

An Evolution of Data Platform Architectures

Lambda, Kappa, Delta, Mesh & Fabric

λ

κ

δ

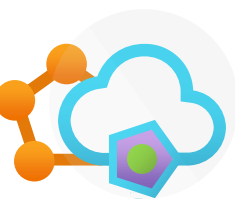


Paul Andrew

Technical Architect | Director



Cloud Formations



Architecture Agenda:

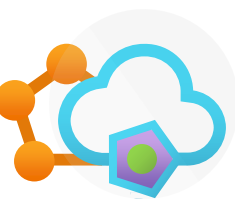
Lambda, Kappa, Delta, Mesh & Fabric

λ

κ

δ





But first, a couple of questions...



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.”



Volume

Velocity

Variety

Veracity

Value

- Buck Woody

@BuckWoodyMSFT

What is the goal of our data solutions?



Data
Sources

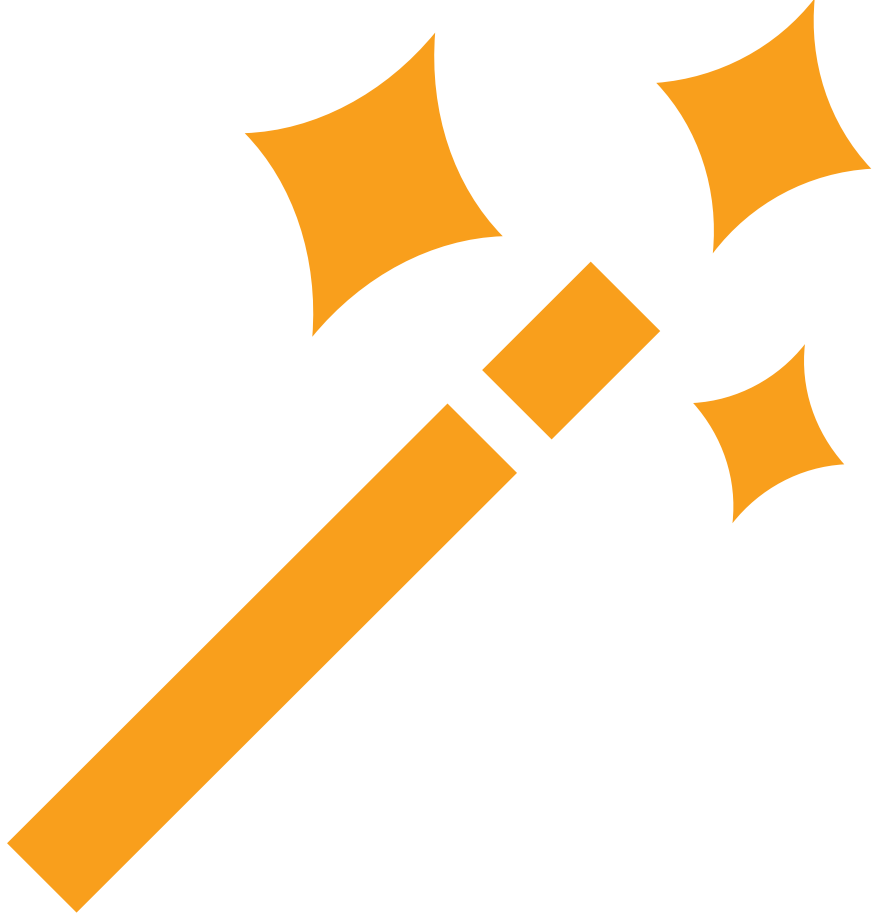
Data
Insight

Data = Information = Knowledge = Power

How do we deliver our data insights?



Data Sources



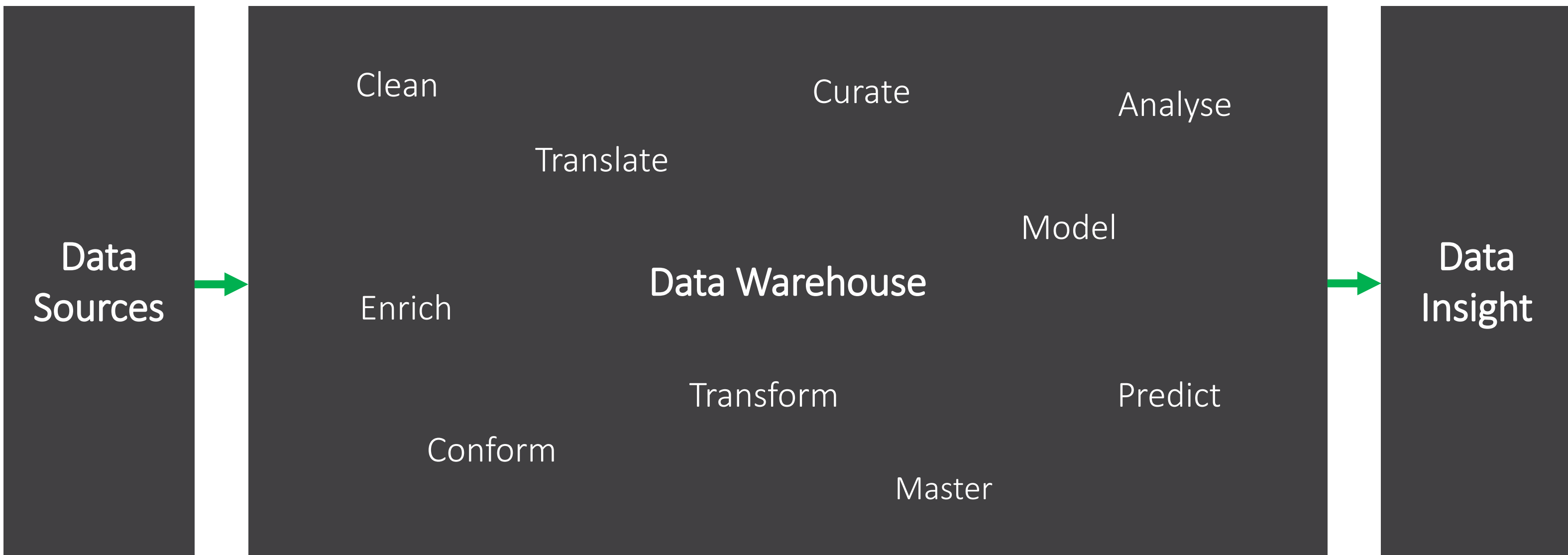
Data Insight

Data = Information = Knowledge = Power

How do we deliver our data insights?

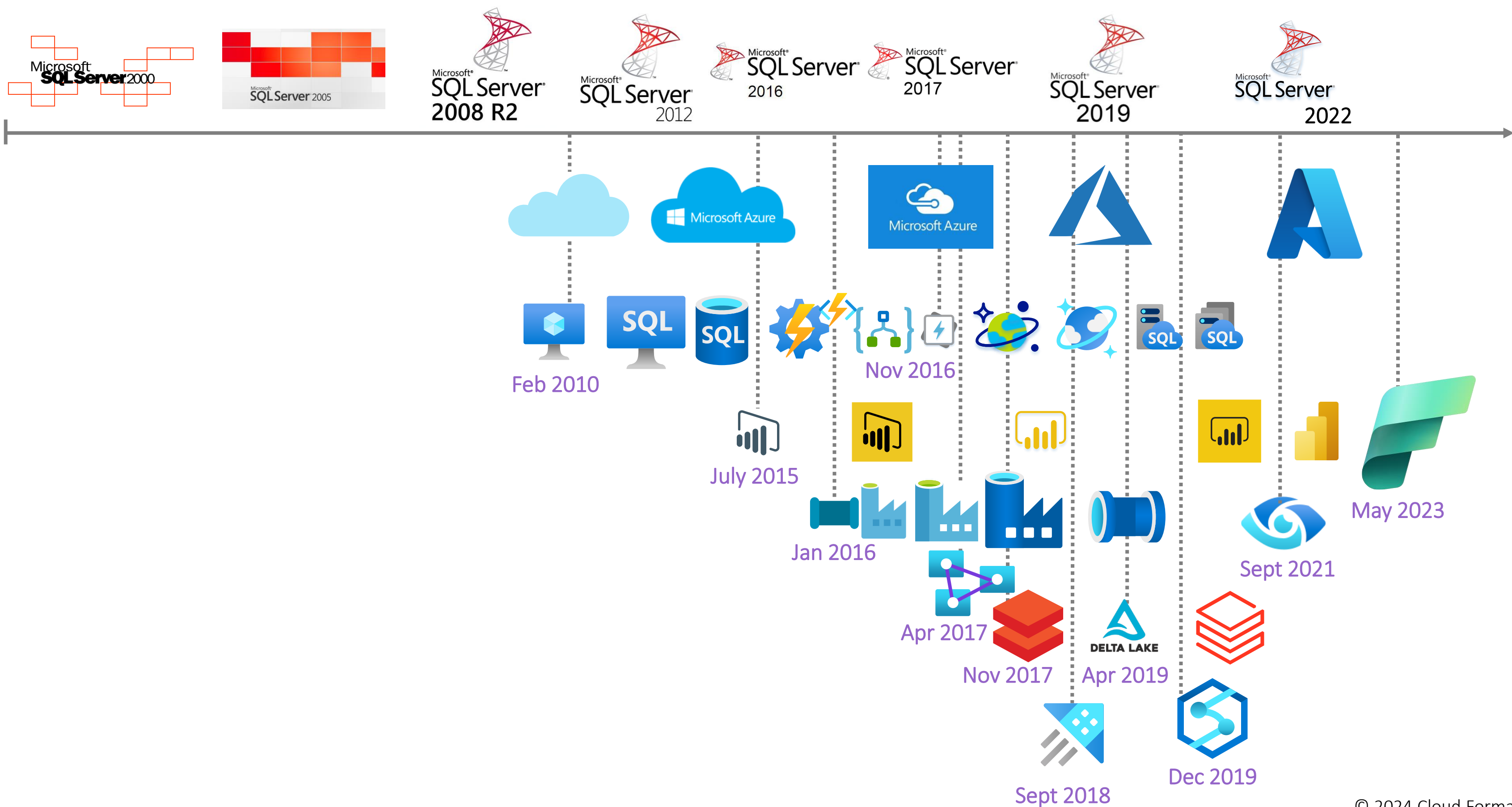


Cloud Formations - Knowledge Transfer & Training



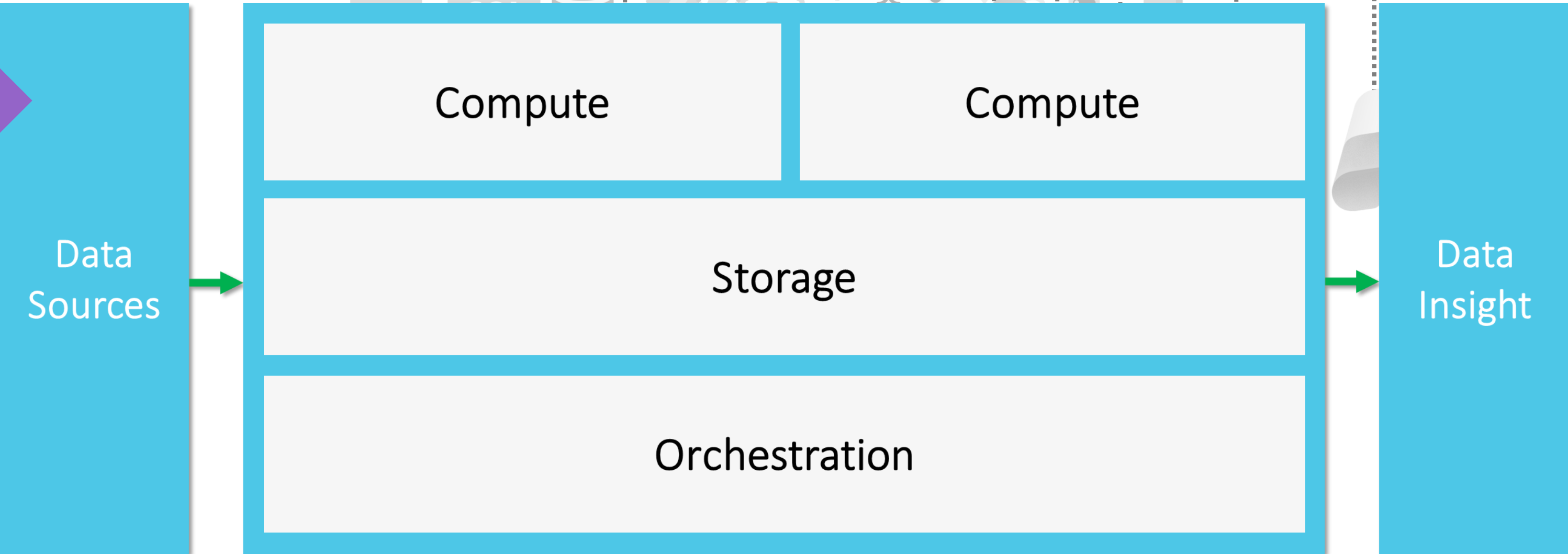
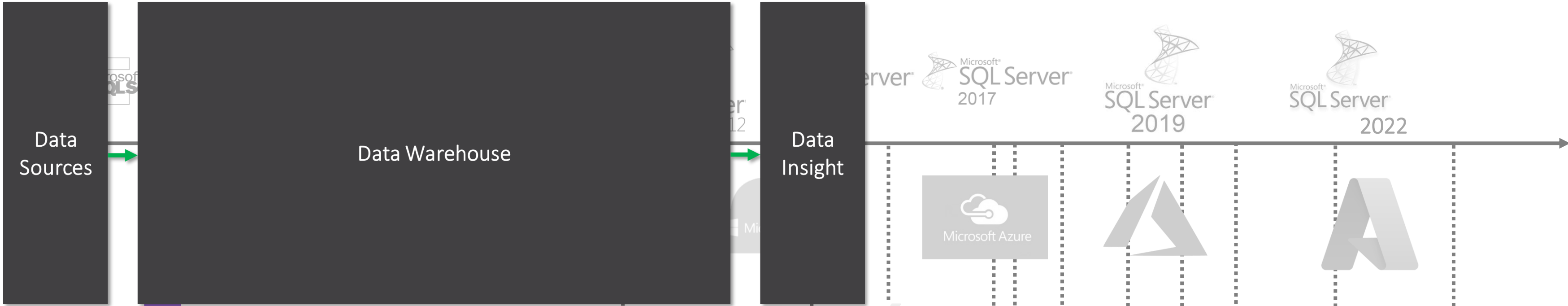
Data = Information = Knowledge = Power

An Evolution of Microsoft Data Technology



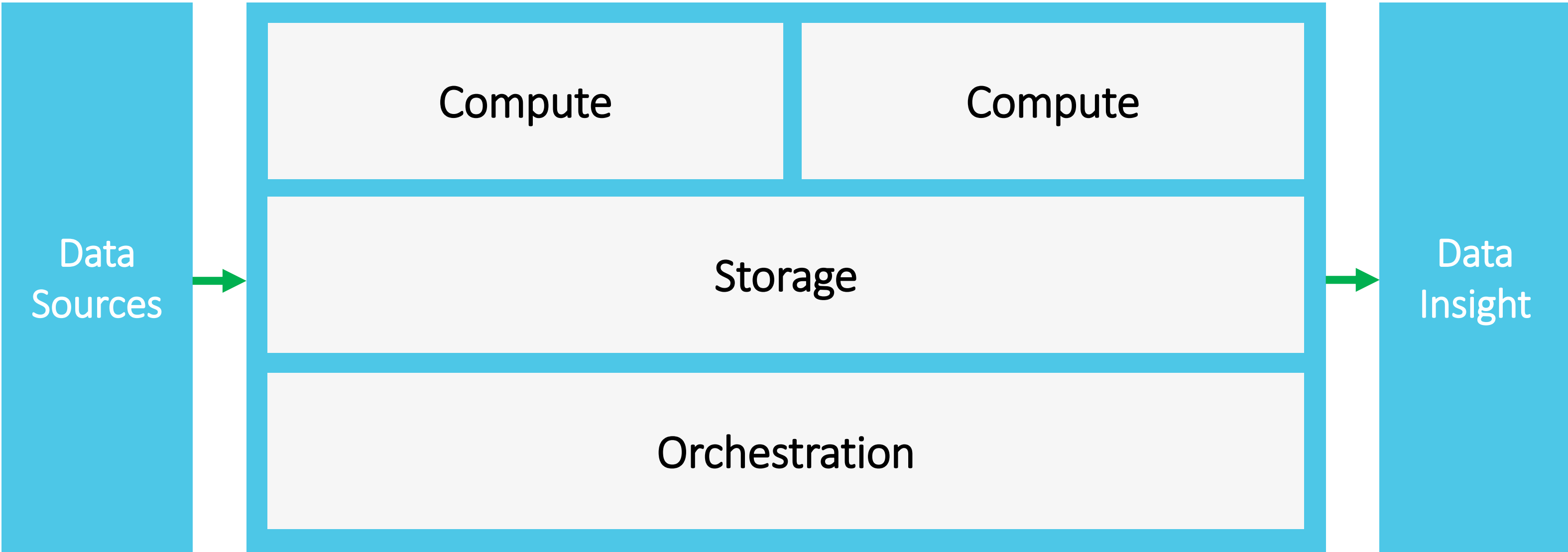
Cloud Formations - Knowledge Transfer & Training

