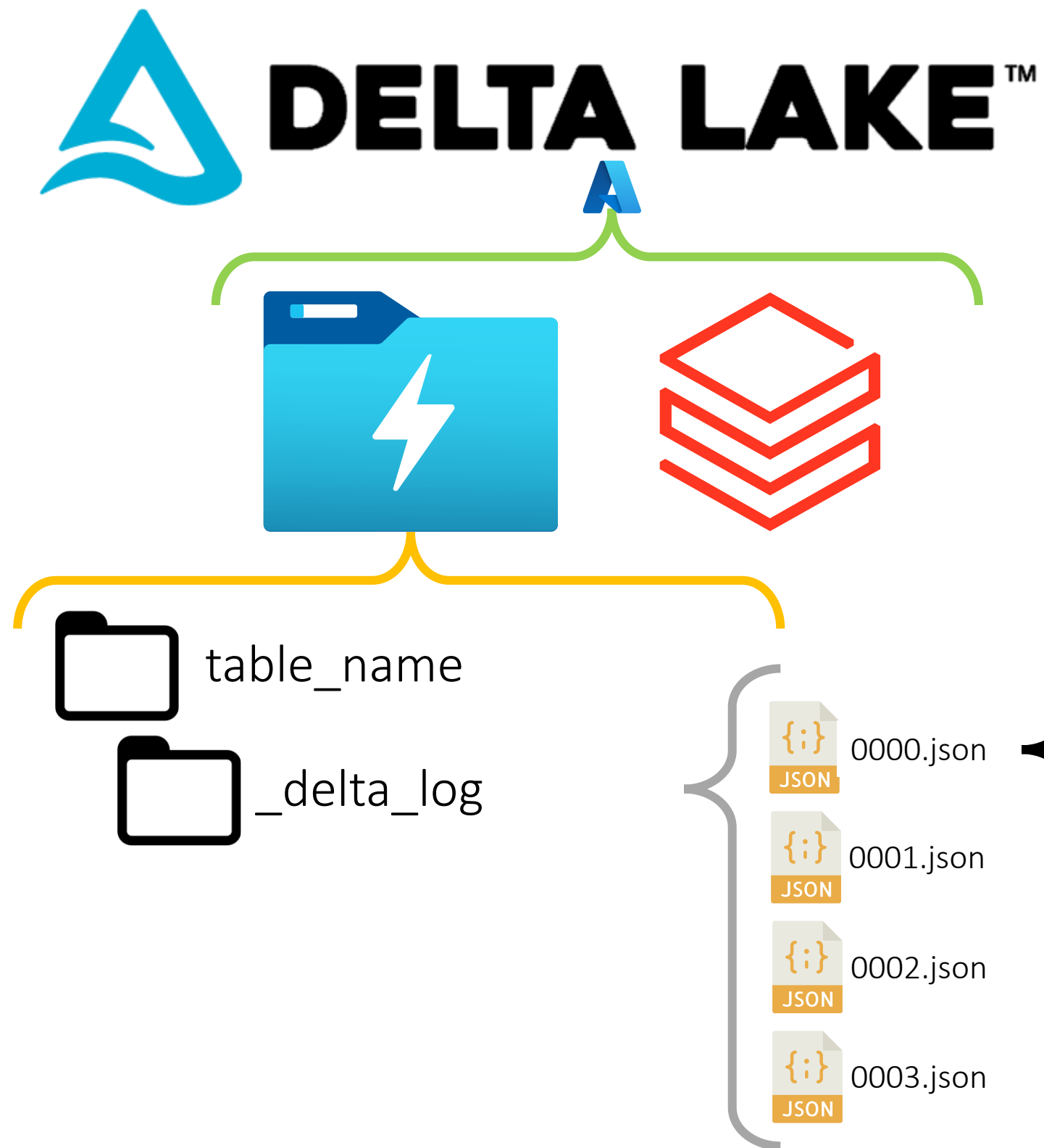*"Delta Lake is an open-source storage layer that brings ACID transactions to Apache Spark™ and big data workloads."*

table_name

_delta_log

{;} 0000.json
JSON

{;} 0001.json
JSON

{;} 0002.json
JSON

{;} 0003.json
JSON
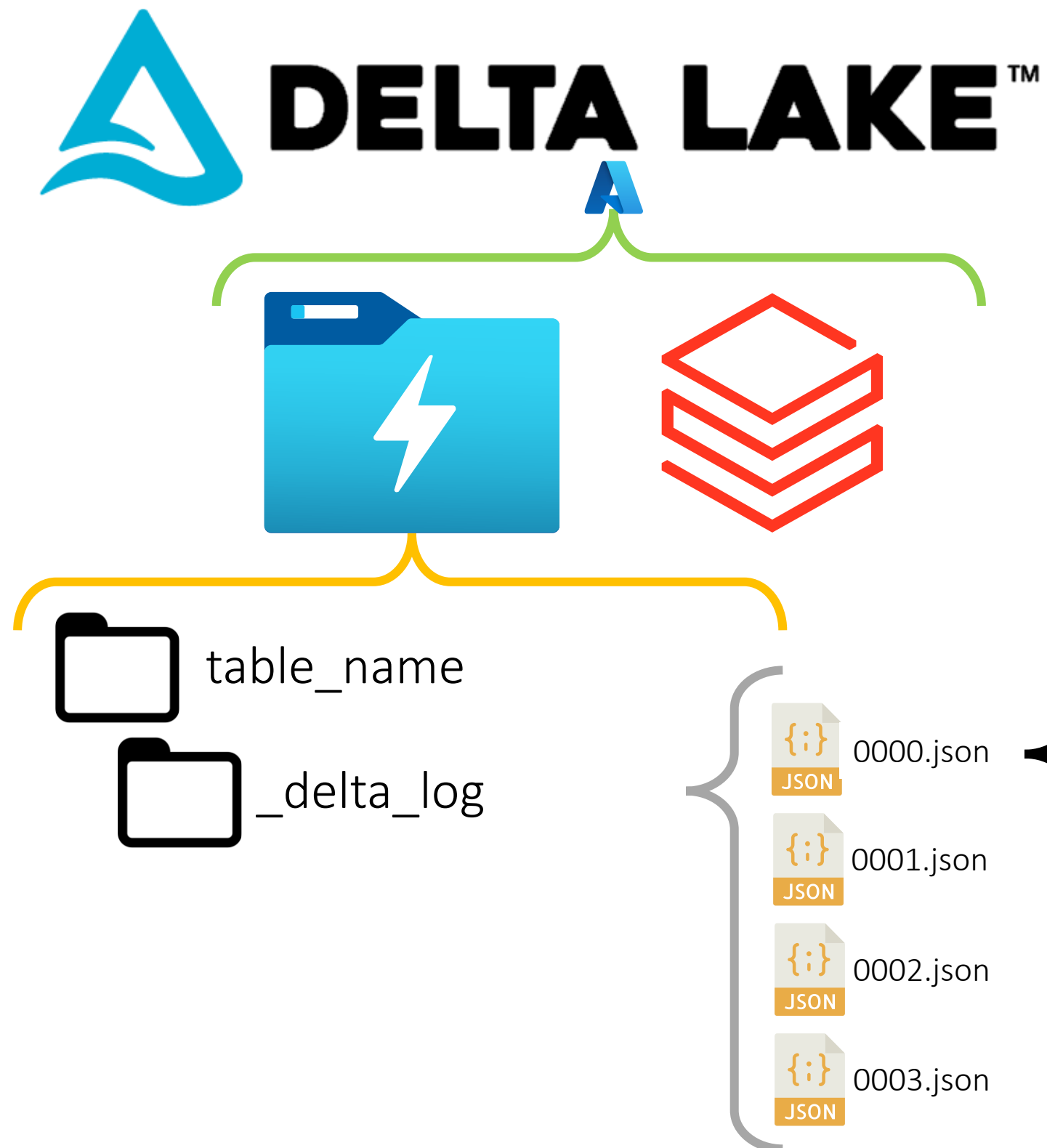
```
{
    "commitInfo":
    {
        "timestamp":1628596034417,
        "operation":"WRITE",
        "operationParameters":
        {
            "mode":"ErrorIfExists",
            "partitionBy":"[]"
        },
        "isBlindAppend":true,
        "operationMetrics":
        {
            "numFiles":"6",
            "numOutputBytes":"2407",
            "numOutputRows":"5"
        }
    }
}
{"protocol":{"minReaderVersion":1,"minWriterVersion":2}}
{
    "metaData":
    {
        "id":"58e5de01-de72-4d5b-a208-d0b4ae919efe",
        "format":
        {
            "provider":"parquet",
            "options":{}
        },
        "schemaString":
            "{\"type\":\"struct\",\"fields\":[{\"name\":\"id\",\"type\":\"long\",
            \"nullable\":true,\"metadata\":{}}]}",
        "partitionColumns":[],
        "configuration":{},
        "createdTime":1628596029470
    }
}
{"add":{"path":"part-00000.snappy.parquet","size":262,"modificationTime":1628596034000}}
{"add":{"path":"part-00001.snappy.parquet","size":429,"modificationTime":1628596034000}}
{"txn":{"appId":"731b2c96-bf64-445c-8ca8-cd6cad6735e2","lastUpdated":1628596094191}}
{"add":{"path":"part-00000.snappy.parquet","size":429,"modificationTime":1628596094000}}
{"add":{"path":"part-00001.snappy.parquet","size":429,"modificationTime":1628596094000}}
{"remove":{"path":"part-00150.snappy.parquet","deletionTimestamp":1628596098597}}
{"remove":{"path":"part-00128.snappy.parquet","deletionTimestamp":1628596098597}}
```

table_name

_delta_log

0000.json

0001.json

0002.json

0003.json

## Breaking Down Transactions Into Atomic Commits

Whenever a user performs an operation to modify a table (such as an INSERT, UPDATE or DELETE), Delta Lake breaks that operation down into a series of discrete steps composed of one or more of the actions below.