An Evolution of Microsoft Data Technology

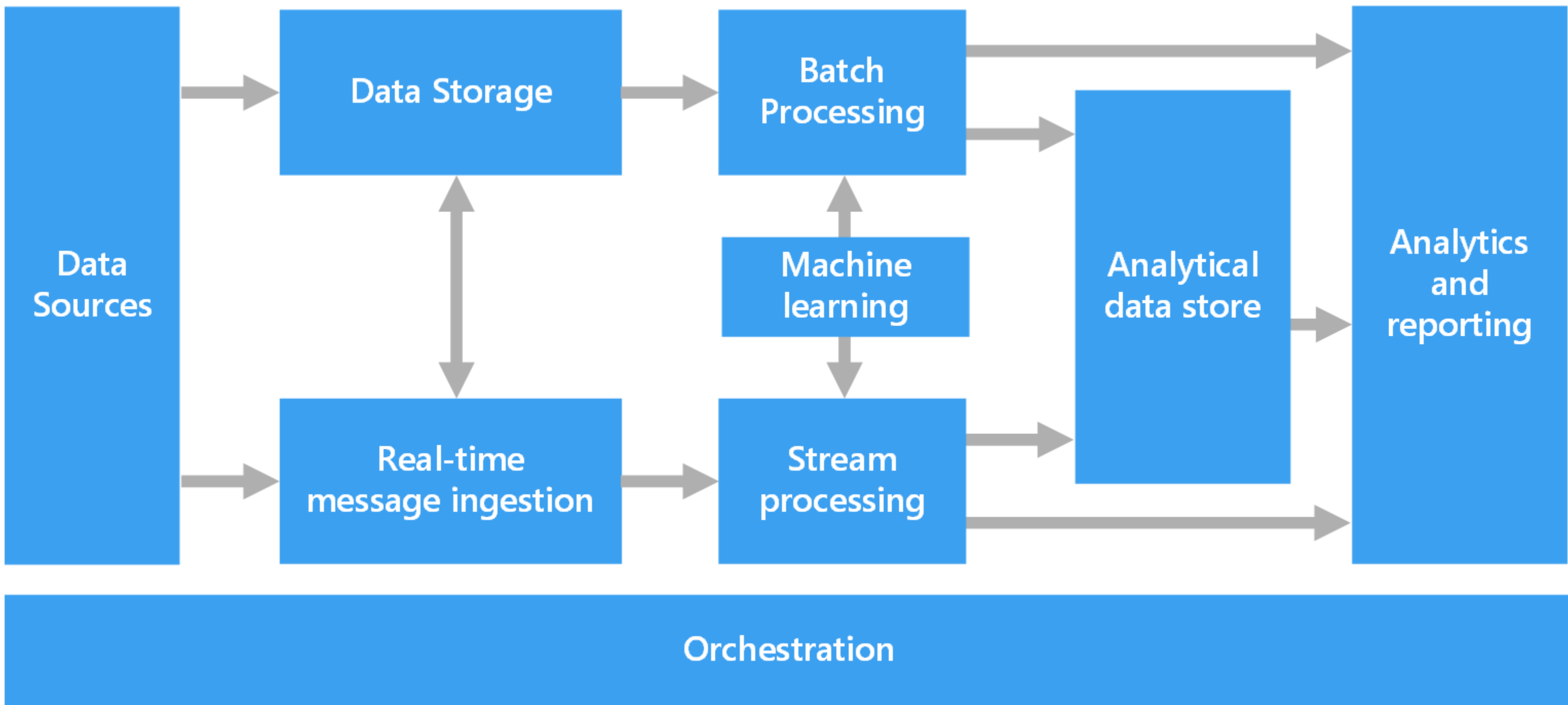


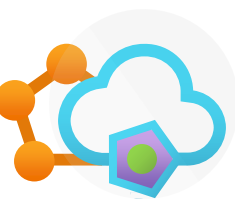
Cloud Formations - Knowledge Transfer & Training

My First Reference Architecture



Microsoft's Components of a Big Data Architecture





Architecture Agenda:

δ

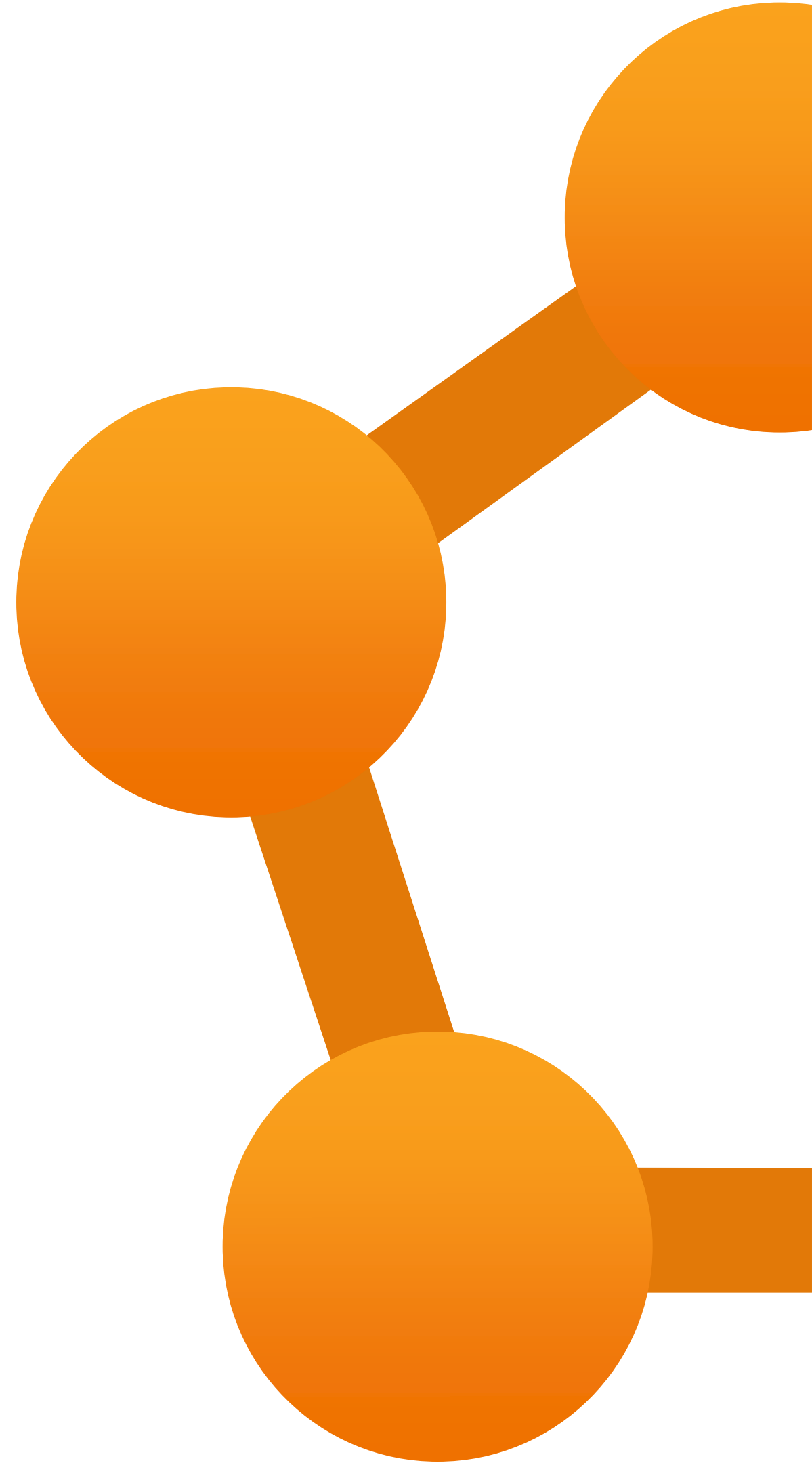
λ

K



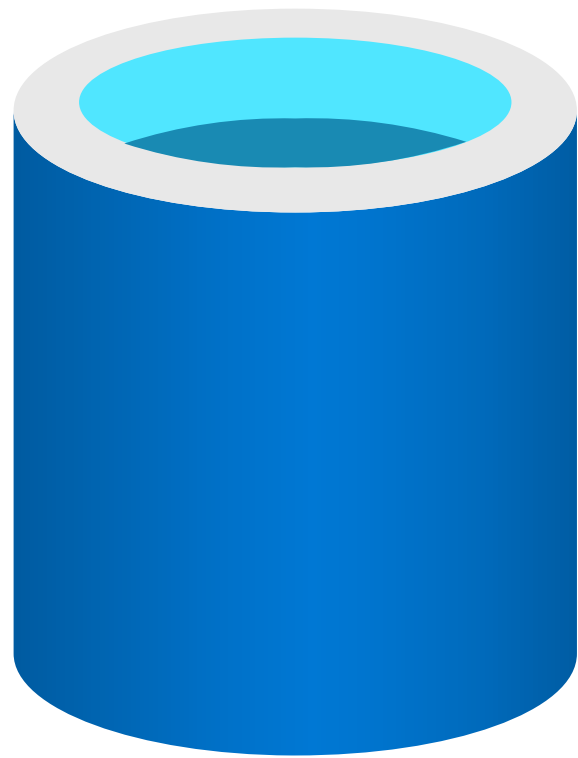
Delta δ

Cloud Formations





DataBase Management System



- Atomicity
- Consistency
- Isolation
- Durability



DataBase Management System

- Atomicity
- Consistency
- Isolation
- Durability

“is a set of properties of database transactions intended to guarantee data validity”

The screenshot shows the Wikipedia article for ACID. The article title is "ACID" and it is categorized as "From Wikipedia, the free encyclopedia". A notice indicates that the article needs additional citations for verification. The main text explains that in computer science, ACID (atomicity, consistency, isolation, durability) is a set of properties of database transactions intended to guarantee data validity despite errors, power failures, and other mishaps. It also mentions that in 1983, Andreas Reuter and Theo Härder coined the acronym ACID, building on earlier work by Jim Gray.

<https://en.wikipedia.org/wiki/ACID>

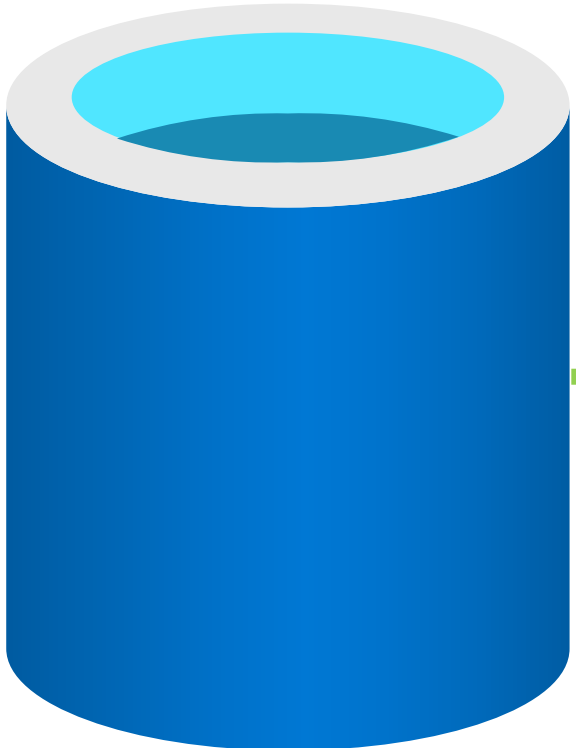
Databases



Creating a Data Warehouse



Online
Line
Transactional
Processing



Application
Data



Extract
Transform
Load



Data
Warehouse



Offline
Analytical
Transactional
Processing

Creating a Data Warehouse

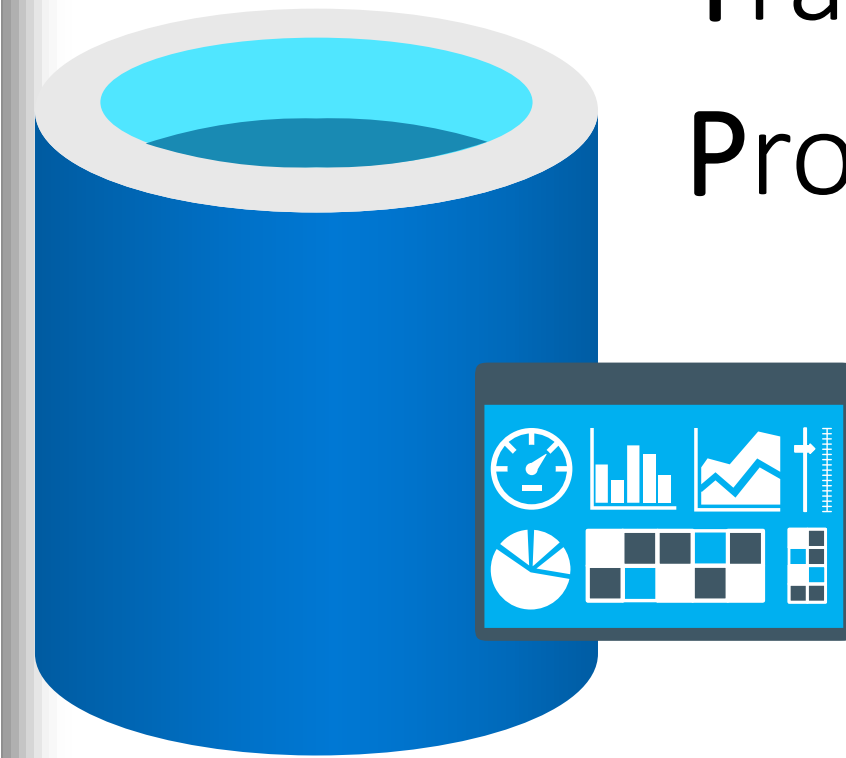


“a system for reporting and data analysis”

Offline
Analytical
Transactional
Processing

The screenshot shows the Wikipedia article for "Data warehouse". The article text states: "In computing, a **data warehouse (DW or DWH)**, also known as an **enterprise data warehouse (EDW)**, is a system used for reporting and data analysis and is considered a core component of business intelligence.^[1] DWs are central repositories of integrated data from one or more disparate sources. They store current and historical data in one single place^[2] that are used for creating analytical reports for workers throughout the enterprise.^[3]"

The article also includes a "Data warehouse overview" diagram and a "The basic architecture of a data warehouse" diagram. The latter diagram shows the flow from Data Sources (Operational systems, Flat files) through a Staging area to a Warehouse (Meta data, Summary data, Raw data), which then feeds into Data Marts (Purchasing, Sales) and finally to Users (Analysis, Reporting).



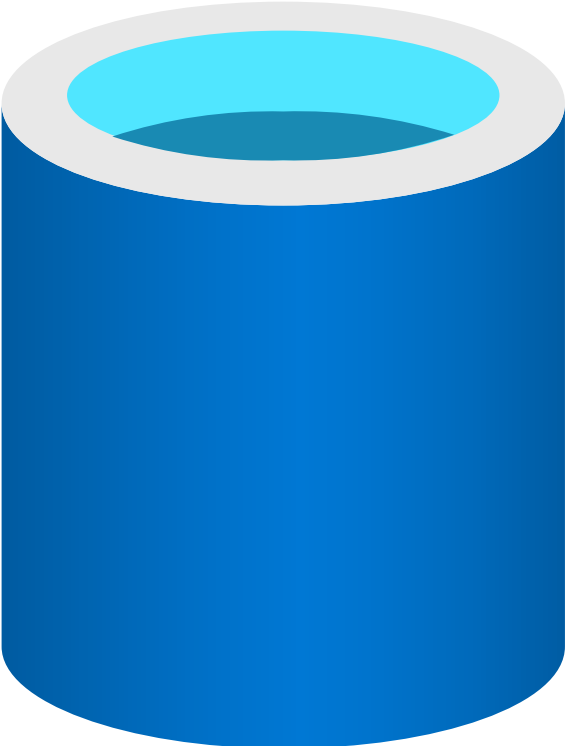
https://en.wikipedia.org/wiki/Data_warehouse





Big Data:

- Volume
- Velocity
- Variety
- Veracity
- Value

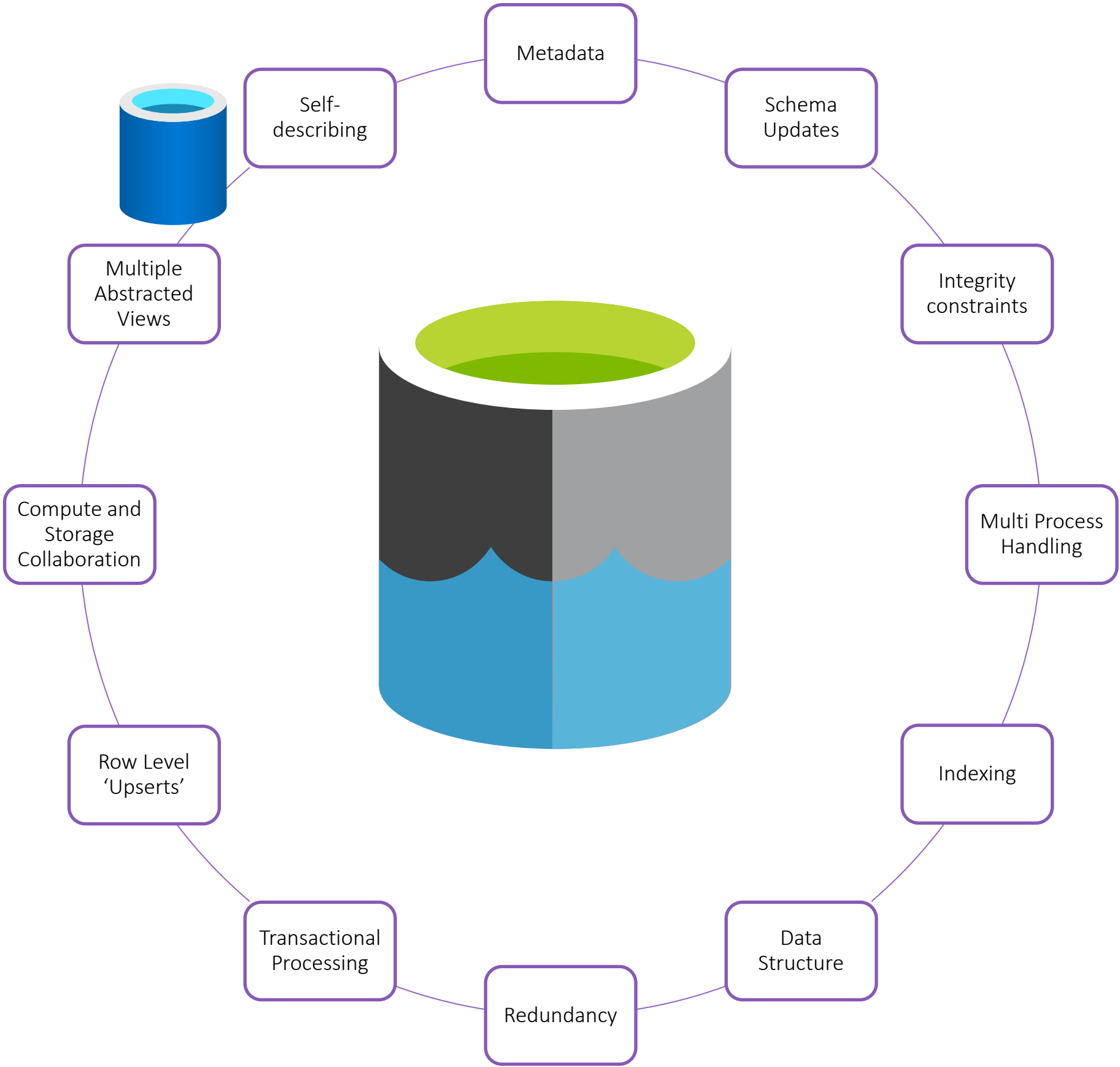


Data Warehouse

Data Lakes

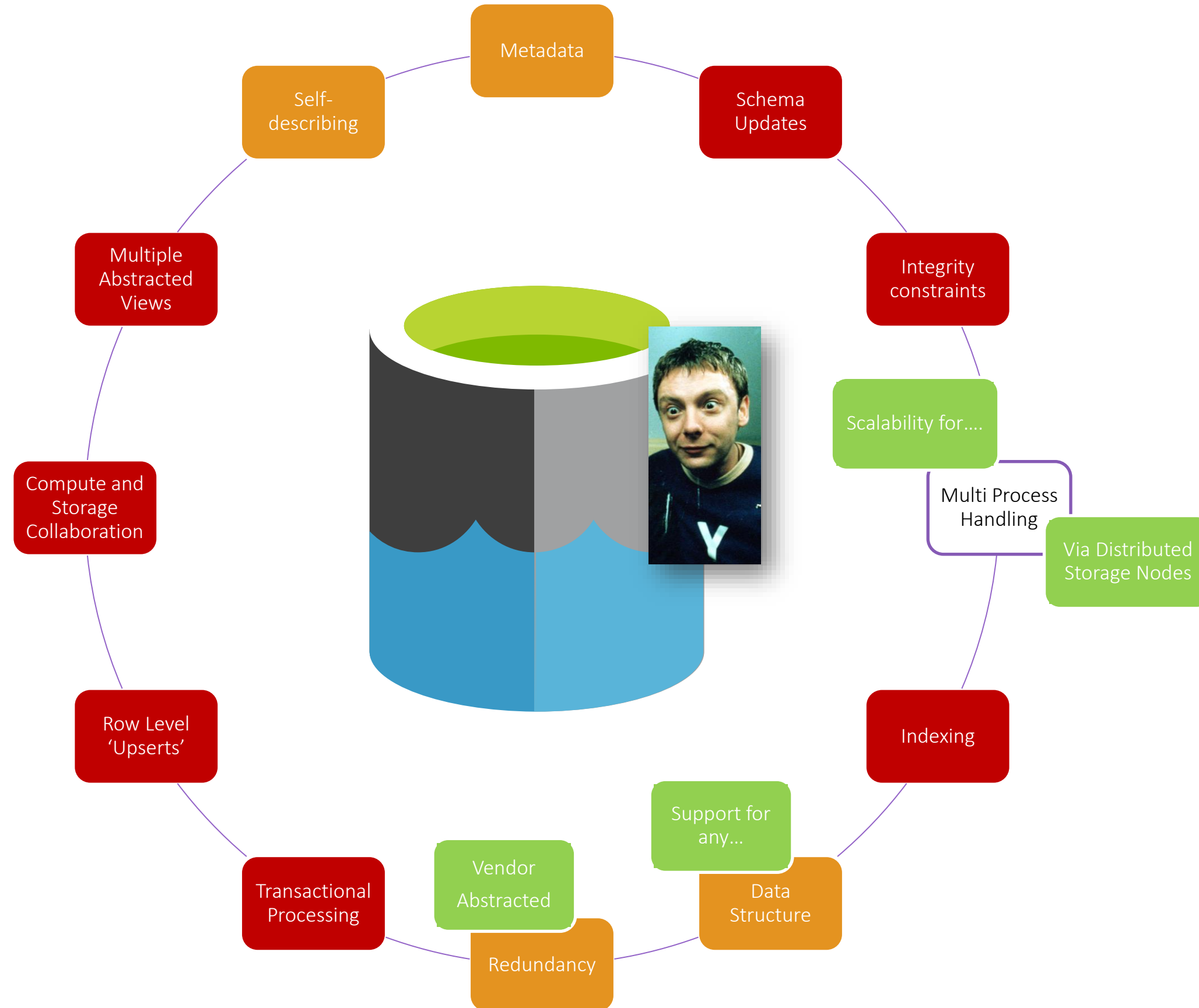


Big Data:
Volume
Velocity
Variety
Veracity
Value





Big Data:
Volume
Velocity
Variety
Veracity
Value



Problem Summary



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

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



VS



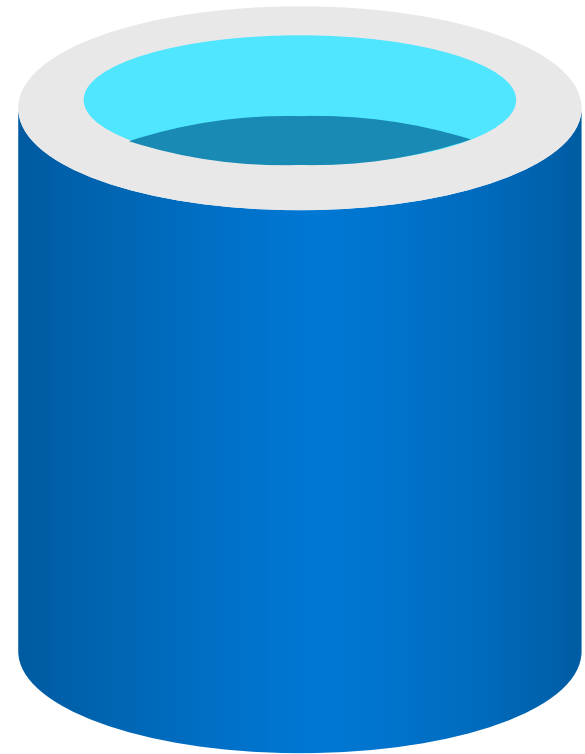
Scales Up	Scales Out
Natural Home for Structured Data	Any Data Structure
Storage Limits	No Storage Limits
Transactional Resilience	No Transactional Handling
Storage & Compute Coupled	Storage & Compute Decoupled

Problem Summary



Data Lakes are good, but they still lack some of the basic ACID 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



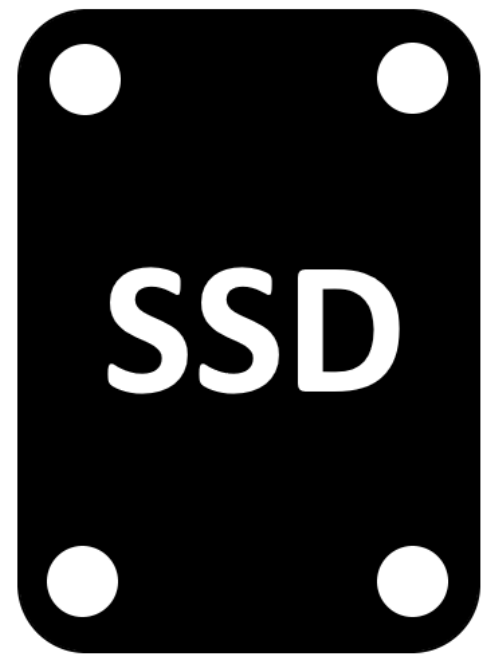
Solution

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

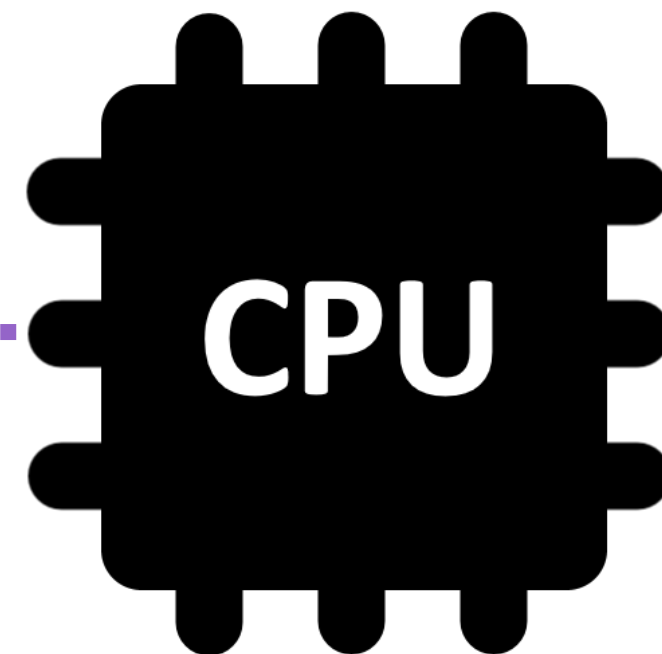
Big Data:

- Volume
- Velocity
- Variety
- Veracity
- Value

Storage

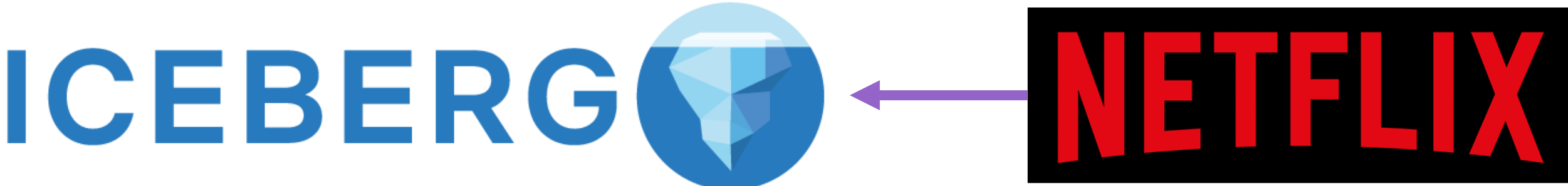
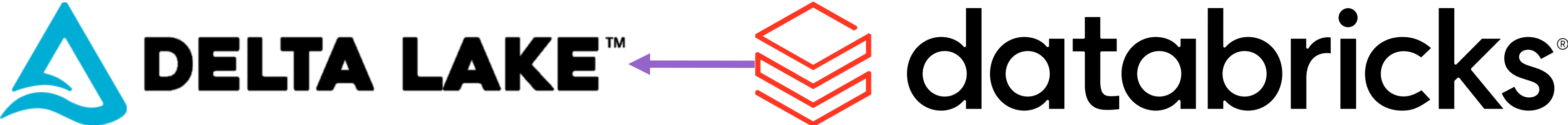


Compute



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

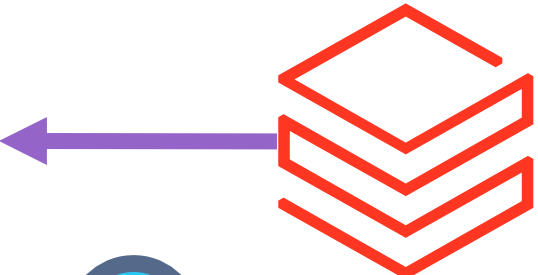
ACID Data Frameworks for Data Lakes



What is Delta Lake?

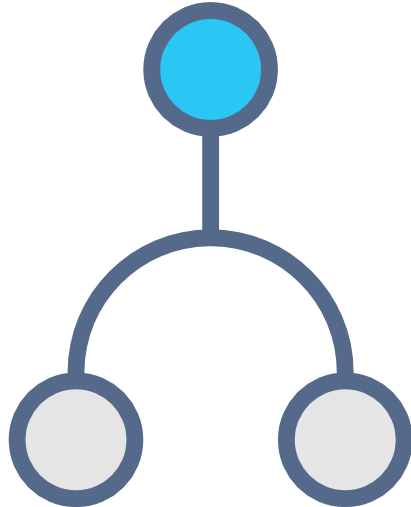


DELTA LAKE™

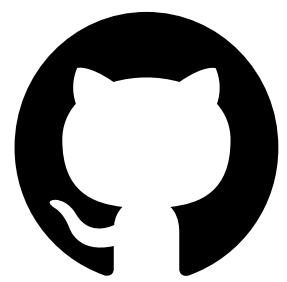


databricks®

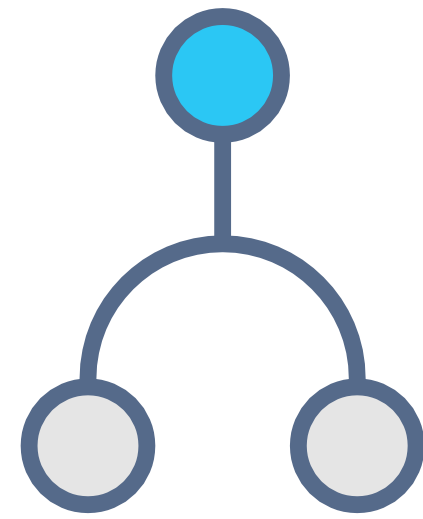
February 2019



What is Delta Lake?