**Add file** – adds a data file.

**Remove file** – removes a data file.

**Update metadata** – Updates the table's metadata (e.g., changing the table's name, schema or partitioning).
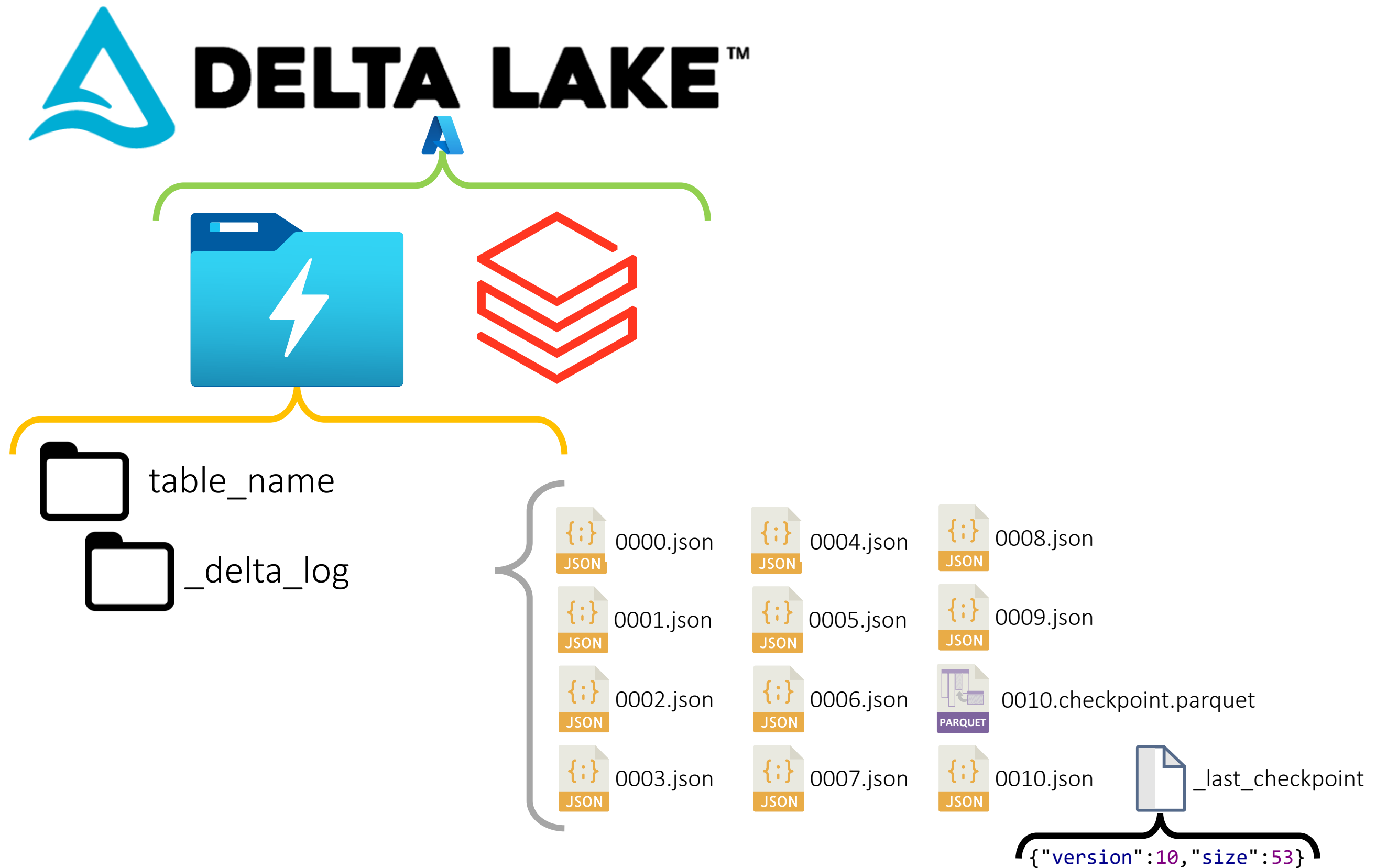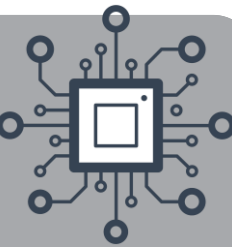
**Set transaction** – Records that a structured streaming job has committed a micro-batch with the given ID.

**Change protocol** – enables new features by switching the Delta Lake transaction log to the newest software protocol.
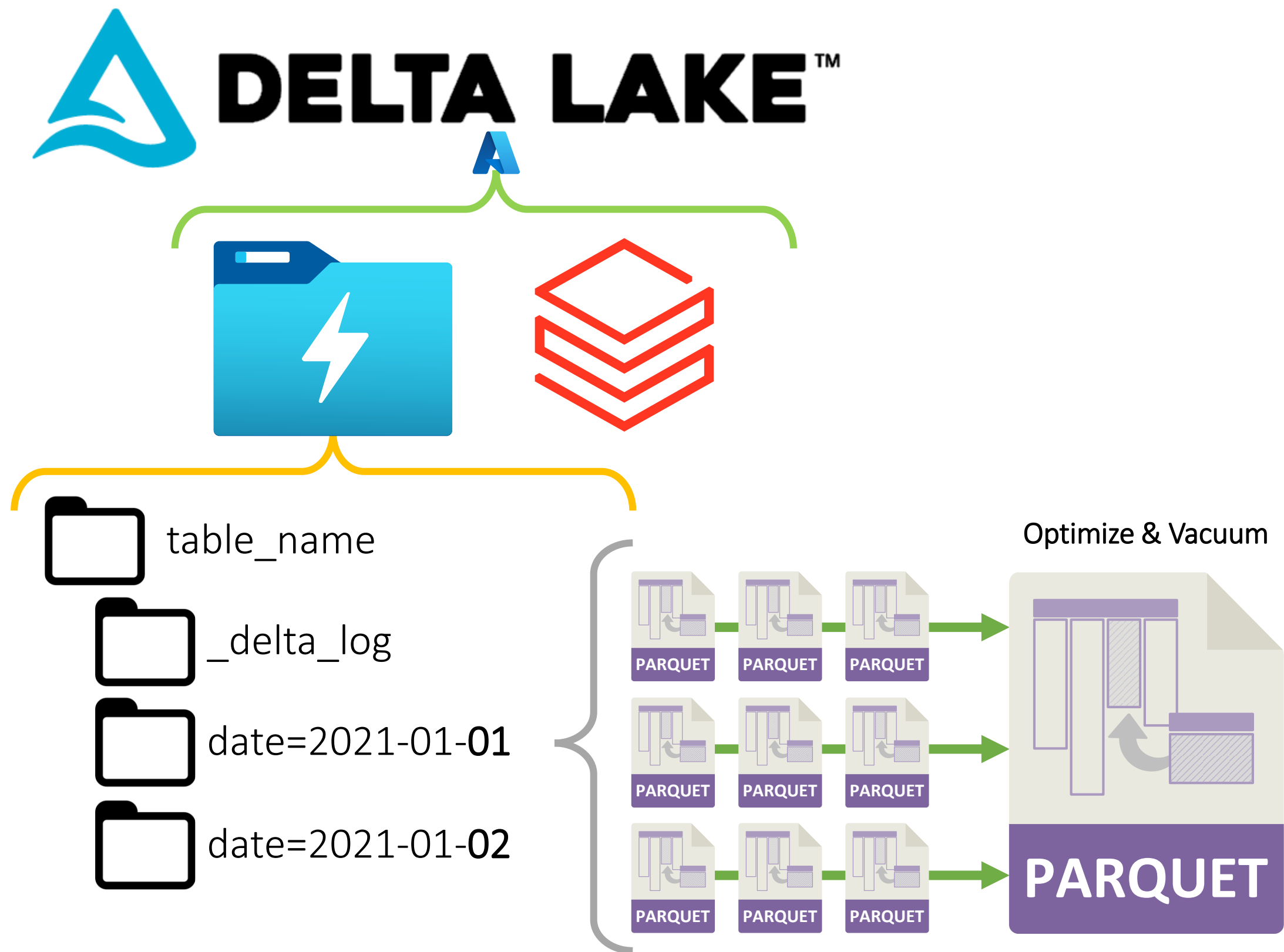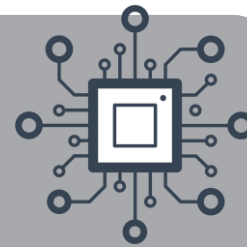
**Commit info** – Contains information around the commit, which operation was made, from where and at what time.
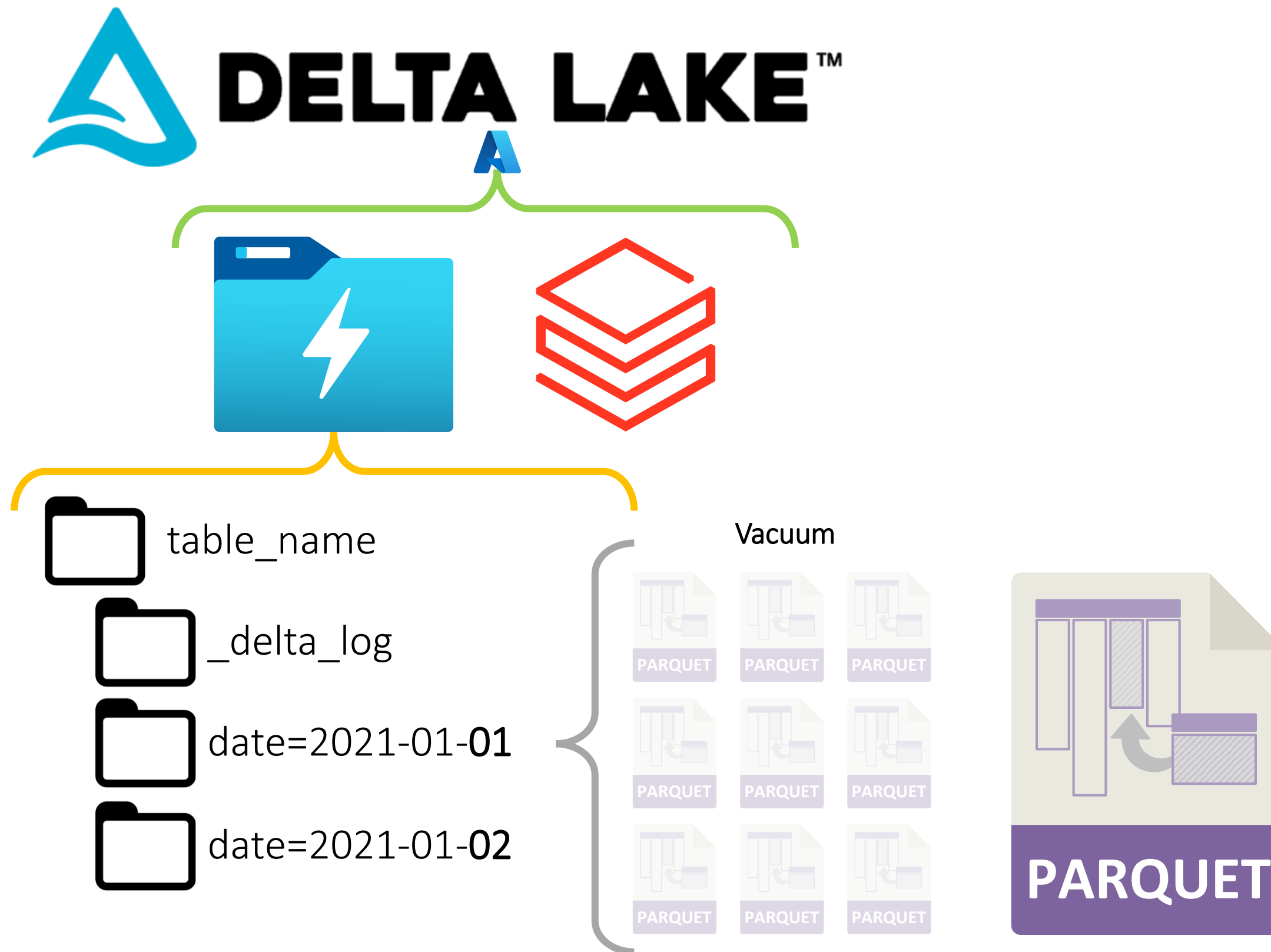
*Source: https://databricks.com/blog*

DELTA LAKE™

table_name

_delta_log

0000.json  0004.json  0008.json

0001.json  0005.json  0009.json

0002.json  0006.json  0010.checkpoint.parquet

0003.json  0007.json  0010.json  _last_checkpoint

{"version":10,"size":53}

table_name

_delta_log

date=2021-01-**01**

date=2021-01-**02**

Optimize & Vacuum

PARQUET

DELTA LAKE™

- table_name
  - _delta_log
  - date=2021-01-**01**
  - date=2021-01-**02**

Vacuum

PARQUET

# Data Warehouse

**O**nline
**L**ine
**T**ransactional
**P**rocessing

**O**ffline
**A**nalytical
**T**ransactional
**P**rocessing

Application
Data

**E**xtract
**T**ransform
**L**oad

Data Collection

Data Sources

Data Warehouse

Data Insight

# Lake House



**O**nline
**L**ine
**T**ransactional
**P**rocessing

**O**ffline
**A**nalytical
**T**ransactional
**P**rocessing

DELTA LAKE

Extract
Transform
Load

Application
Data

Lake
House

Spark

# Lake House

WIKIPEDIA
The Free Encyclopedia

Article   Talk

Read   Edit   View history

Search Wikipedia

## The Lake House (film)

From Wikipedia, the free encyclopedia

**The Lake House** is a 2006 American fantasy romantic drama film directed by Alejandro Agresti, starring Keanu Reeves and Sandra Bullock (who had previously appeared together in the box office hit *Speed*). It was written by David Auburn.[2] A remake of the South Korean motion picture *Il Mare* (2000), it centers on an architect living in 2004 and a doctor living in 2006 who meet via letters left in a mailbox at the lake house where they have lived at separate points in time. They carry on correspondence over two years, remaining separated by their original difference of two years.[3]

### The Lake House

Theatrical release poster

| | |
|---|---|
| Directed by | Alejandro Agresti |
| Written by | David Auburn |
| Based on | *Il Mare* by Kim Eun-jeong Kim Mi-yeong |
| Produced by | Doug Davison Roy Lee |
| Starring | Keanu Reeves |

## Plot   [ edit ]

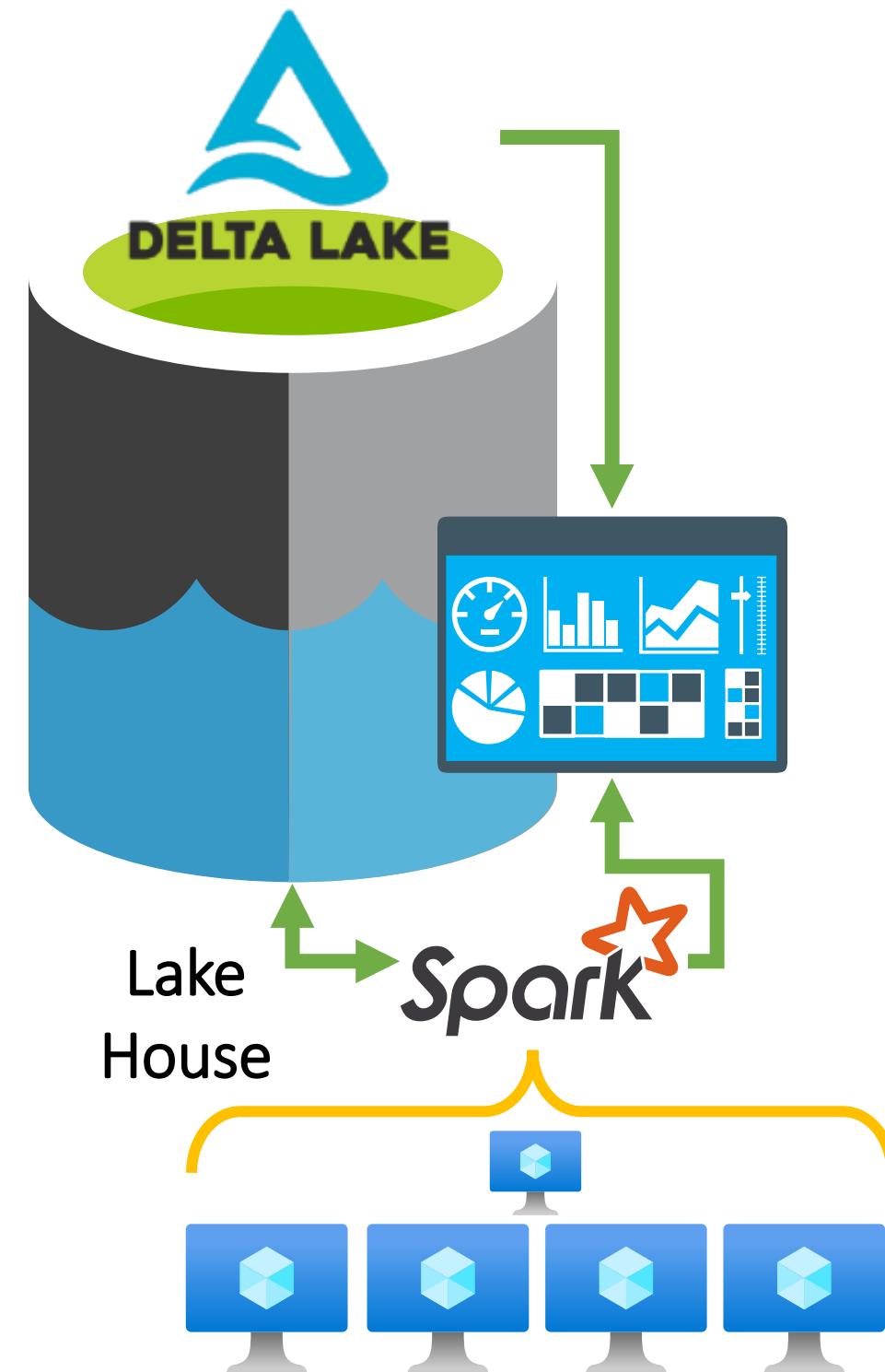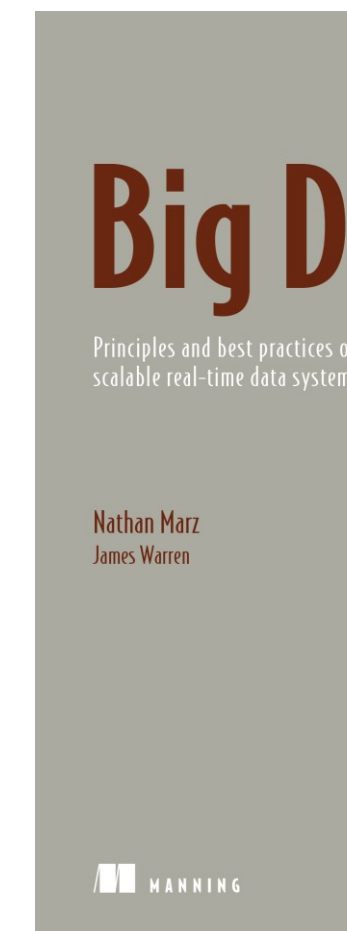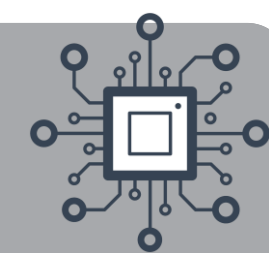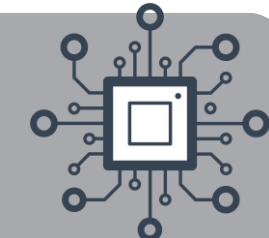In 2006, Dr. Kate Forster (Sandra Bullock) is leaving a lake house that she has been renting in Chicago. Kate leaves a note in the mailbox for the next tenant to forward her mail, adding that the paint-embedded pawprints on the path leading to the house were already there when she arrived.

### Navigation sidebar

Main page
Contents
Current events
Random article
About Wikipedia
Contact us
Donate

**Contribute**

Help
Learn to edit
Community portal
Recent changes
Upload file

**Tools**

What links here
Related changes
Special pages
Permanent link
Page information
Cite this page
Wikidata item

**Print/export**

Download as PDF
Printable version

**Languages**

العربية
Deutsch
Español
Français

DELTA LAKE

Lake House

Spark

# Lambda* & Kappa

λ κ

Batch →

Speed →

Serve

← Query

**Big Data**

Principles and best practices of
scalable real-time data systems

Nathan Marz
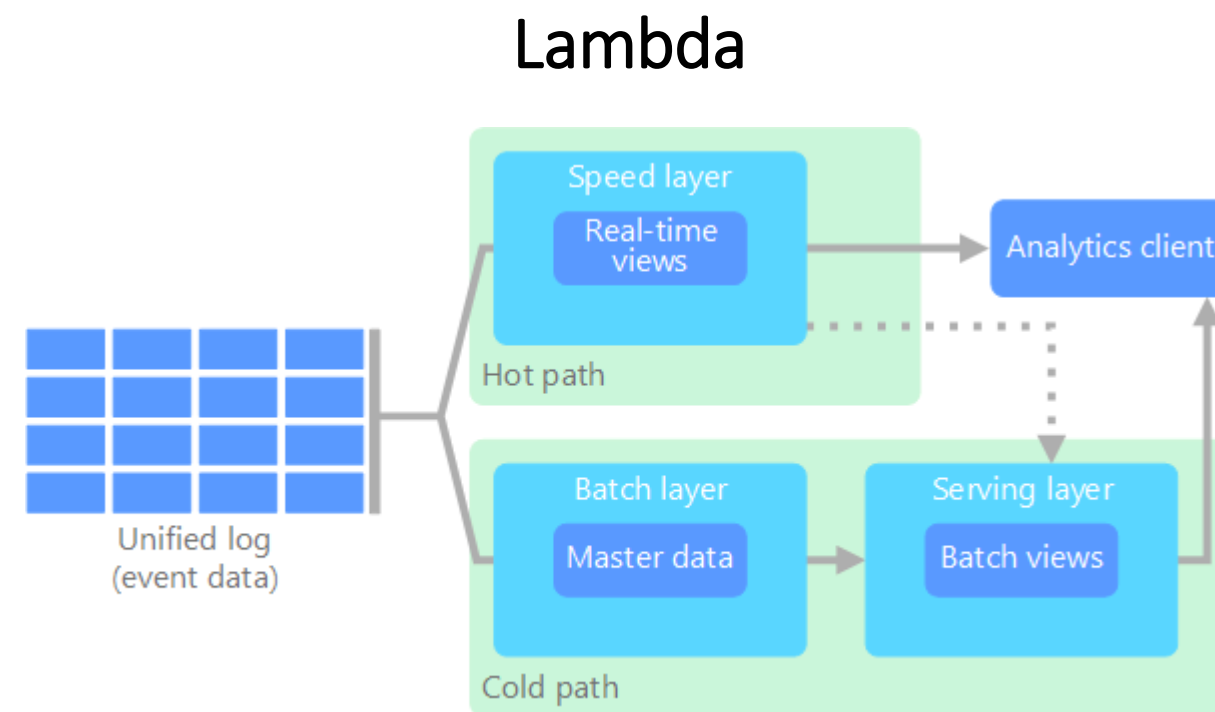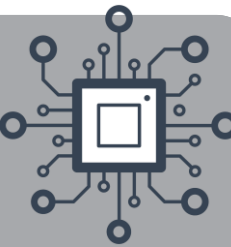James Warren