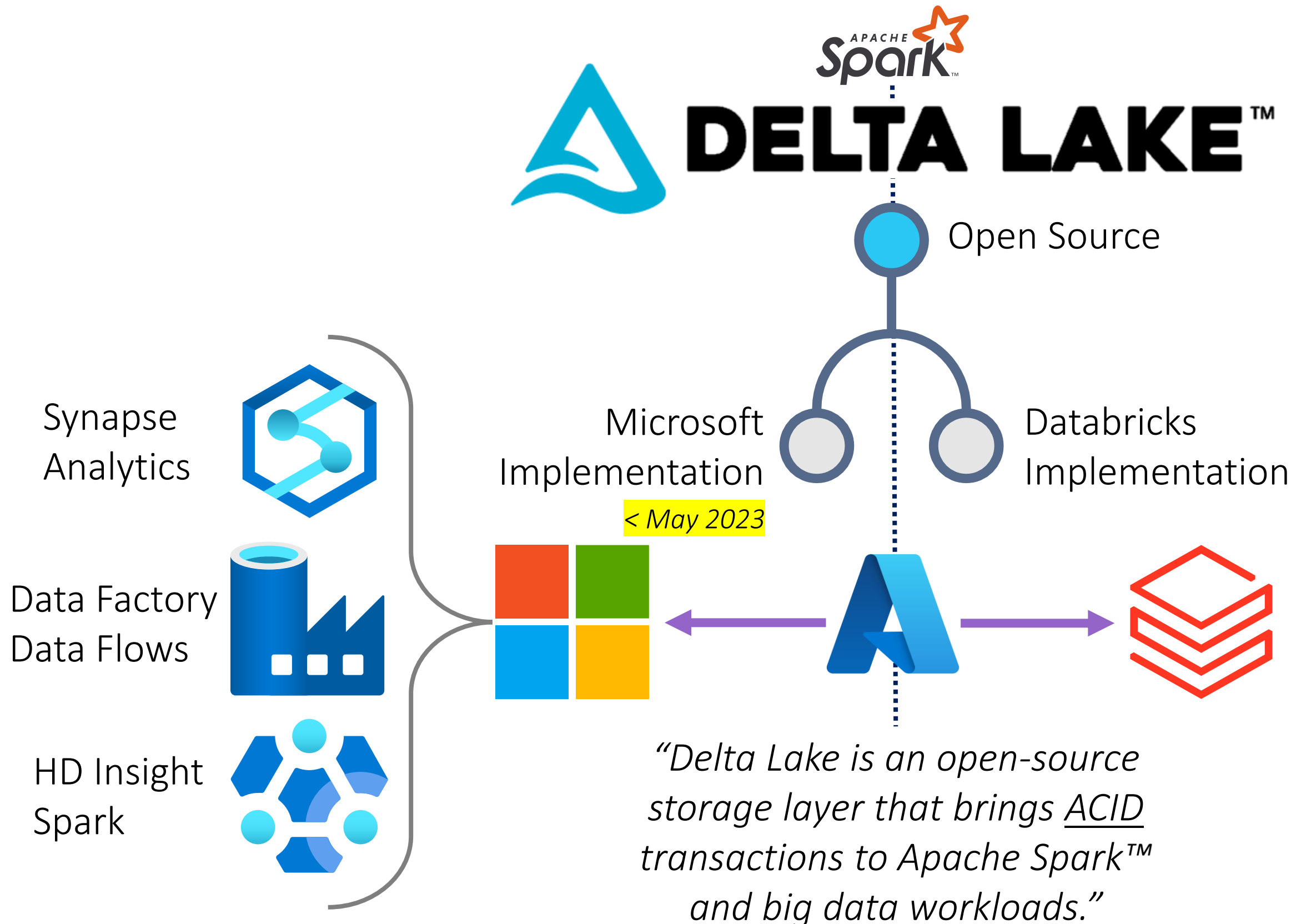


<https://delta.io>
<https://github.com/delta-io/delta>
April 2019

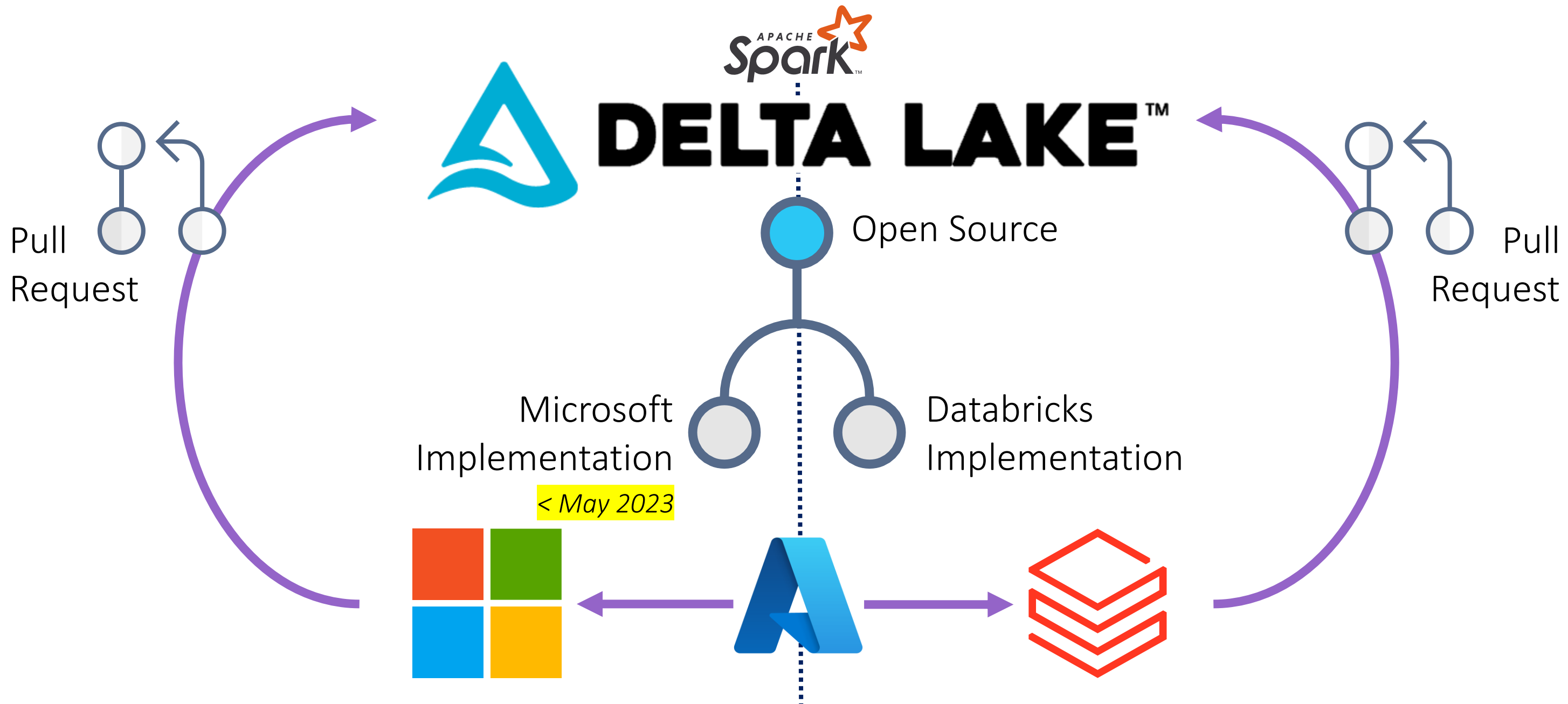


databricks®

What is Delta Lake?

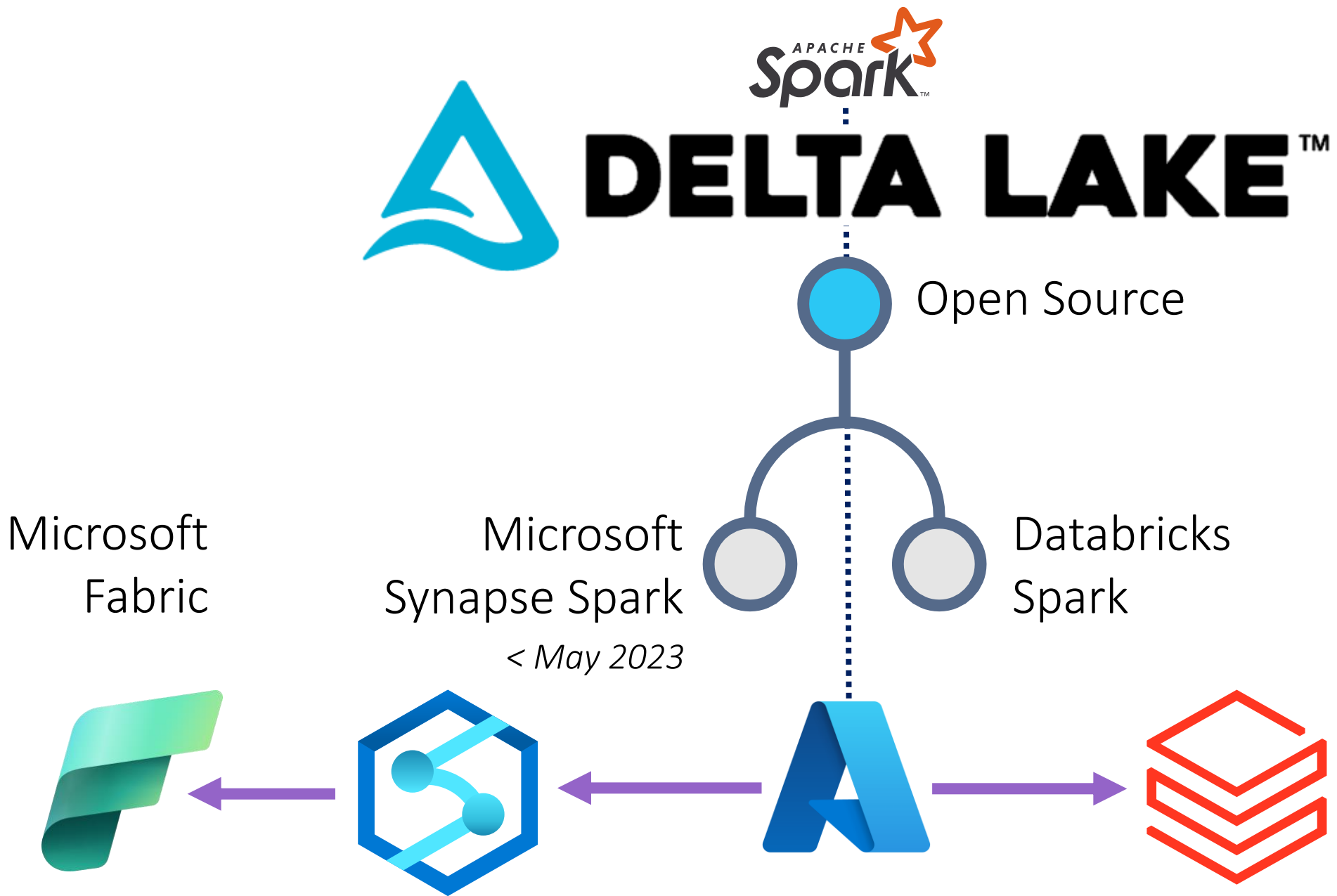


What is Delta Lake?



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

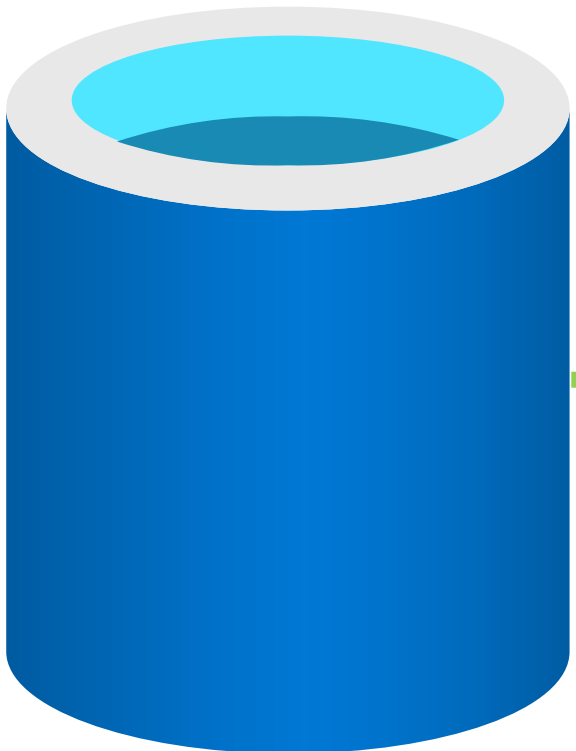
Which Spark Implementation is Better?



Data Warehouse



Online
Line
Transactional
Processing



Application
Data



Extract
Transform
Load



Data
Warehouse



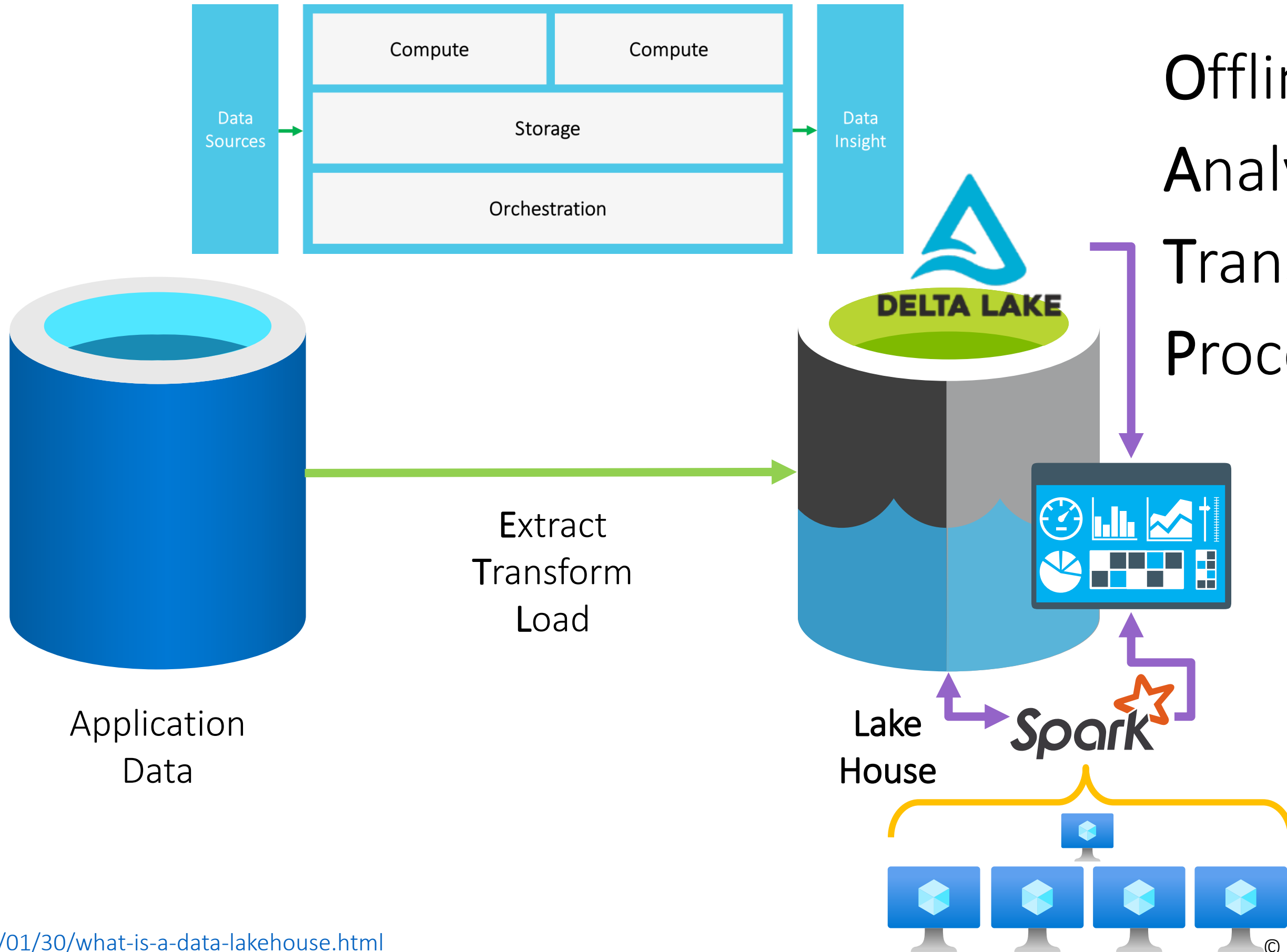
Offline
Analytical
Transactional
Processing

Lake House (Data-Ware-Lake-Delta-Beach-House-Lakes)