MANNING

~2018

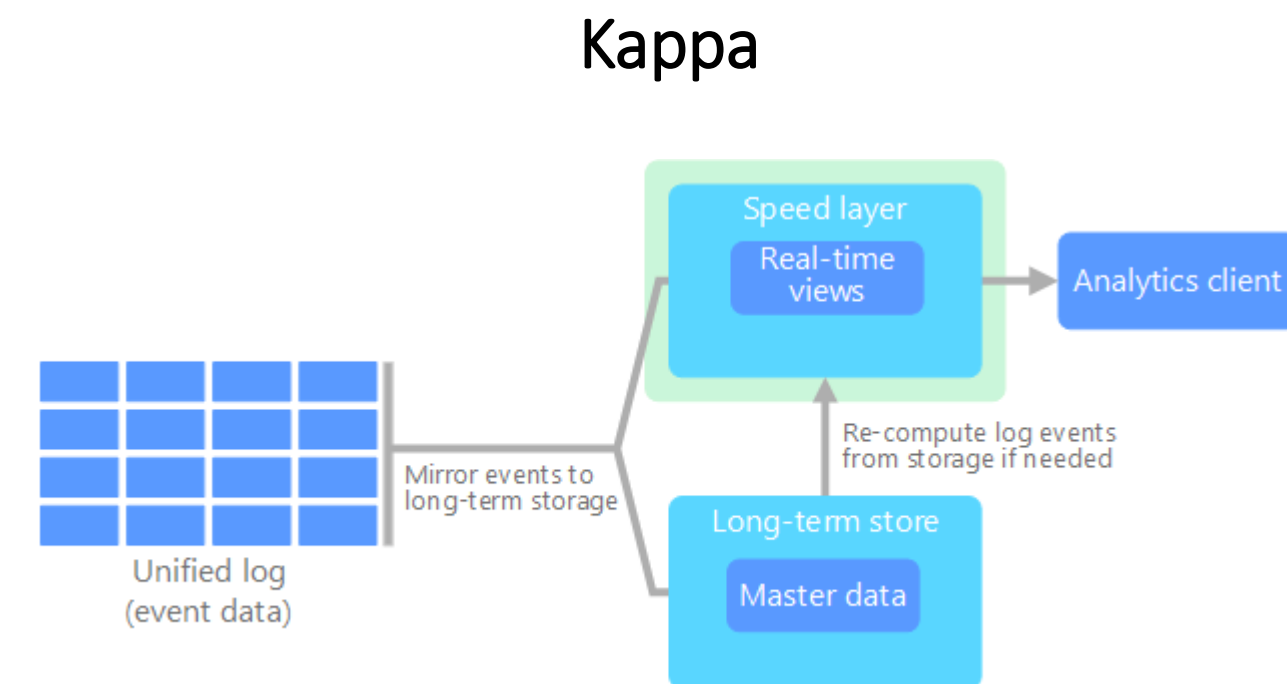Batch →

Speed →

Serve

Query →

# Lambda & Kappa Architectures

## Lambda



## Kappa



"The **lambda architecture**, first proposed by Nathan Marz, addresses this problem by creating two paths for data flow. All data coming into the system goes through these two paths:

A **batch layer** (cold path) stores all of the incoming data in its raw form and performs batch processing on the data. The result of this processing is stored as a **batch view**.

A **speed layer** (hot path) analyzes data in real time. This layer is designed for low latency, at the expense of accuracy."
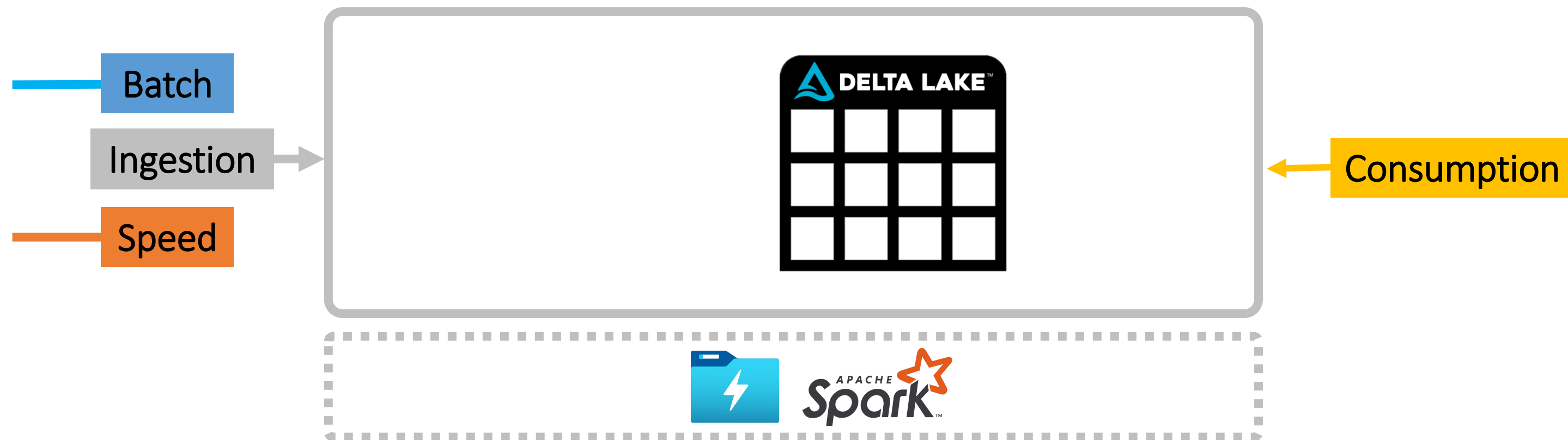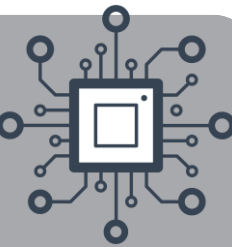
"A drawback to the lambda architecture is its ==complexity==. ==Processing logic appears in two different places== — the cold and hot paths — using different frameworks. This leads to duplicate computation logic and the complexity of managing the architecture for both paths.
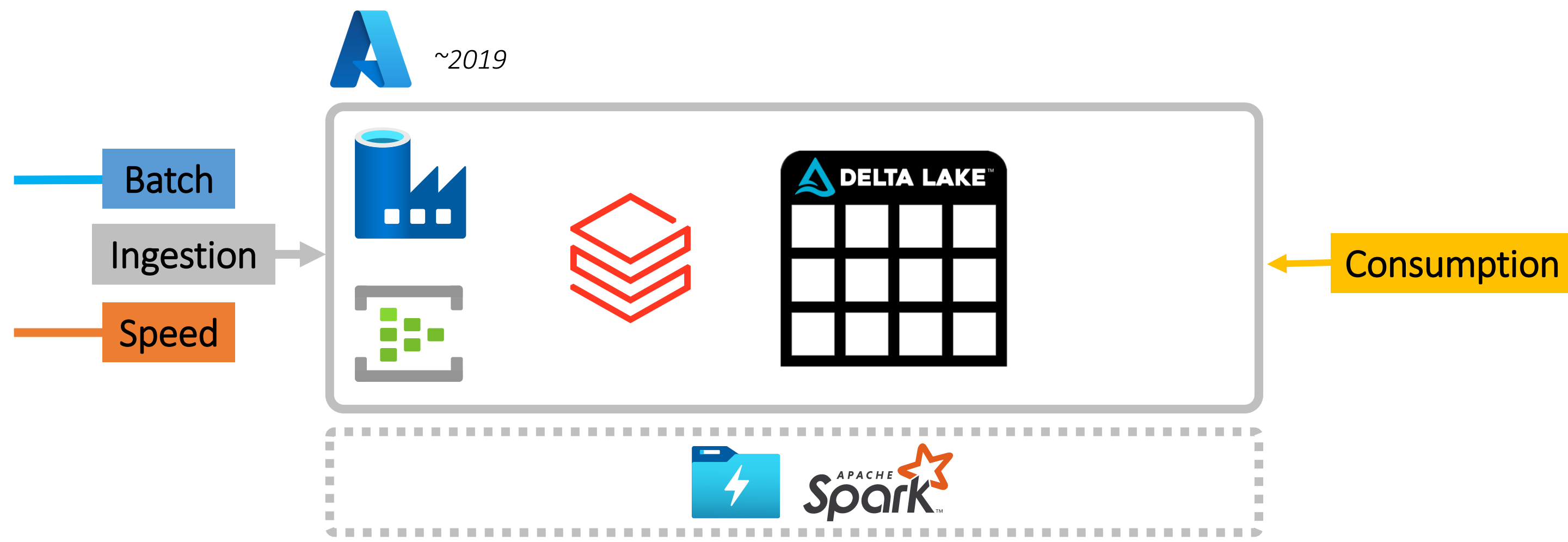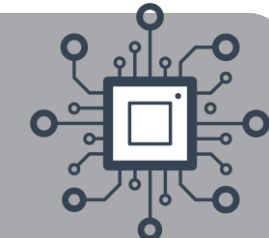
The **kappa architecture** was proposed by Jay Kreps as an alternative to the lambda architecture. It has the same basic goals as the lambda architecture, but with an important distinction: All data flows through a single path, using a stream processing system."
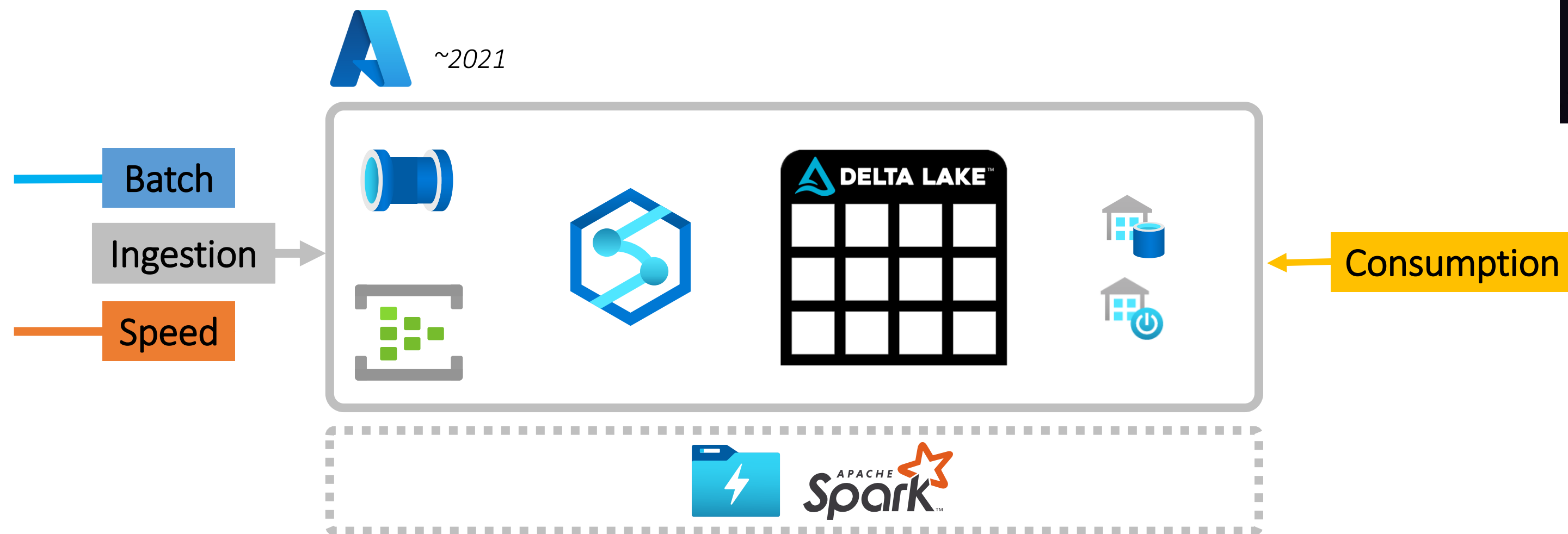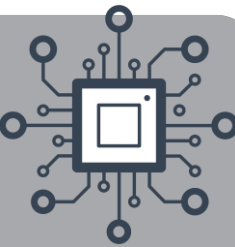
https://docs.microsoft.com/en-us/azure/architecture/data-guide/big-data/

# Lambda & Kappa Architectures

~2021

Batch

Ingestion

Speed

Consumption

DELTA LAKE

APACHE Spark
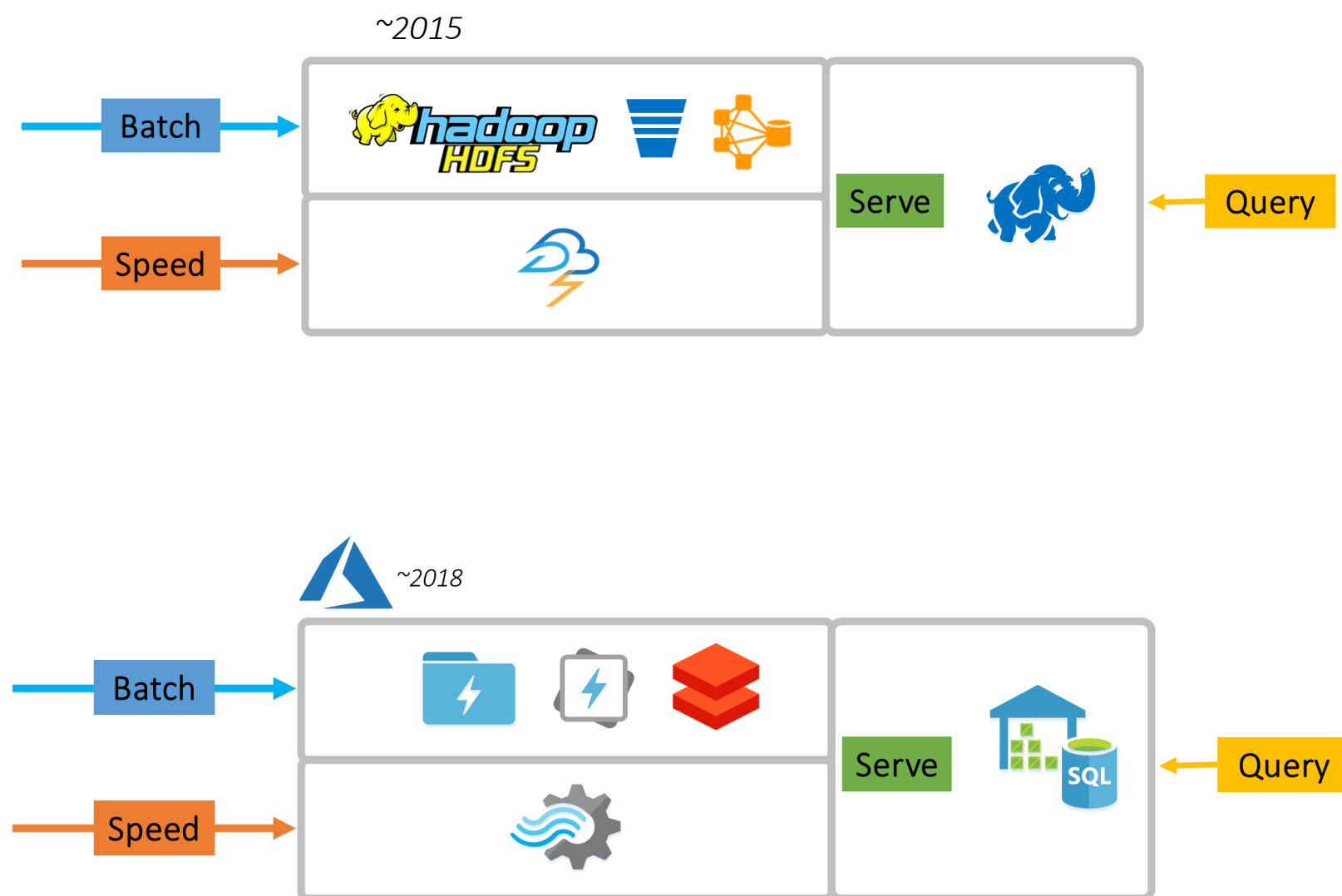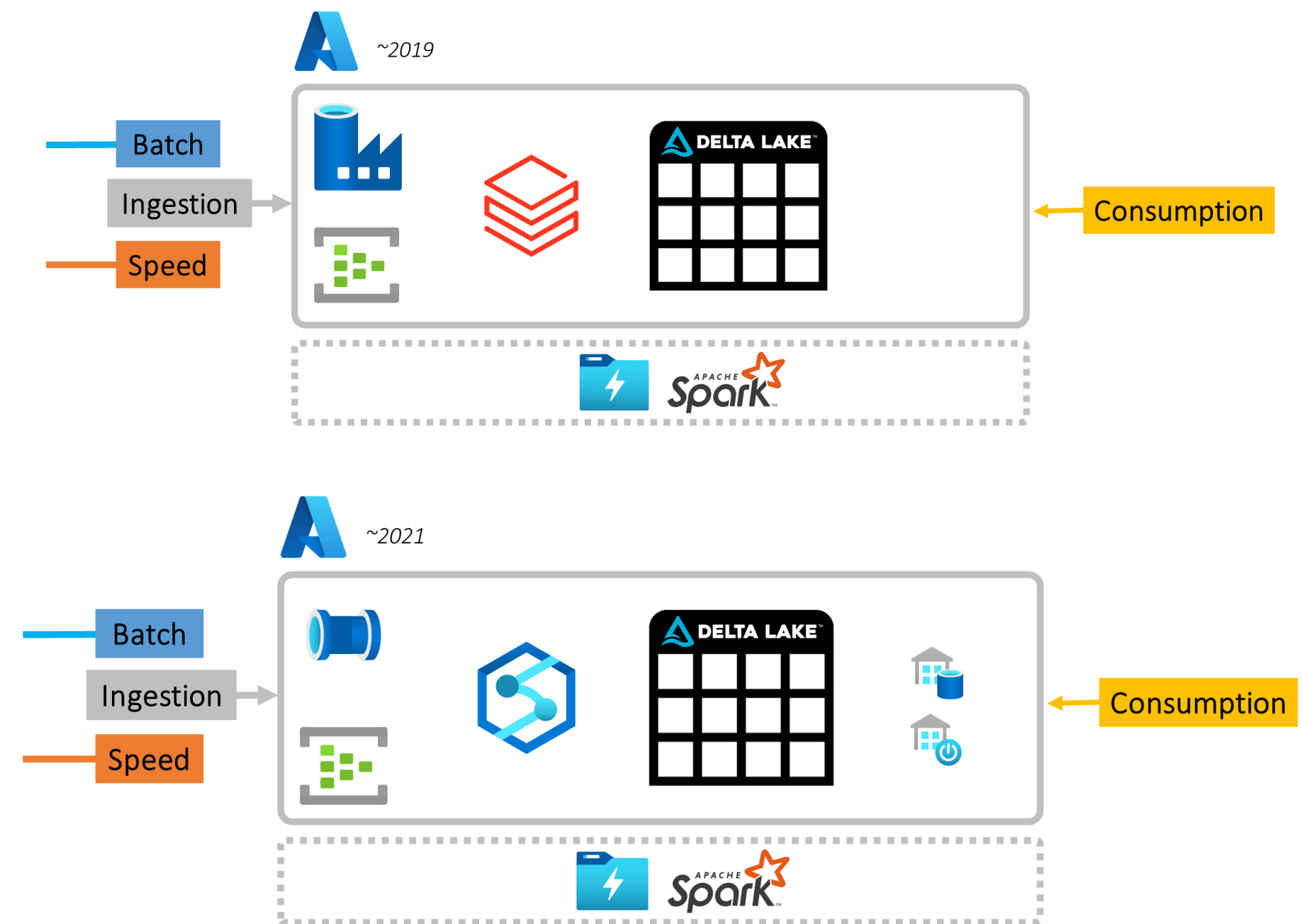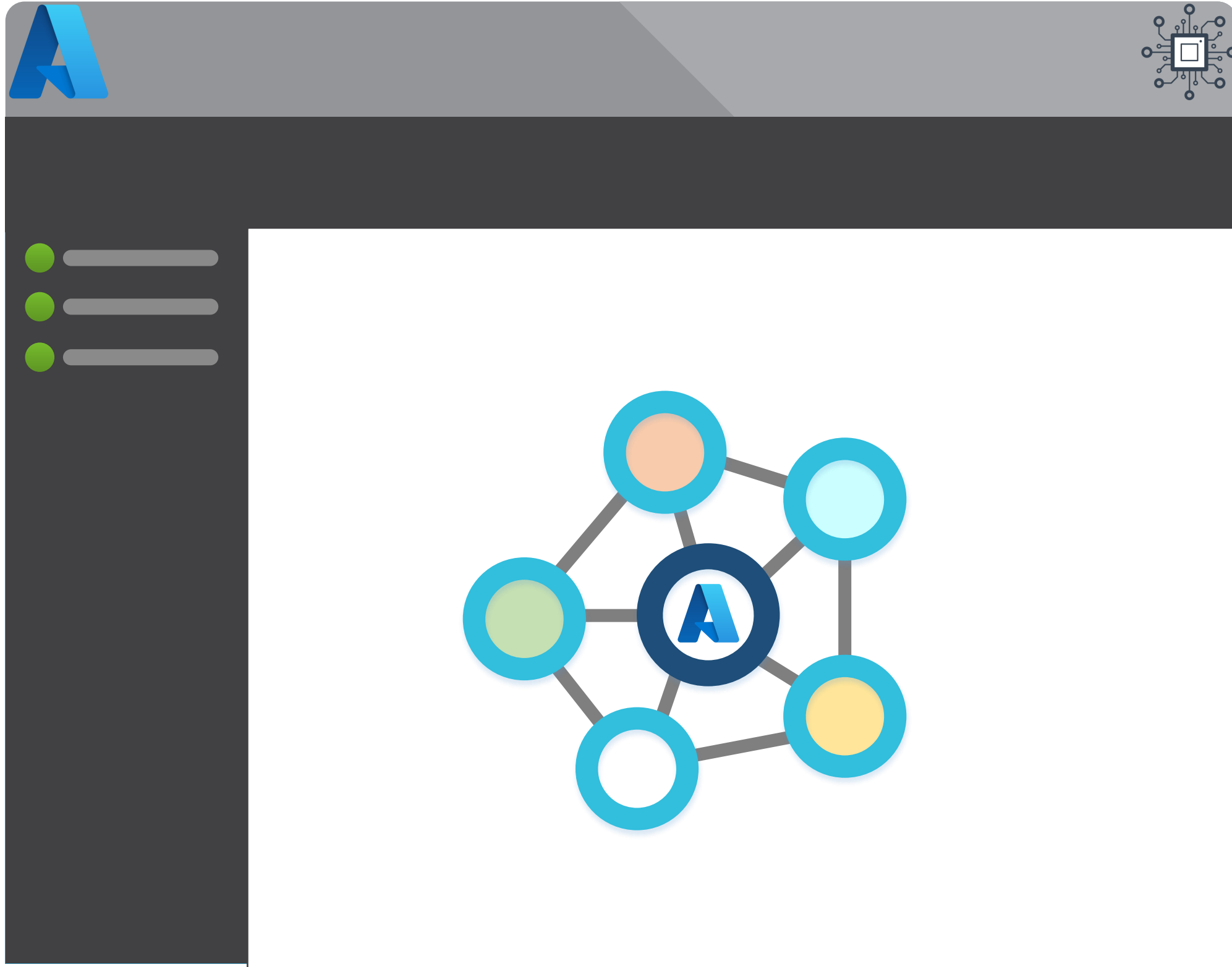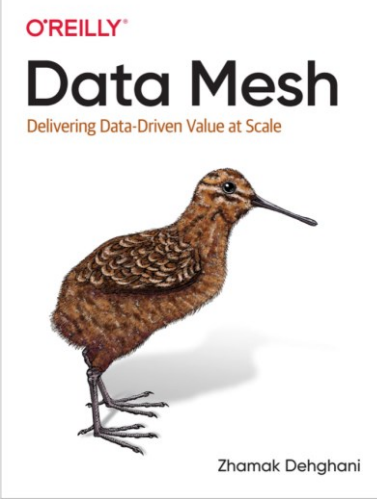
Delta Lake in the Context of Lambda & Kappa
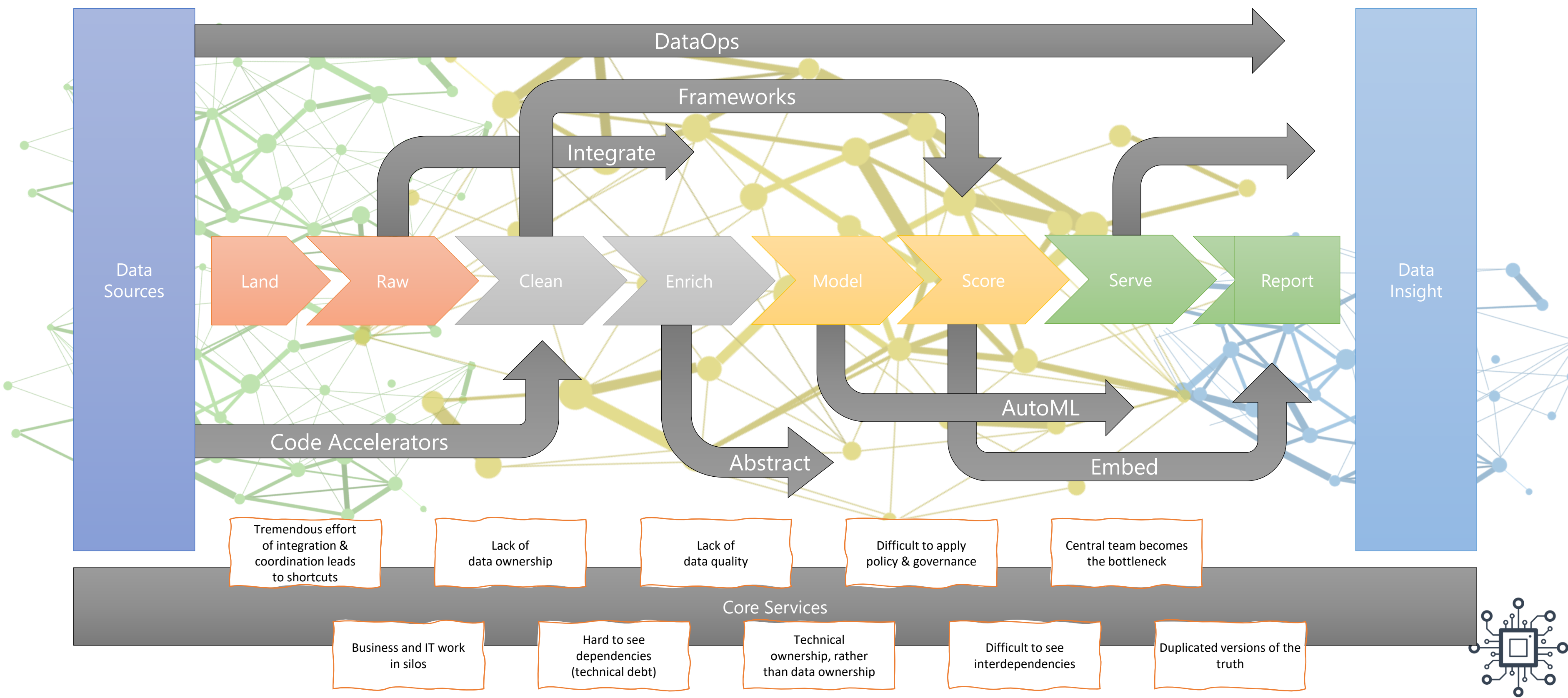
# Data Mesh