Cloud Formations - Knowledge Transfer & Training

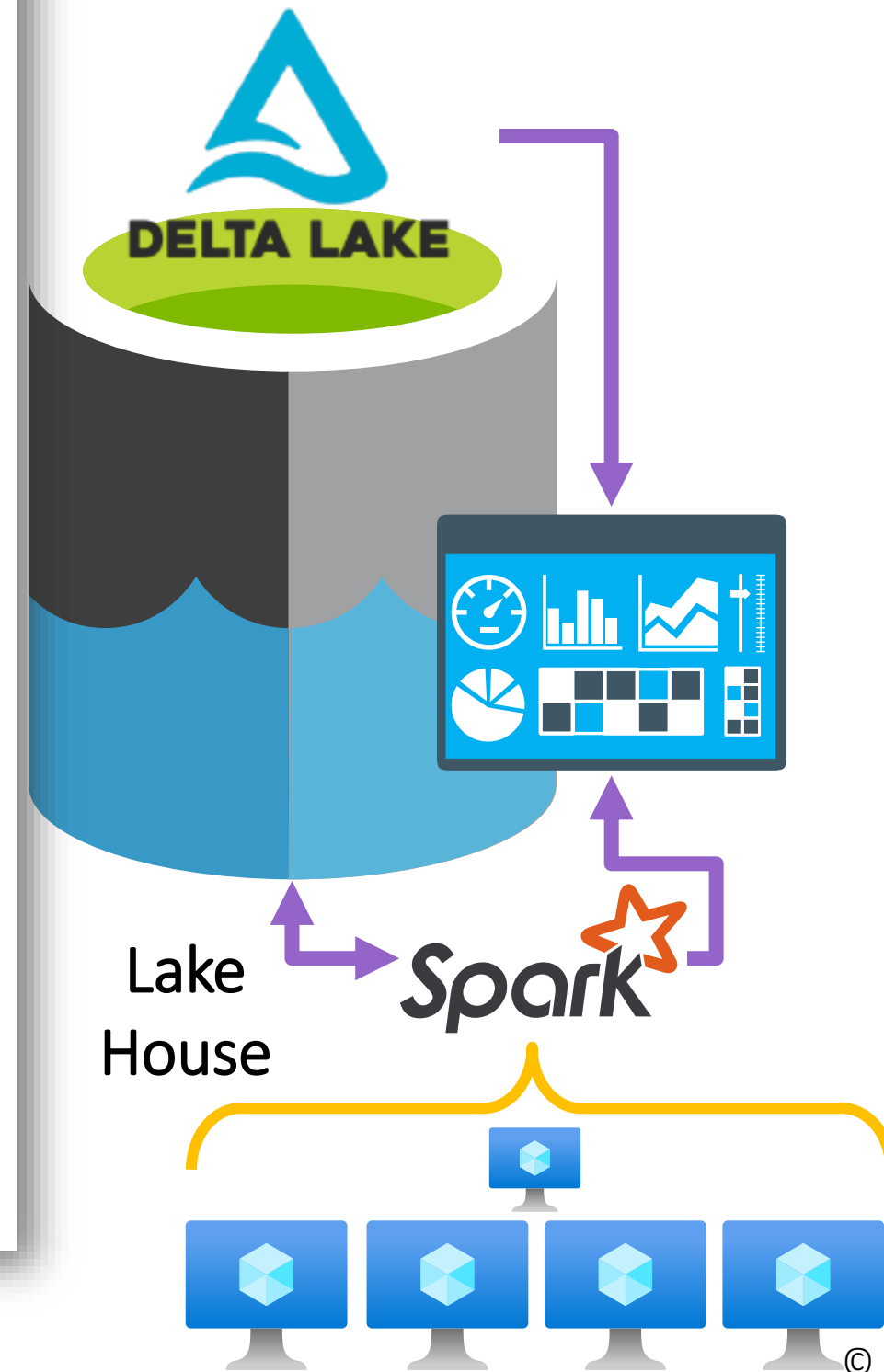
Online
Line
Transactional
Processing



Offline
Analytical
Transactional
Processing



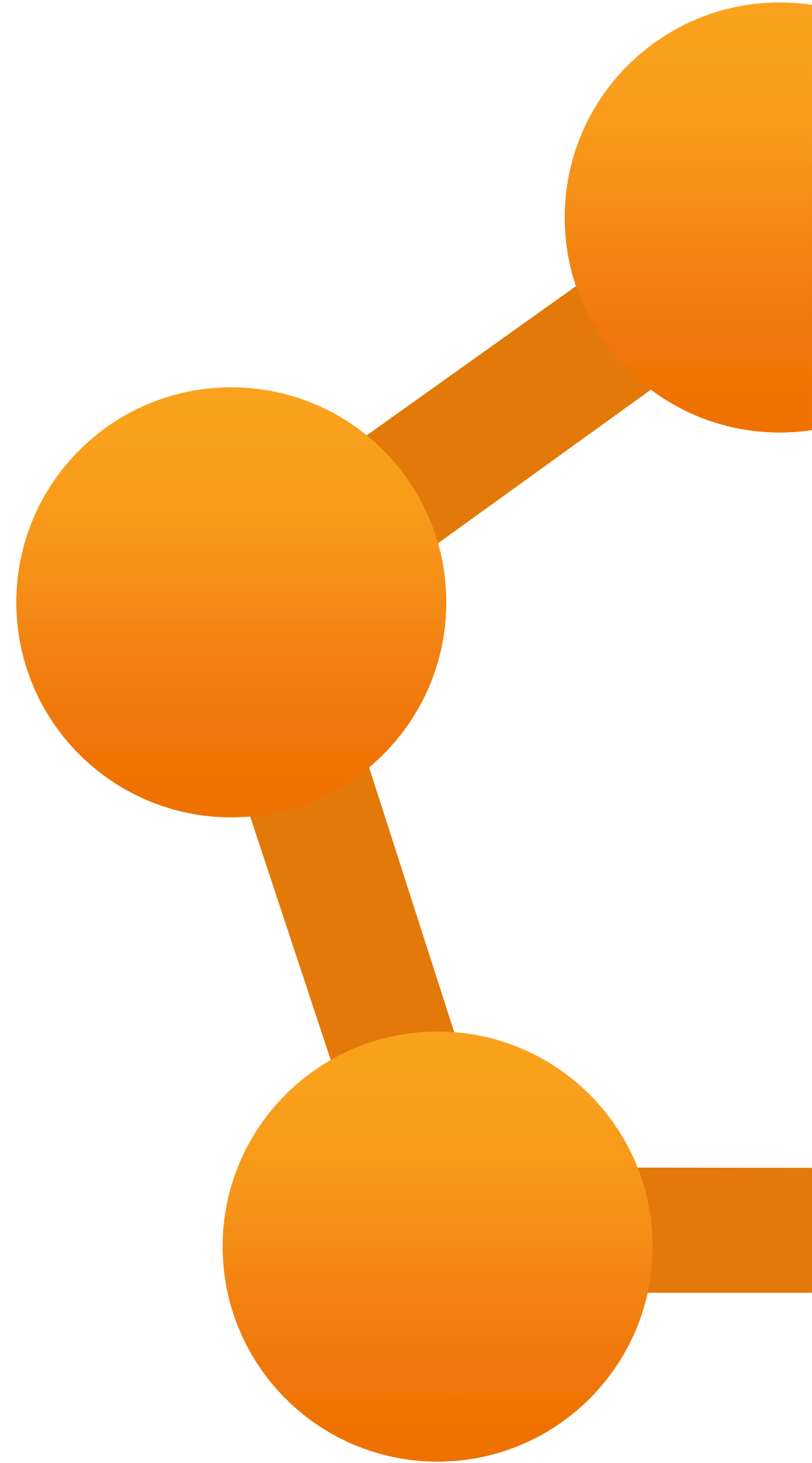
The screenshot shows the Wikipedia article for "The Lake House (film)". At the top, it says "Not logged in" with links for "Talk", "Contributions", "Create account", and "Log in". Below this are tabs for "Article" and "Talk", and a search bar. The article title is "The Lake House (film)" with a subtitle "From Wikipedia, the free encyclopedia". A yellow box contains a warning: "This article includes a list of general references, but it remains largely unverified because it lacks sufficient corresponding inline citations. Please help to improve this article by introducing more precise citations. (October 2017) (Learn how and when to remove this template message)". The main text describes the film as a 2006 American fantasy romantic drama directed by Alejandro Agresti, starring Keanu Reeves and Sandra Bullock. A "Contents" table of contents is visible on the left, listing sections like Plot, Cast, Production, Music, Reception, and References. A "Theatrical release poster" is shown on the right, featuring Keanu Reeves and Sandra Bullock. Below the poster is a list of credits: Directed by Alejandro Agresti, Written by David Auburn, Based on Il Mare by Kim Eun-jeong and Kim Mi-yeong, Produced by Doug Davison and Roy Lee, and Starring Keanu Reeves.



Lambda & Kappa

λ K

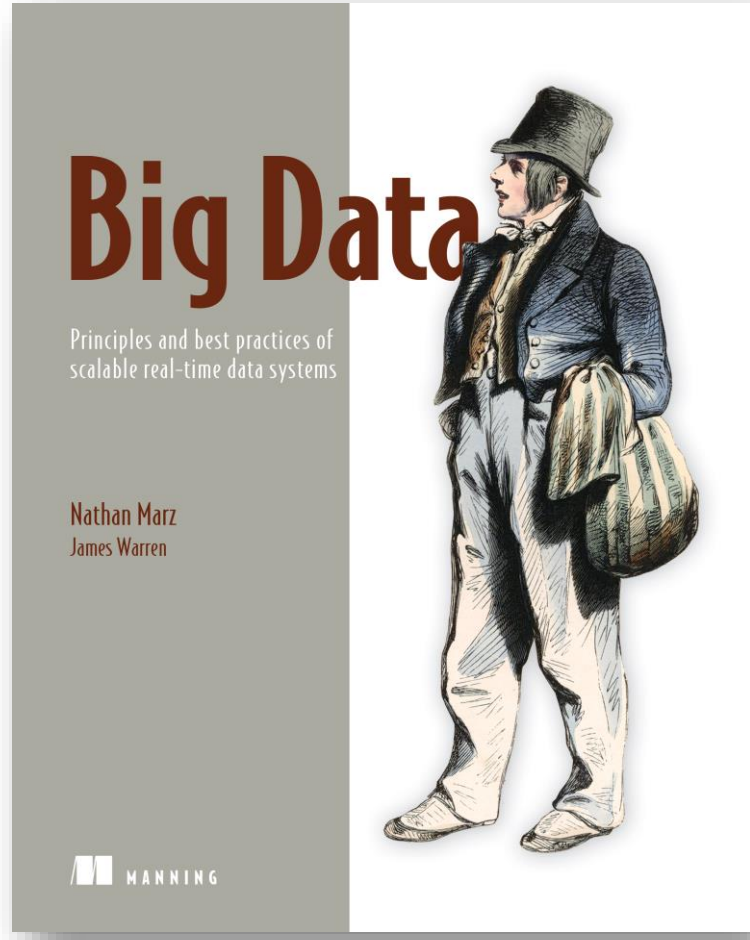
Cloud Formations



Lambda & Kappa Architectures



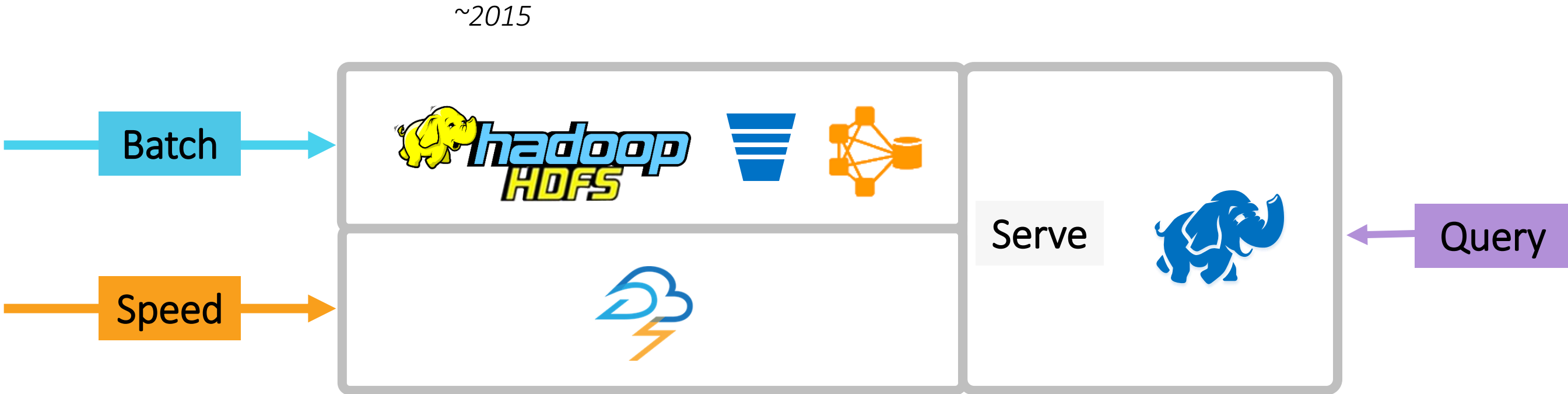
Cloud Formations - Knowledge Transfer & Training



Big Data:
Volume
Velocity
Variety
Veracity
Value



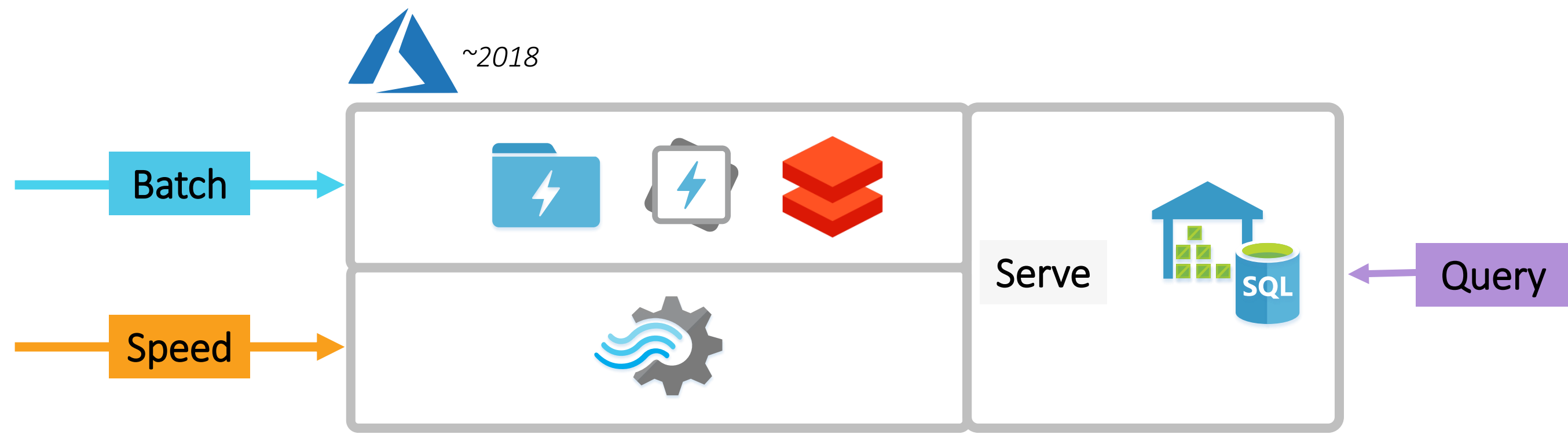
Lambda & Kappa Architectures



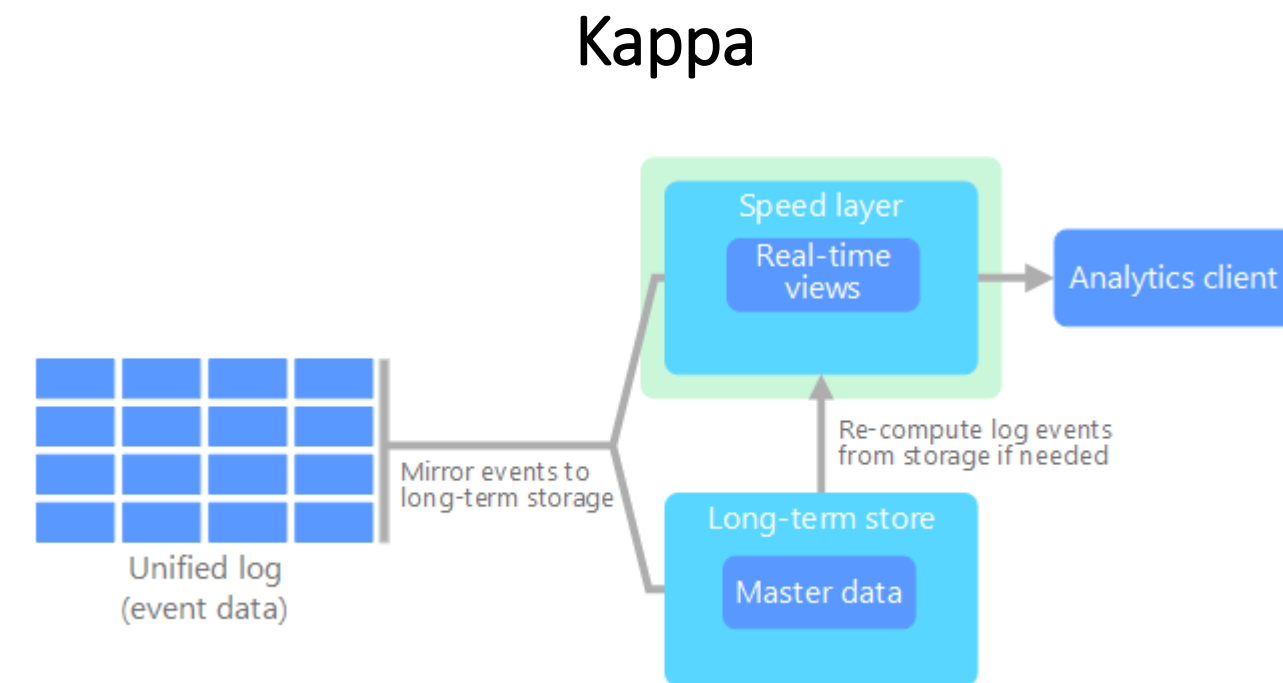
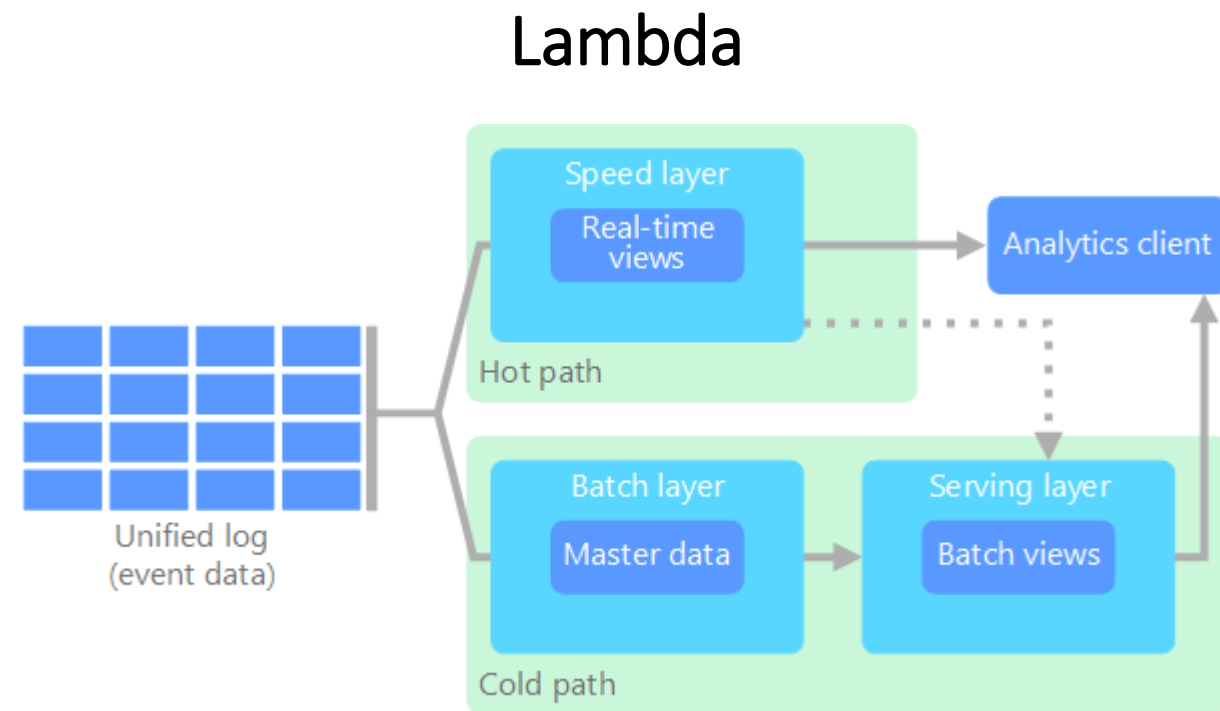
Lambda & Kappa Architectures



Cloud Formations - Knowledge Transfer & Training



Lambda & Kappa Architectures



“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:

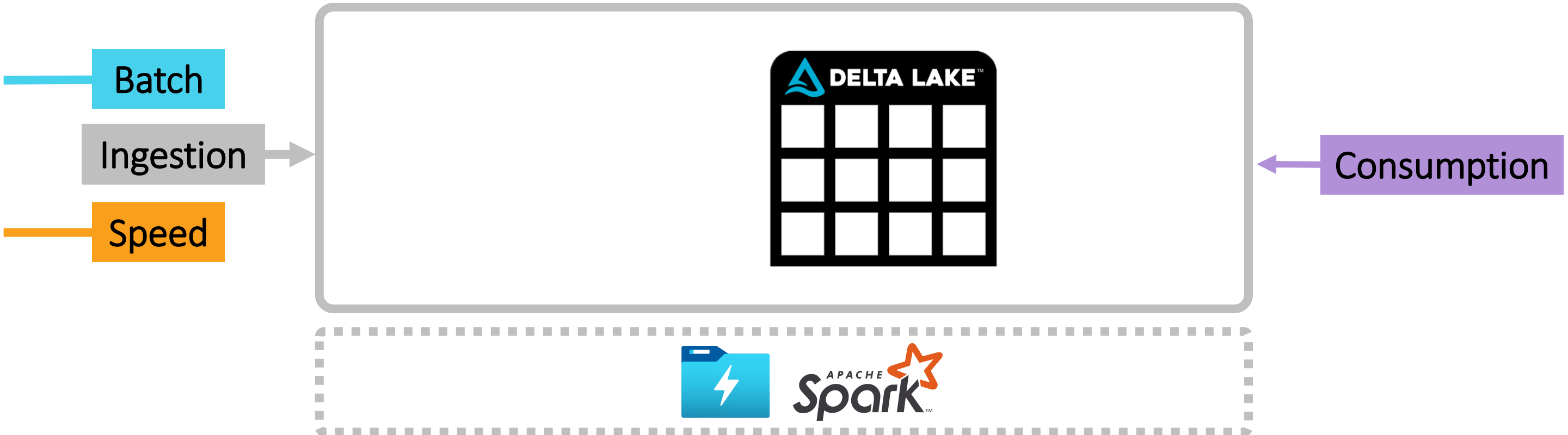
A **batch layer** (cold path) stores all of the incoming data in its raw form and performs batch processing on the data. The result of this processing is stored as a **batch view**.

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.”

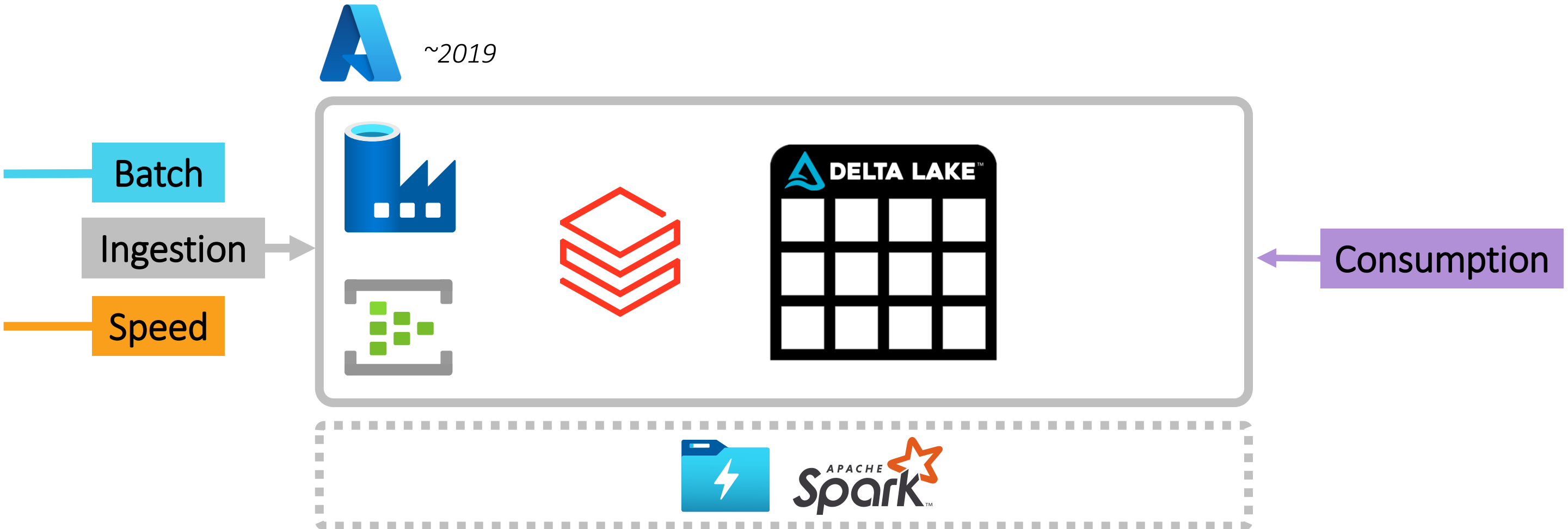
Lambda & Kappa Architectures



Lambda & Kappa Architectures



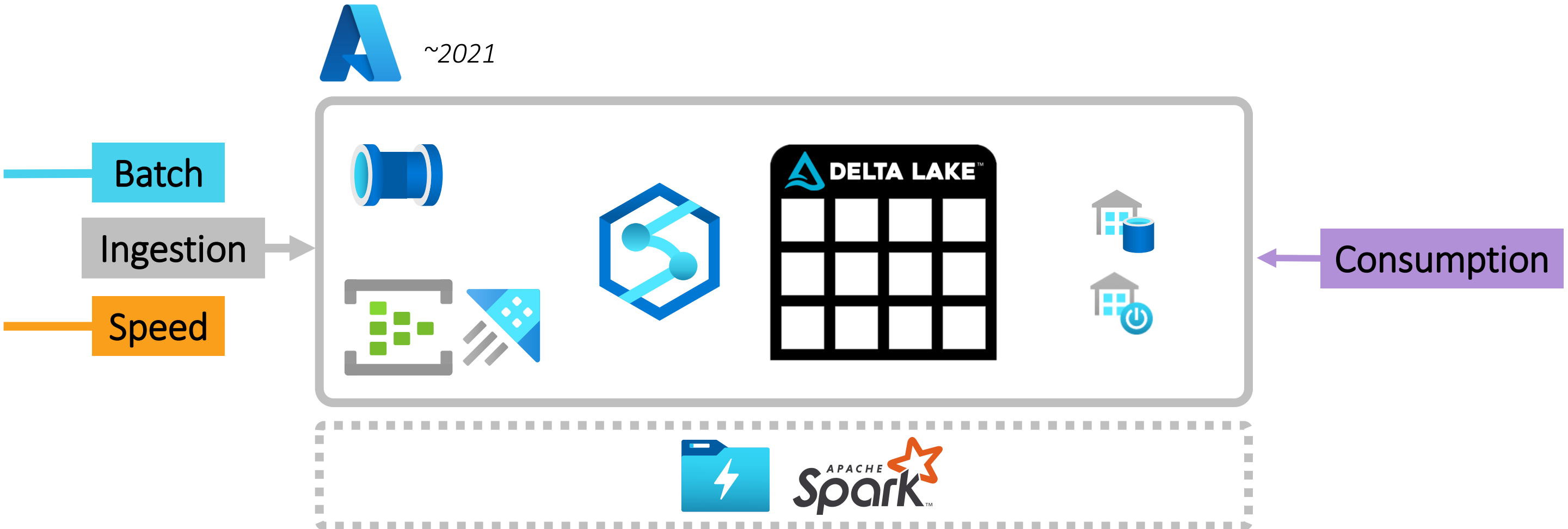
Cloud Formations - Knowledge Transfer & Training



Lambda & Kappa Architectures



Cloud Formations - Knowledge Transfer & Training

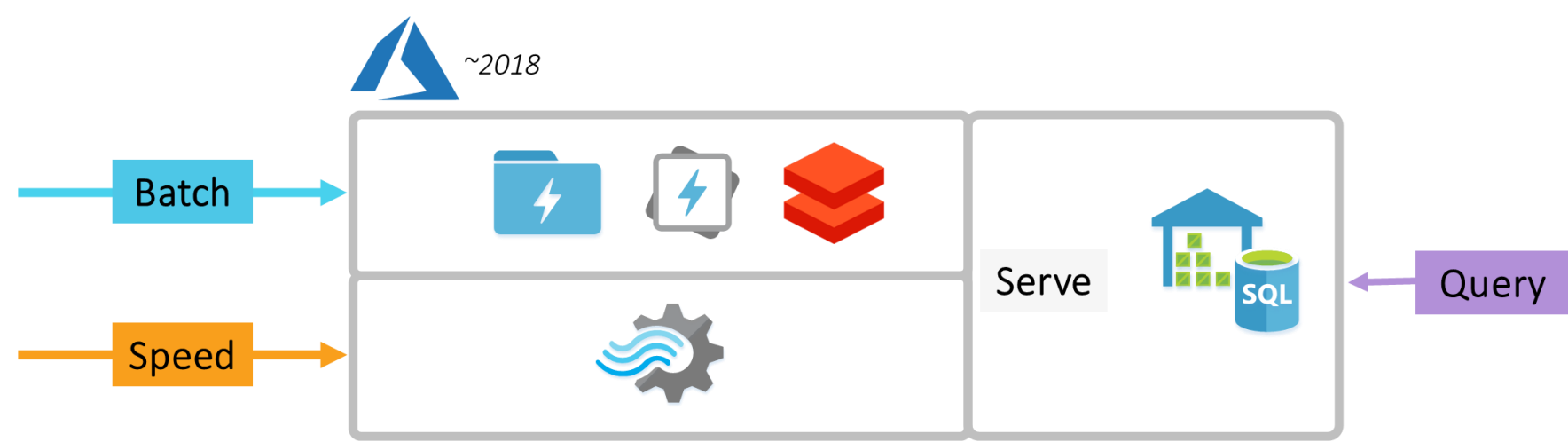
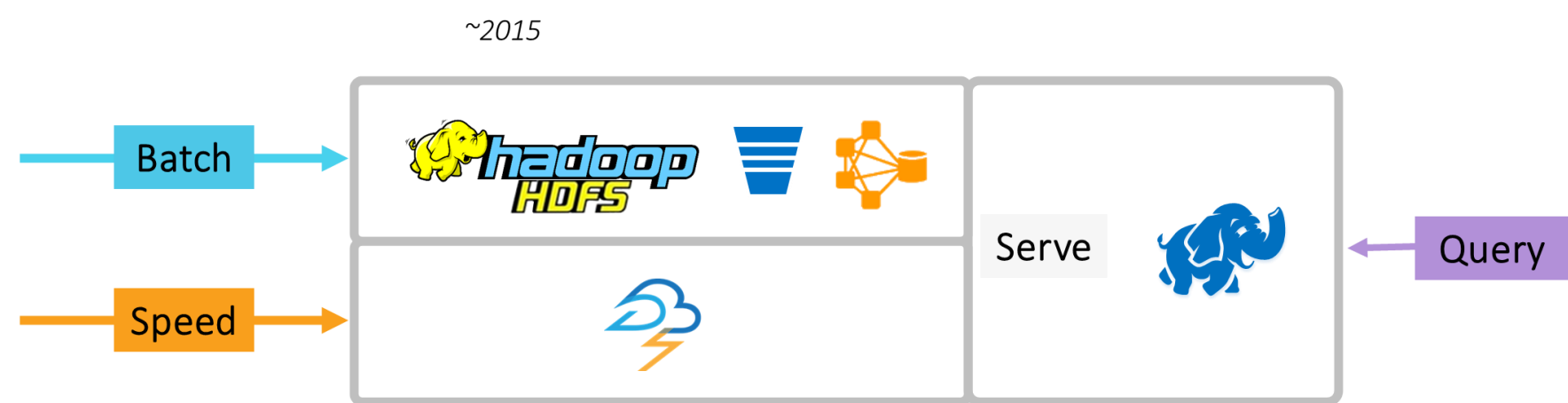


Lambda & Kappa Architectures vs Technology

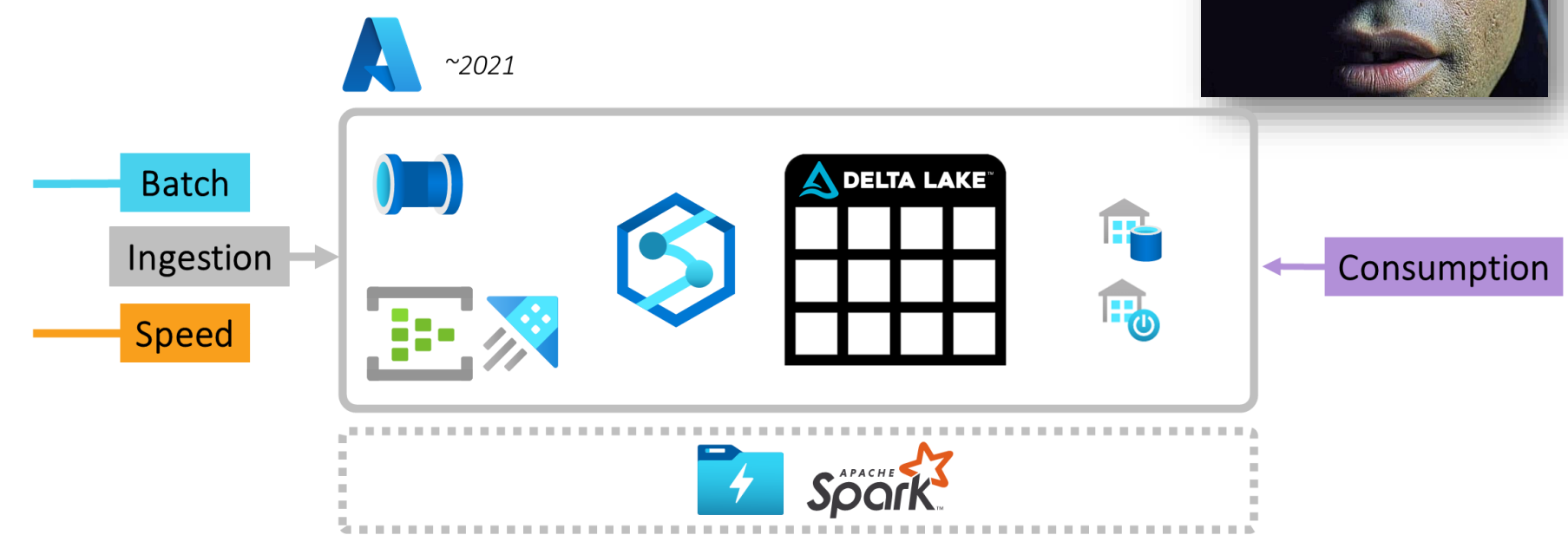
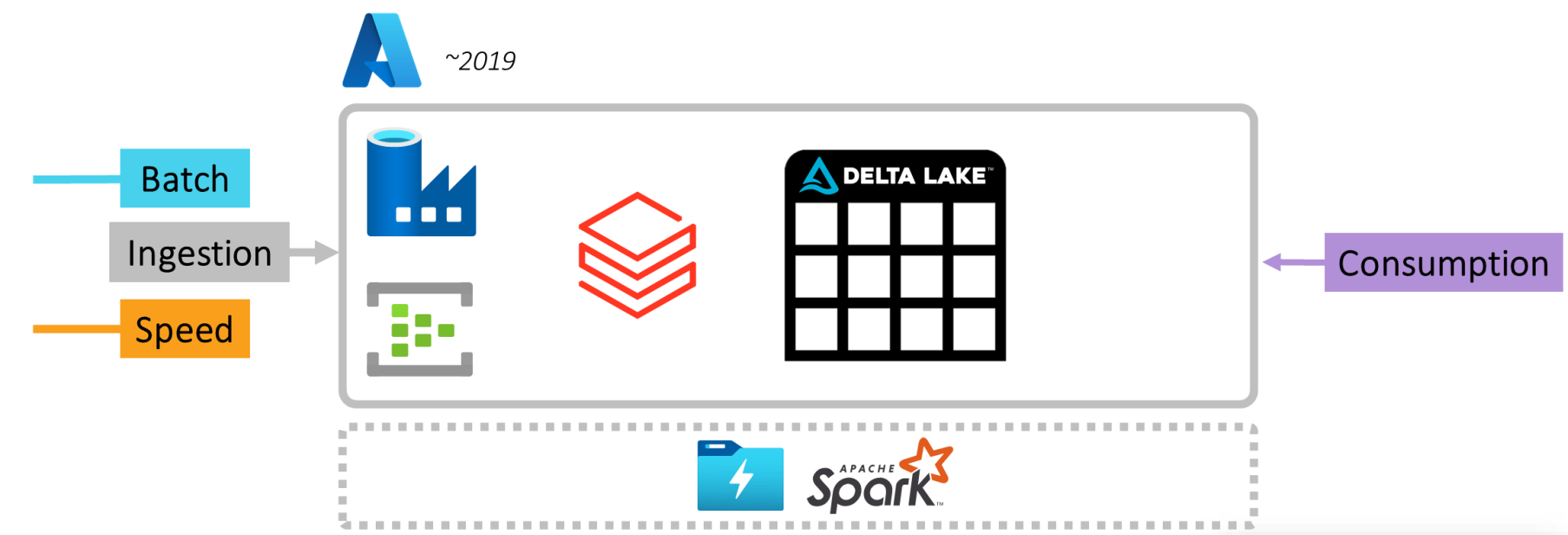