*- Zhamak Dehghani*

*@zhamakd*

1. Domain-oriented decentralised data ownership and architecture.

2. Data as a product.

3. Self-serve data infrastructure as a platform.

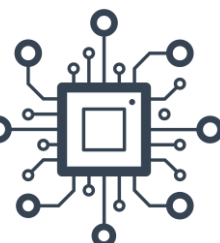4. Federated computational governance.

# Data Mesh – *Why should we build it?*

Using a **traditional centralised approach**, enhanced with cloud scale technologies to create a modern data analytics platform.



Data Sources

DataOps

Frameworks

Integrate

Land — Raw — Clean — Enrich — Model — Score — Serve — Report

Code Accelerators

Abstract

AutoML

Embed

Data Insight

Tremendous effort of integration & coordination leads to shortcuts

Lack of data ownership

Lack of data quality

Difficult to apply policy & governance

Central team becomes the bottleneck

Core Services

Business and IT work in silos

Hard to see dependencies (technical debt)

Technical ownership, rather than data ownership

Difficult to see interdependencies
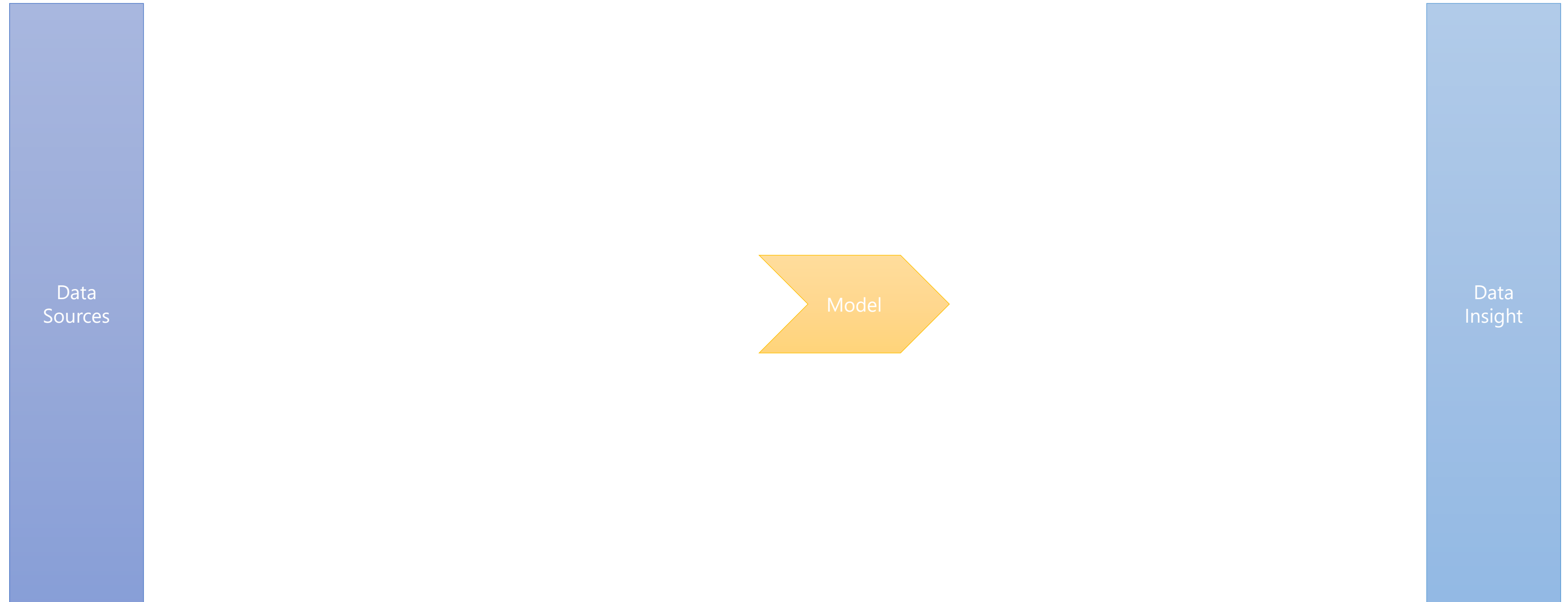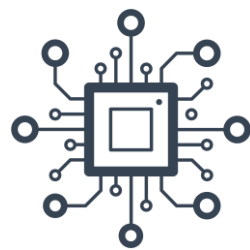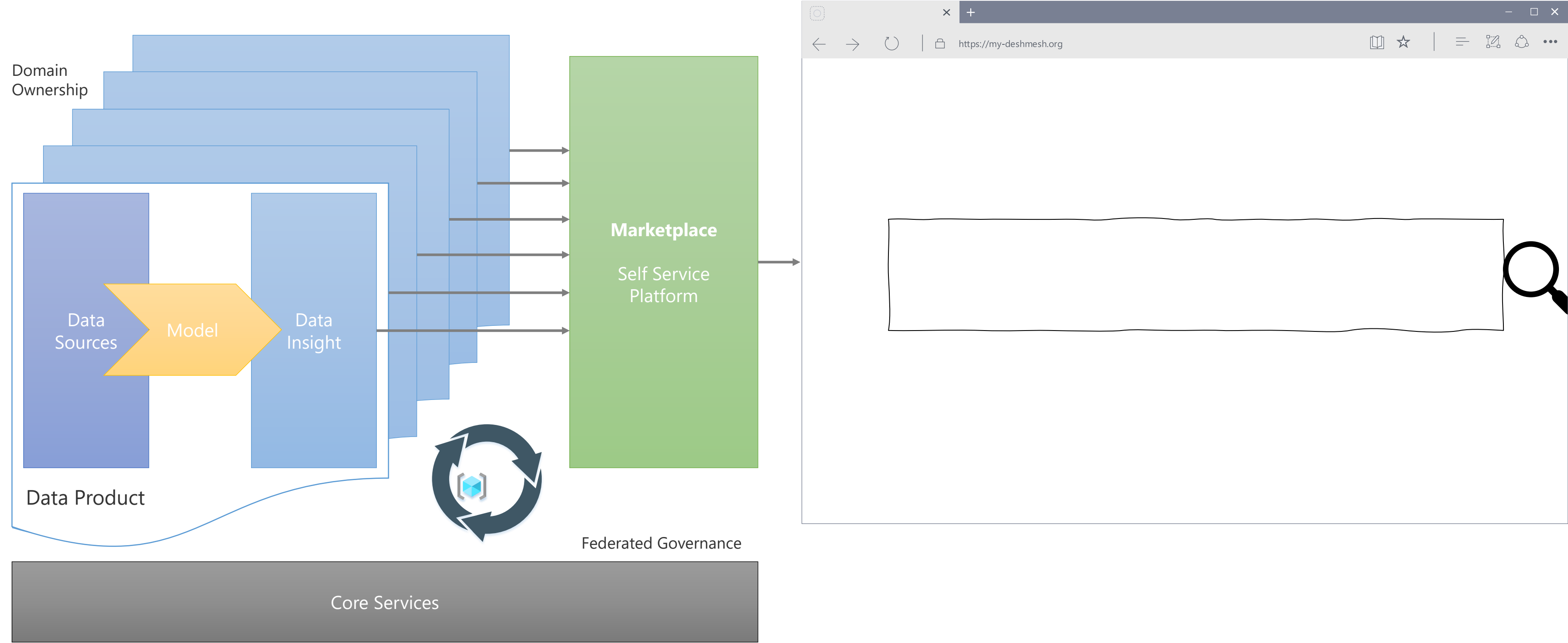
Duplicated versions of the truth