Cloud Formations - Knowledge Transfer & Training

Lambda

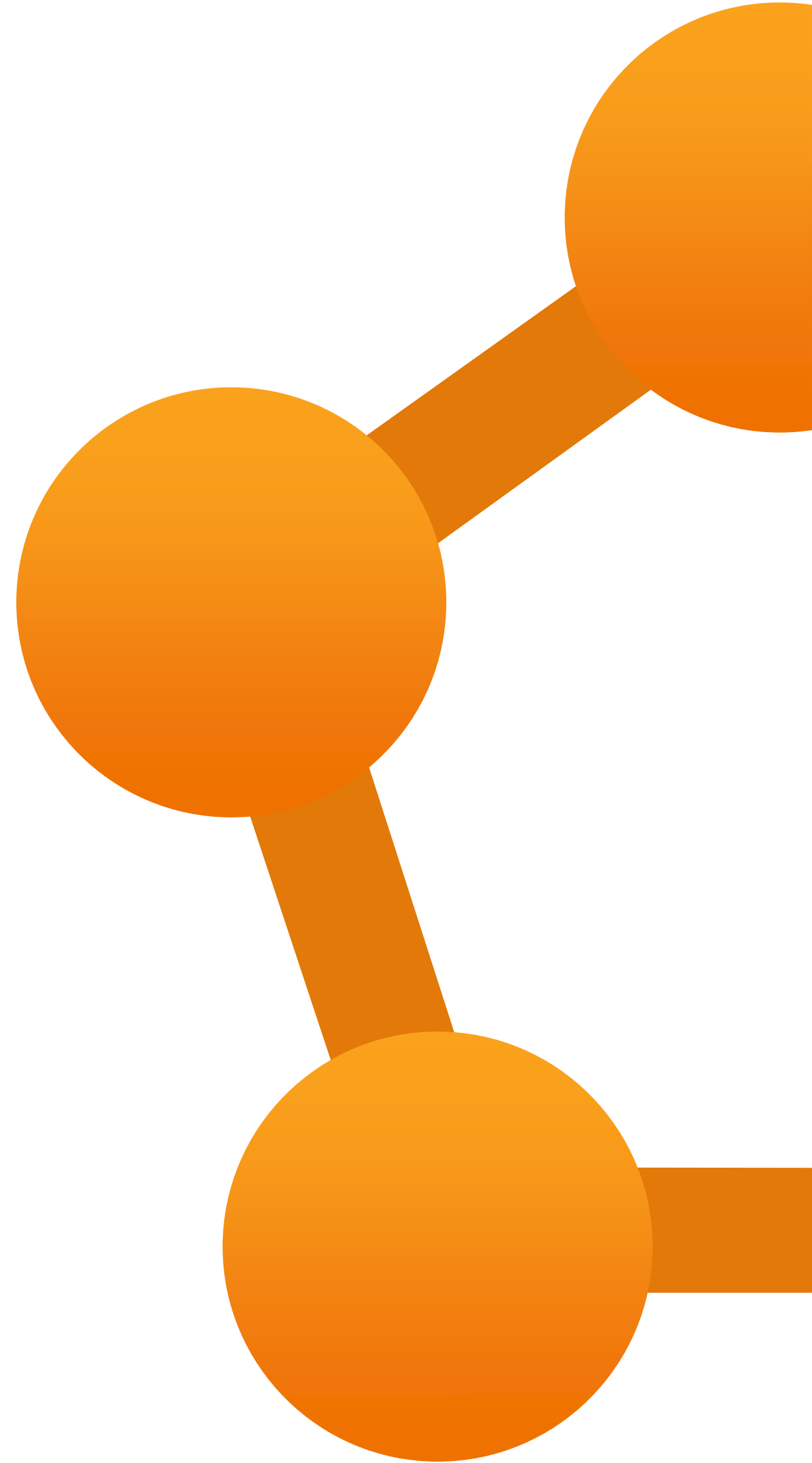


Kappa





Cloud Formations



Different Types of Fabric



Different Types of Fabric



The screenshot shows the Gartner website's Information Technology Glossary page for 'Data Fabric'. The page title is 'Gartner Glossary' under 'Information Technology'. The breadcrumb trail is 'Gartner Glossary > Information Technology Glossary > D > Data Fabric'. The main heading is 'Data Fabric'. The definition states: 'A data fabric is an emerging data management design for attaining flexible, reusable and augmented data integration pipelines, services and semantics. A data fabric supports both operational and analytics use cases delivered across multiple deployment and orchestration platforms and processes. Data fabrics support a combination of different data integration styles and leverage active metadata, knowledge graphs, semantics and ML to augment data integration design and delivery.'

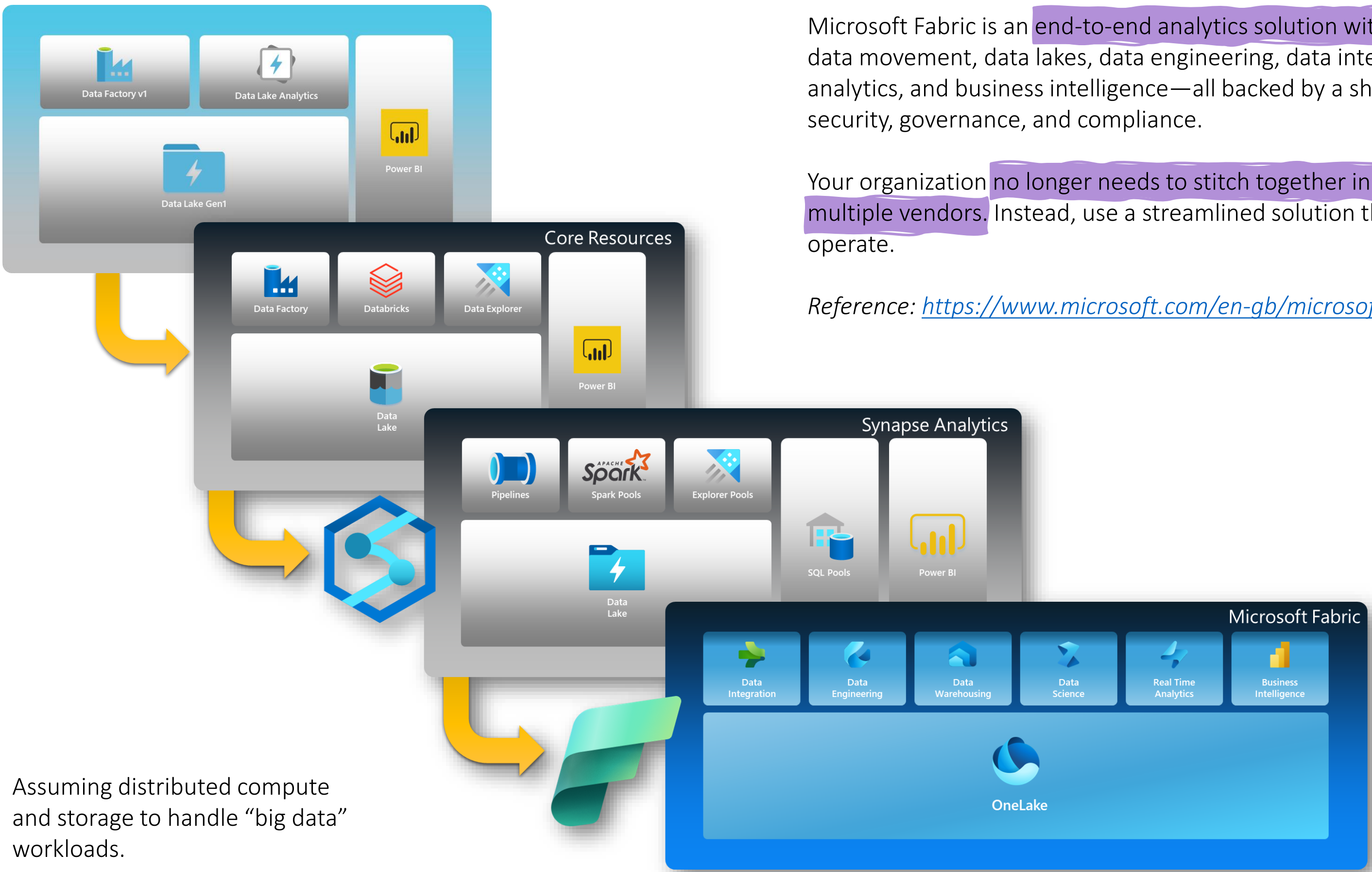
Ref: <https://www.gartner.com/en/information-technology/glossary/data-fabric>

The screenshot shows the Microsoft documentation page titled 'What is Microsoft Fabric?'. The breadcrumb trail is 'Learn / Microsoft Fabric / Get started /'. The article is dated 11/15/2023 and has 6 contributors. A 'Feedback' button is visible. The 'In this article' section lists: 'SaaS foundation', 'Components of Microsoft Fabric', 'OneLake and lakehouse - the unification of lakehouses', 'Fabric solutions for ISVs', and 'Next steps'. The main text describes Microsoft Fabric as an 'all-in-one analytics solution' for enterprises covering data movement, data science, Real-Time Analytics, and business intelligence. It offers a 'comprehensive suite of services' including data lake, data engineering, and data integration. The text notes that with Fabric, users don't need to piece together services from multiple vendors, but instead enjoy a highly integrated, end-to-end, and easy-to-use product. The platform is built on a foundation of Software as a Service (SaaS), which takes 'simplicity and integration to a whole new level.'

Ref: <https://www.gartner.com/en/information-technology/glossary/data-fabric>



What is Microsoft Fabric? – Vision and Stack Evolution



Microsoft Fabric is an **end-to-end analytics solution with full-service capabilities** including data movement, data lakes, data engineering, data integration, data science, real-time analytics, and business intelligence—all backed by a shared platform providing robust data security, governance, and compliance.

Your organization **no longer needs to stitch together individual analytics services from multiple vendors.** Instead, use a streamlined solution that's easy to connect, onboard, and operate.

Reference: <https://www.microsoft.com/en-gb/microsoft-fabric>

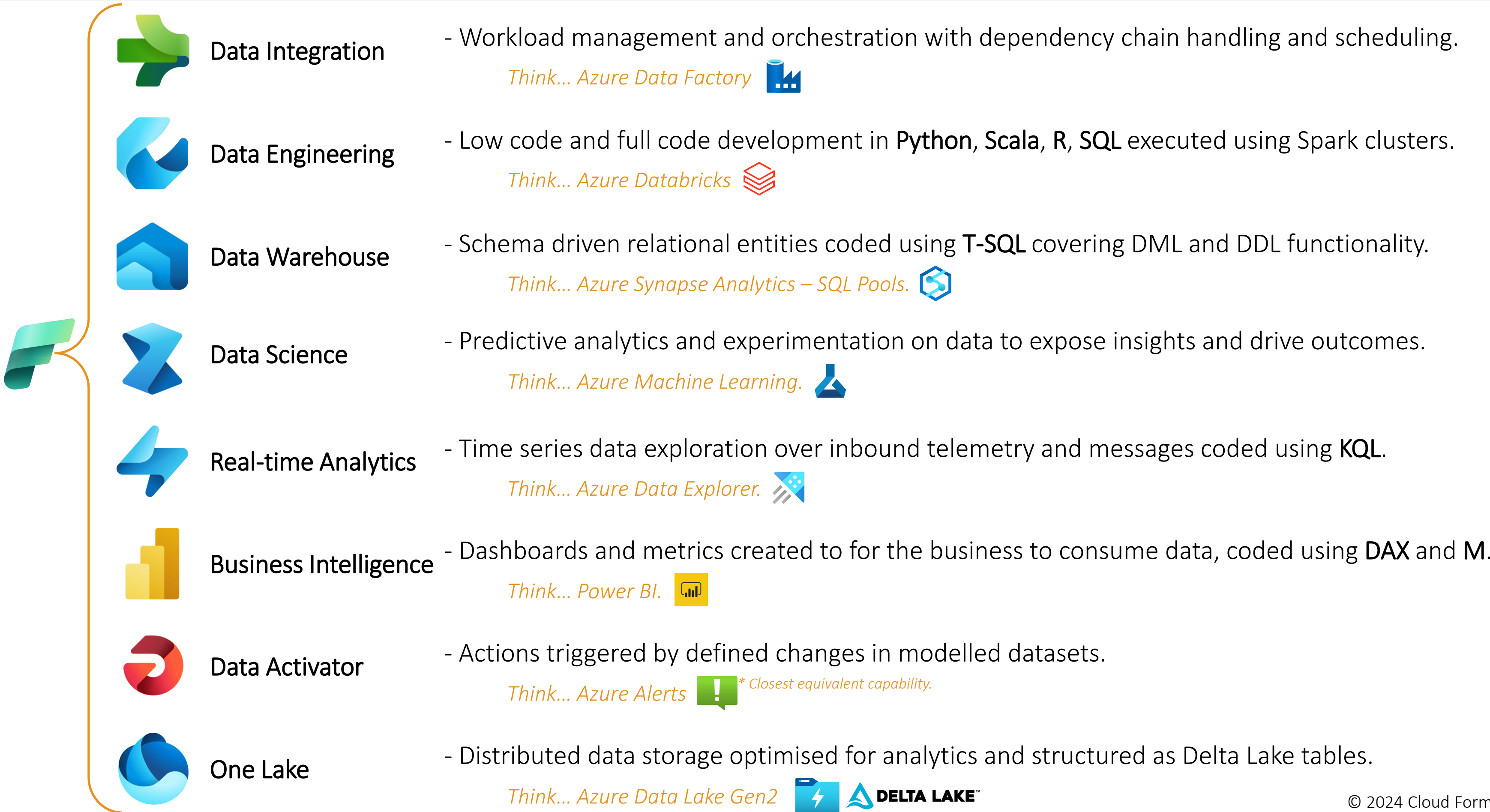
Cloud Formations - Knowledge Transfer & Training

Assuming distributed compute and storage to handle “big data” workloads.

What is Microsoft Fabric?