# Data Mesh – *Why should we build it?*

Using a **traditional centralised** approach, enhanced with cloud scale
technologies to create a modern data analytics platform.
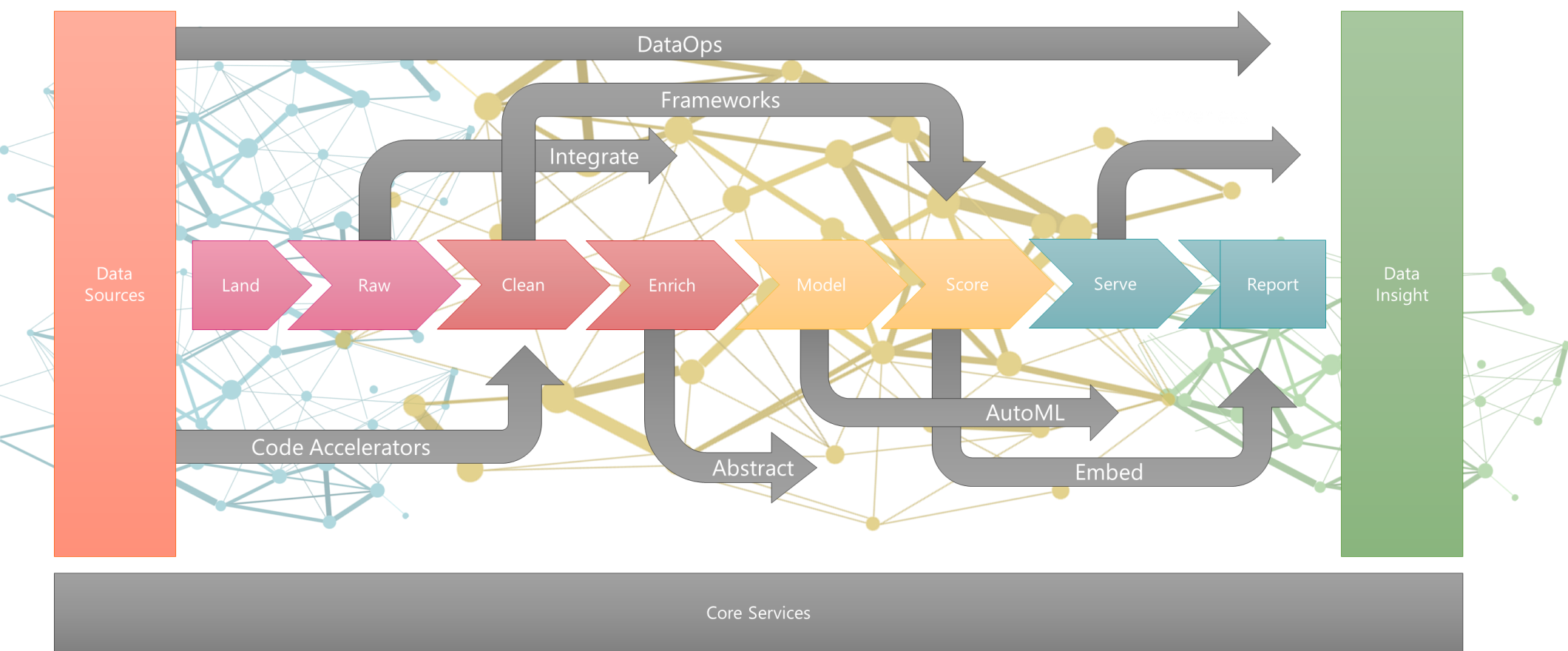
Data
Sources

Model

Data
Insight

# Data Mesh – *Why should we build it?*

Using a **de-centralised** approach to cloud scale analytics, empowering users to rapidly gain insights to make strategic business decisions.

Domain Ownership

Data Sources

Model

Data Insight

Data Product

**Marketplace**

Self Service Platform

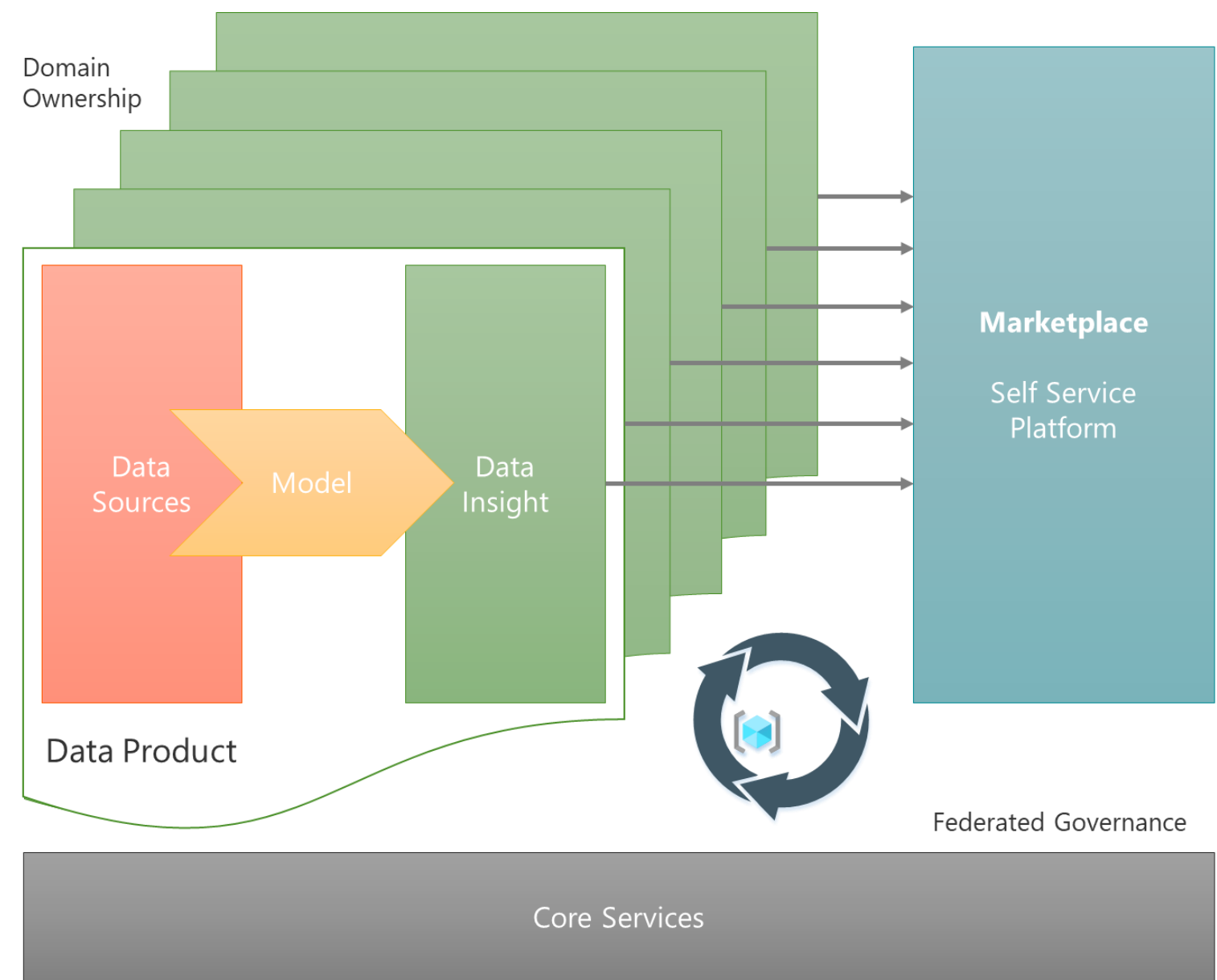Federated Governance

Core Services

https://my-deshmesh.org

# Time to Insight

Using a **traditional centralised approach**, enhanced with cloud scale technologies to create a modern data analytics platform.
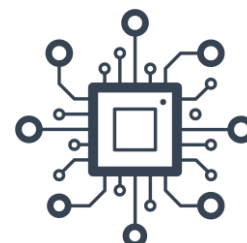
Using a **de-centralised** approach to cloud scale analytics, empowering users to rapidly gain insights to make strategic business decisions.



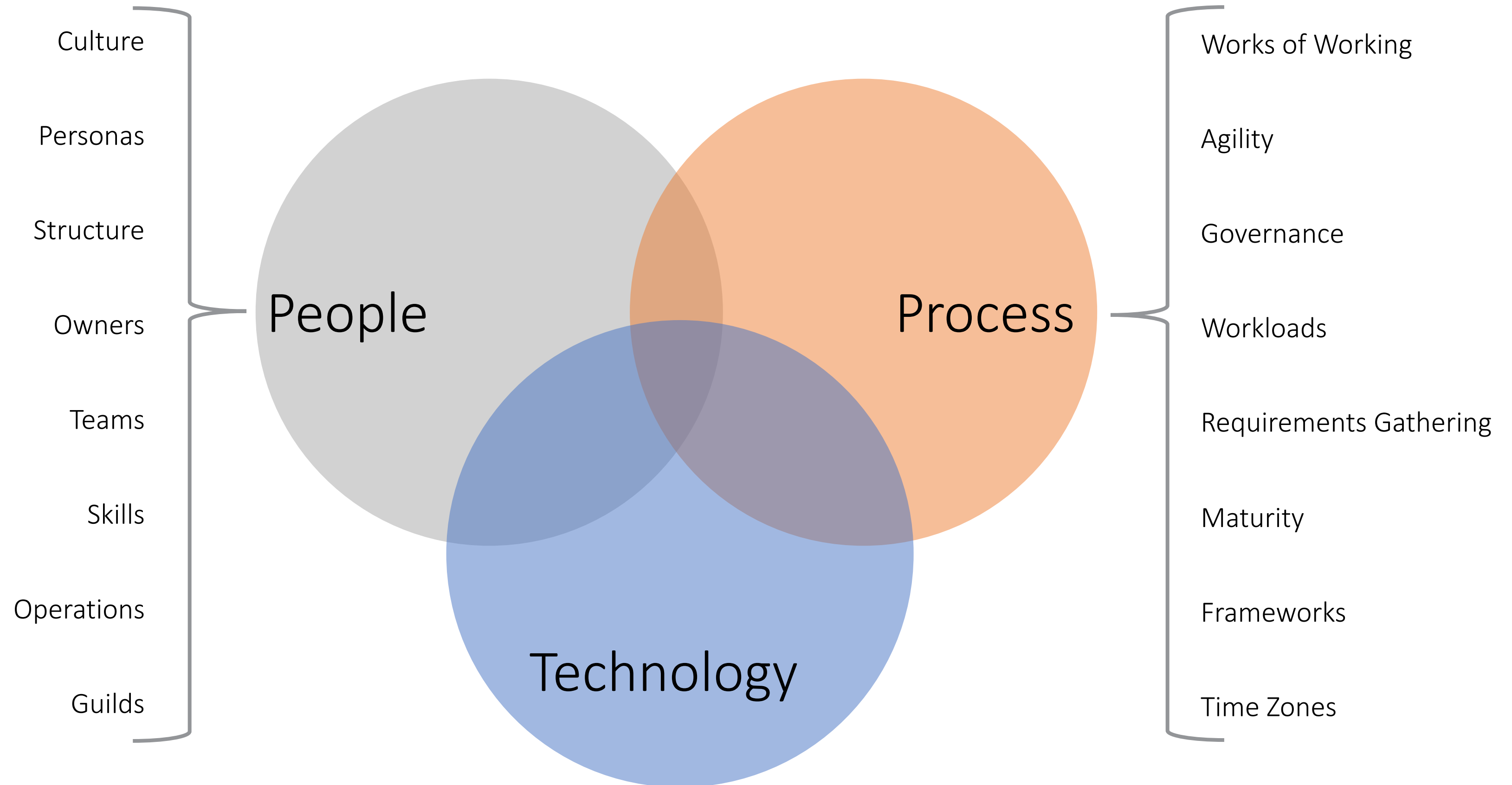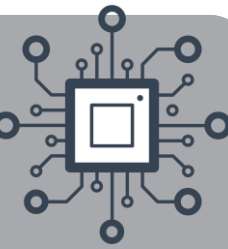**... Weeks/Months**

**... Hours/Days**

# Introducing the Data Mesh
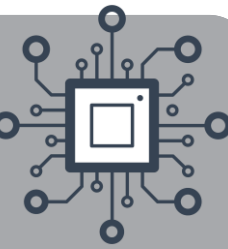


Culture

Personas

Structure

Owners

Teams

Skills

Operations

Guilds

People

Process

Technology

Works of Working

Agility

Governance

Workloads

Requirements Gathering

Maturity

Frameworks

Time Zones

Technology

# Data Products

# Data Products in Azure

**Diagram 1 (top):**

Data Collection / Data Sources → [ Compute | Compute / Storage / Orchestration ] → Data Insight

**Diagram 2 (left):**

Data Collection / Data Sources → [ Compute | Compute / Storage / Orchestration ] → Data Insight

**Diagram 3 (right):**

Data Collection / Data Sources → [ Compute | Compute / Storage / Orchestration ] → Data Insight

**Diagram 4 (bottom):**

Data Collection / Data Sources → [ Compute | Compute / Storage / Orchestration ] → Data Insight

# Data Products in Azure with Interfaces

**Tertiary Interfaces -** resource connectivity.

Service & Private Endpoints

Front Door

Peered VNet

Data Product **Core**

Express Route Circuits

**Secondary Interfaces -** operational reporting and logging.

Workbooks

Alerting

Log Analytics

Event Hub

Data Product **Core**

Storage

Activity Log

Security Center

Azure Monitor

Subscriber Events

Service Queue

Batch Integration

Function Wrapper

Data Product **Core**

REST/OData Endpoints

Management API

Message Stream

SQL Endpoint

**Primary Interfaces –** data integration and exchange.

Data Domain 1

Data Domains in Azure

Data Domain 1

Federated computational governance

Computational policies

Data product as architecture quantum

Multi-plane data platform

Core Services

Policies & Platform Wrappers

Data Product

Data Product

Data Product

PaaS Plane

IaaS Plane

Multi Plane Components

Analytics | Virtualisation | Governance

Operational Platform

IDE's

Governance & Documentation

Build & Deployment Agents

Supplementary Assets

Templates

Versioning

Infrastructure Wrappers

https://mrpaulandrew.com/tag/data-mesh-vs-azure/

**Mesh Supervision Plane**
Capabilities that are accessible more conveniently at mesh level

Manage security policies of DPs

Discovery and explore DPs

Manage emergent graphs of DPs

Read DP

Version DP

Secure DP

Declaratively create DP

Build, deploy, monitor DP

Polyglot big data storage

Data transformation orchestration

CI/CD

Networking

Access control

**Data Product Developer Experience Plane**

The higher level abstraction of data infrastructure designed to support the common data product developer journey

**Data Infrastructure Plane**
Providing the underlying infrastructure required to build, run, monitor data products

Core Services

Policies & Platform Wrappers

Data Product

Data Product

Data Product

SaaS Plane

PaaS Plane

PaaS Plane

IaaS Plane

IaaS Plane

DELTA LAKE

**Multi Plane Components**

Analytics | Virtualisation | Governance

Operational Platform
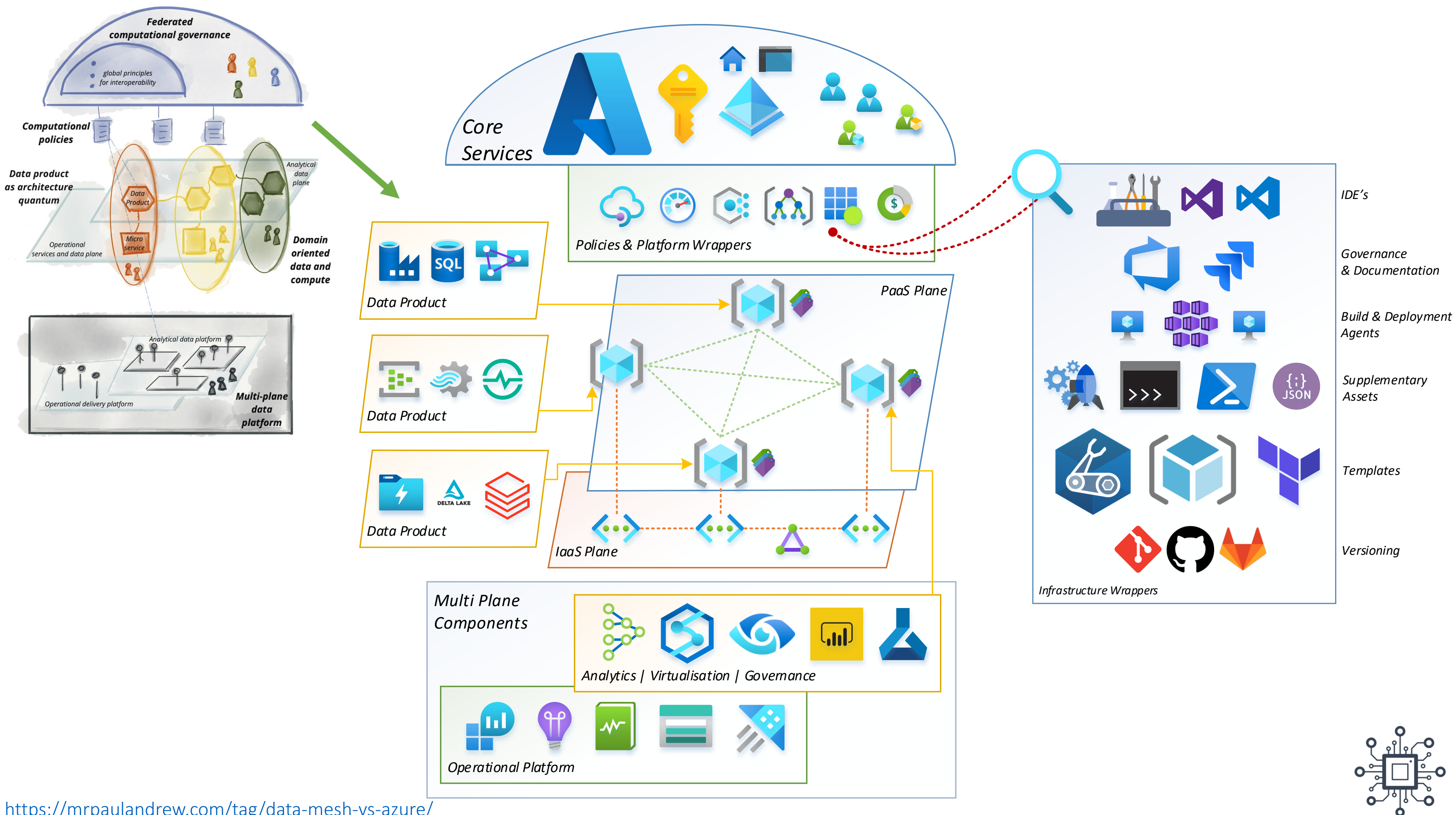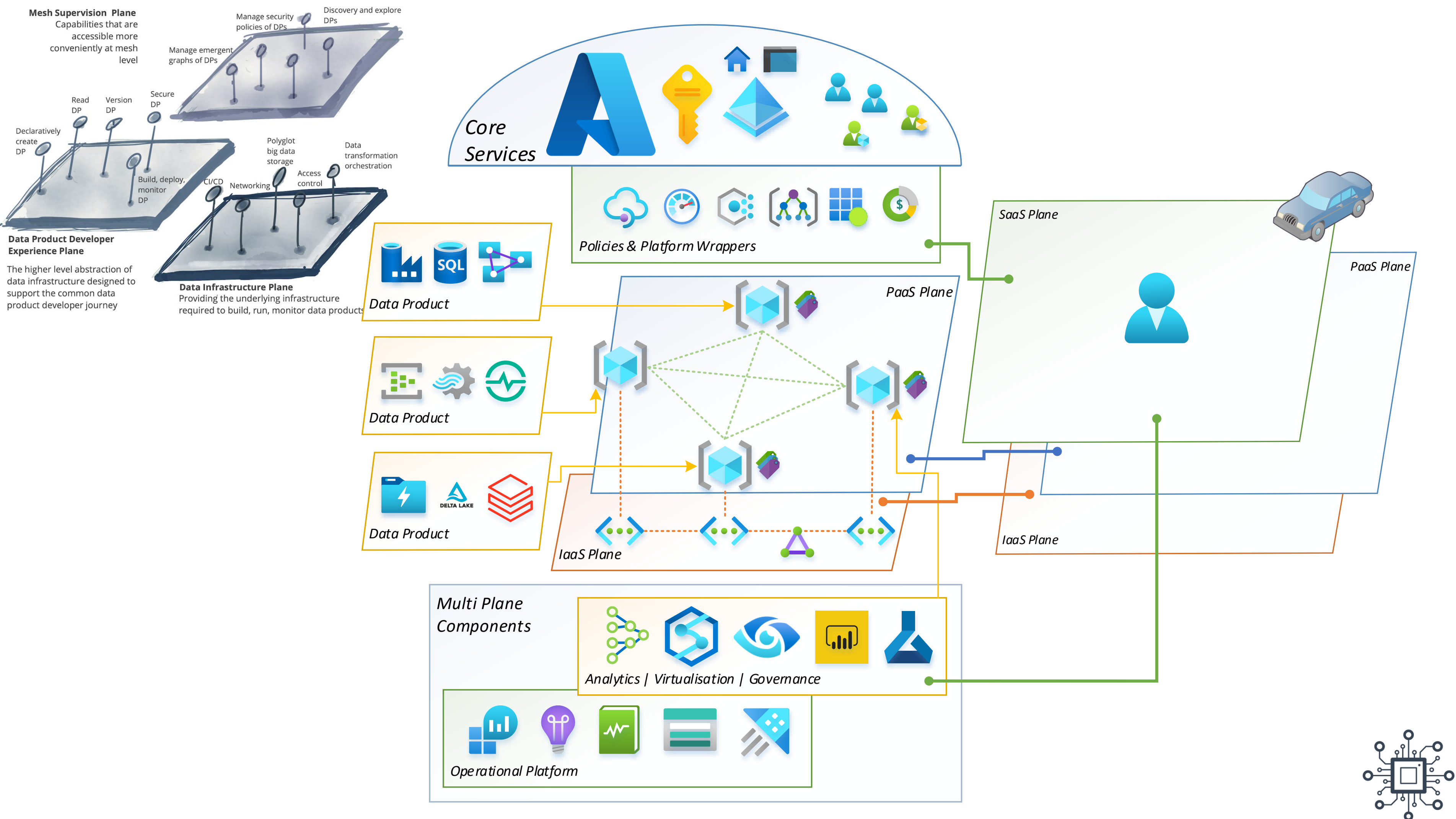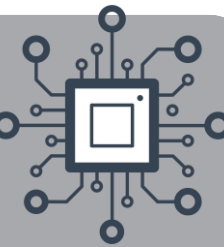
# An Evolution of Data Platform Architectures in Azure

Lambda, Kappa, Delta, Data Mesh

$\lambda$     $\kappa$     $\delta$