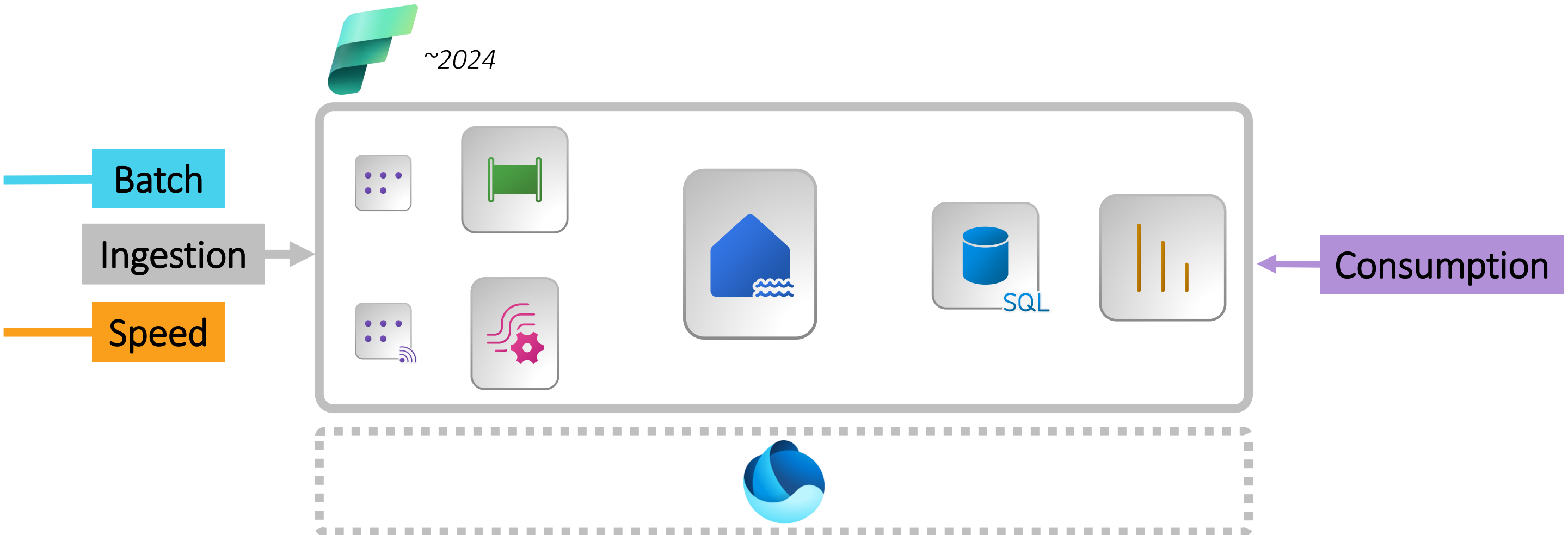


What is Microsoft Fabric? - Experiences vs Technical Capabilities





Microsoft Fabric vs a Kappa Architecture

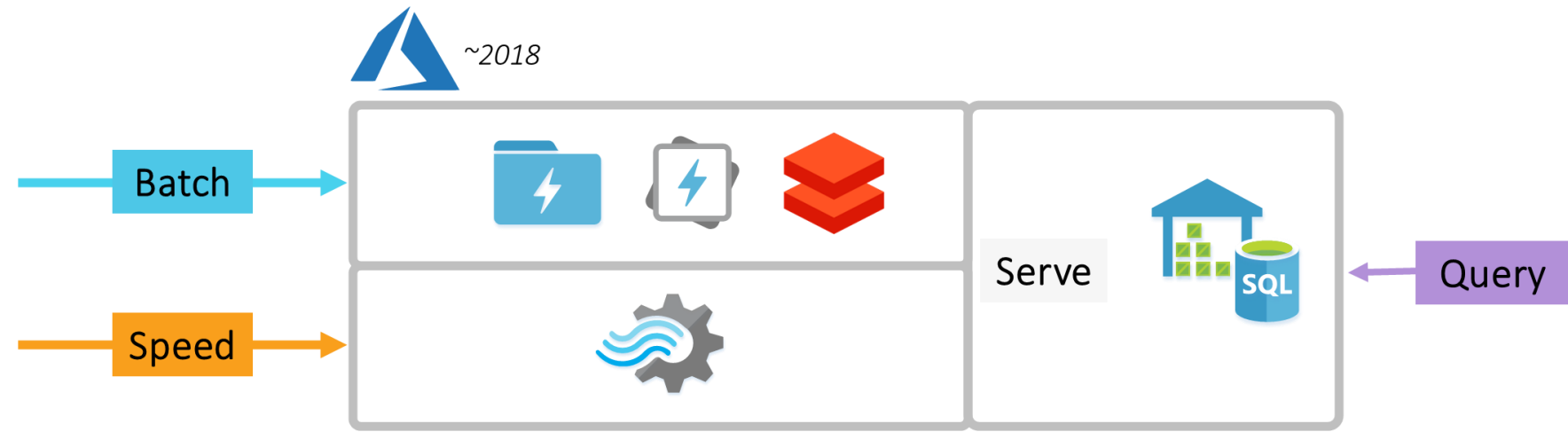
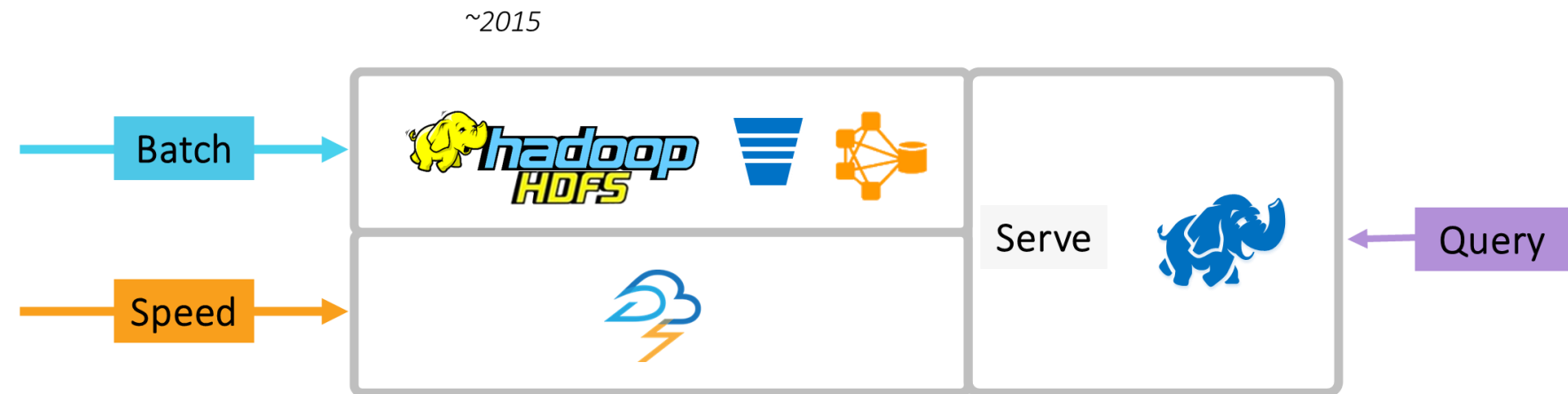


Lambda & Kappa Architectures vs Technology

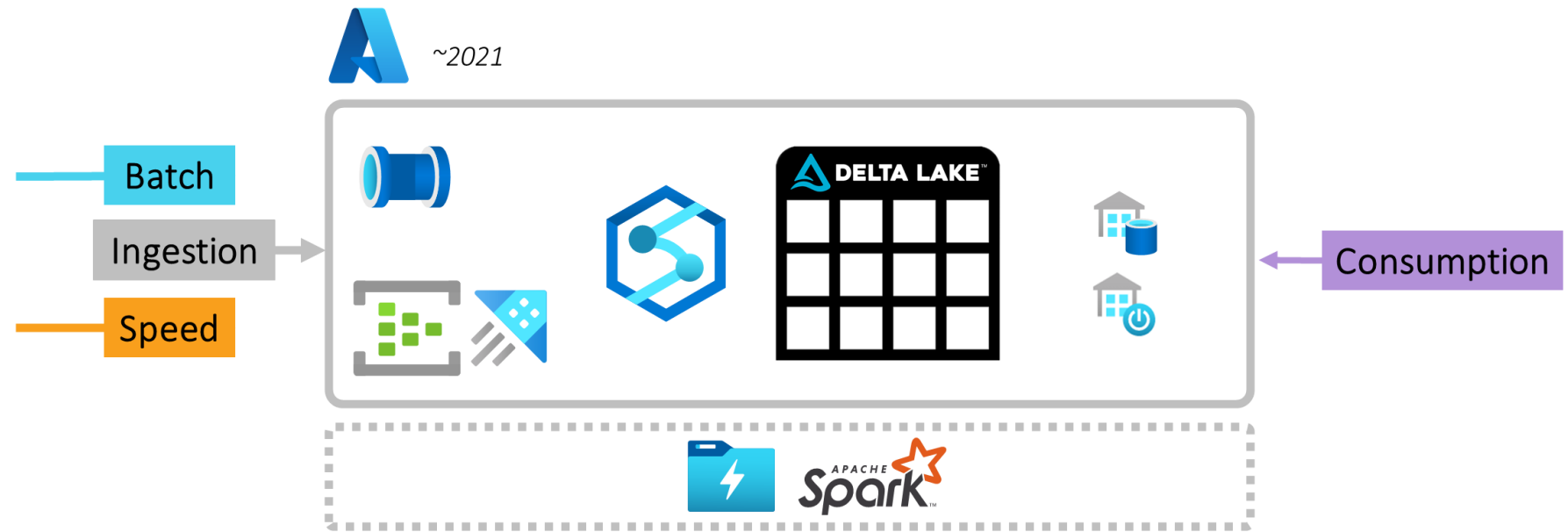
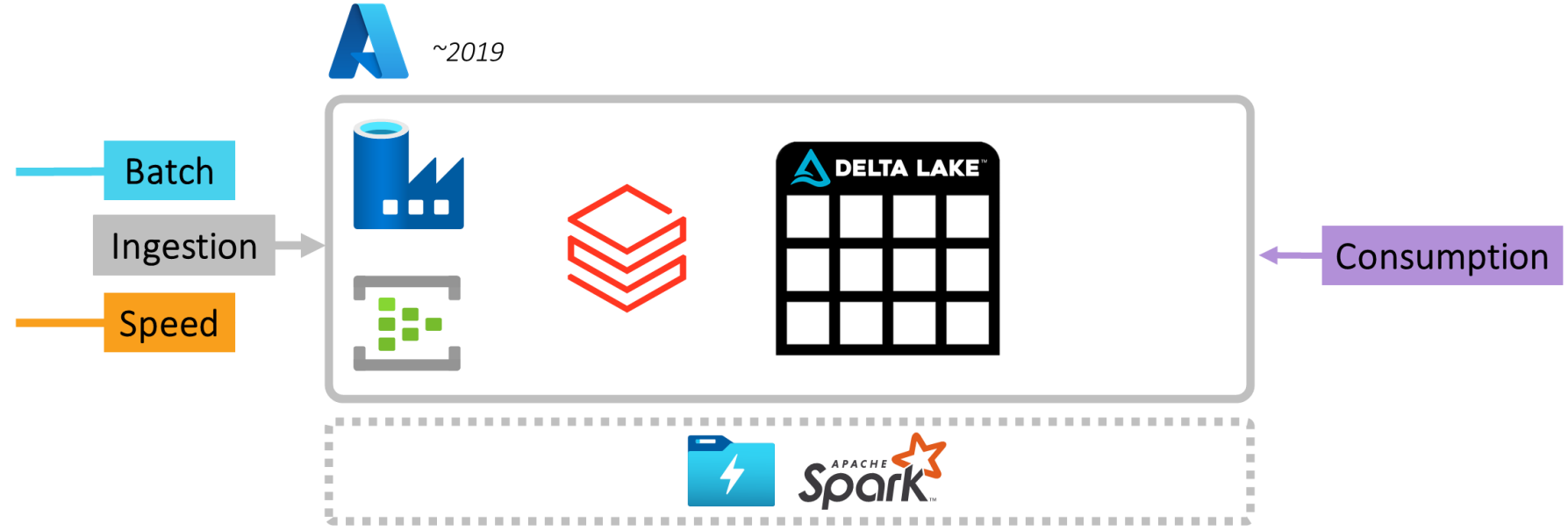


Cloud Formations - Knowledge Transfer & Training

Lambda



Kappa

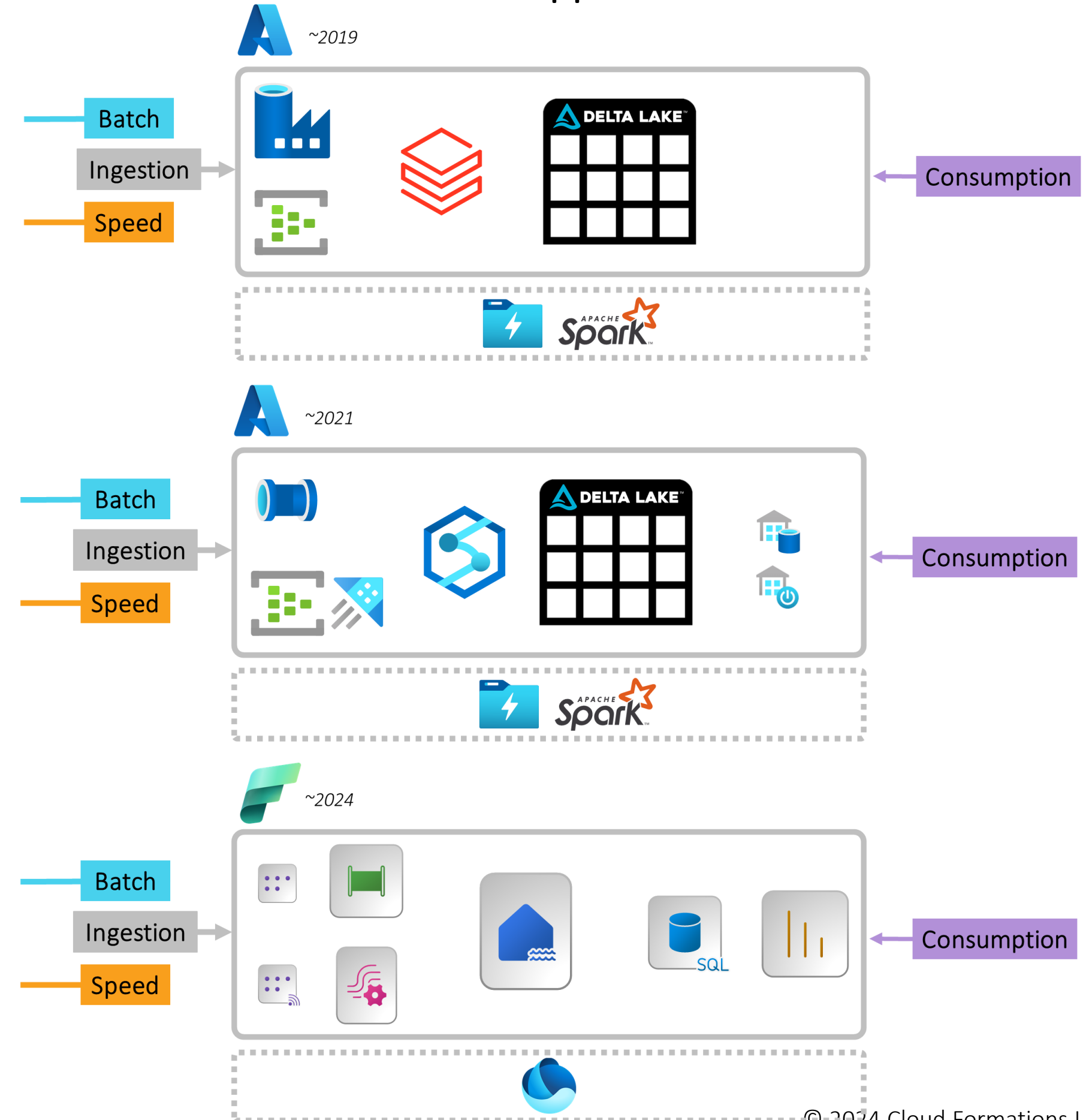
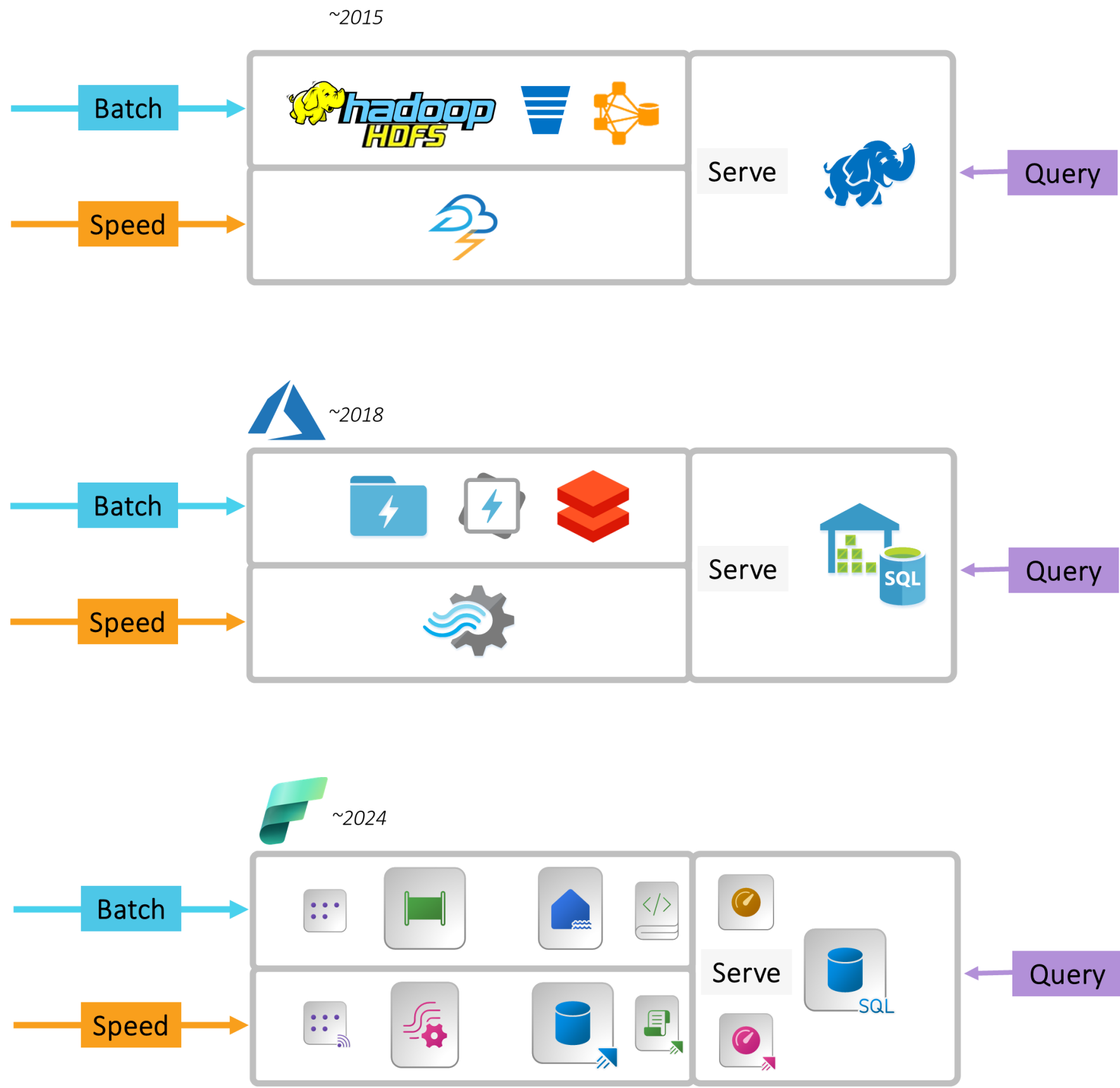


Lambda & Kappa Architectures vs Technology



Lambda

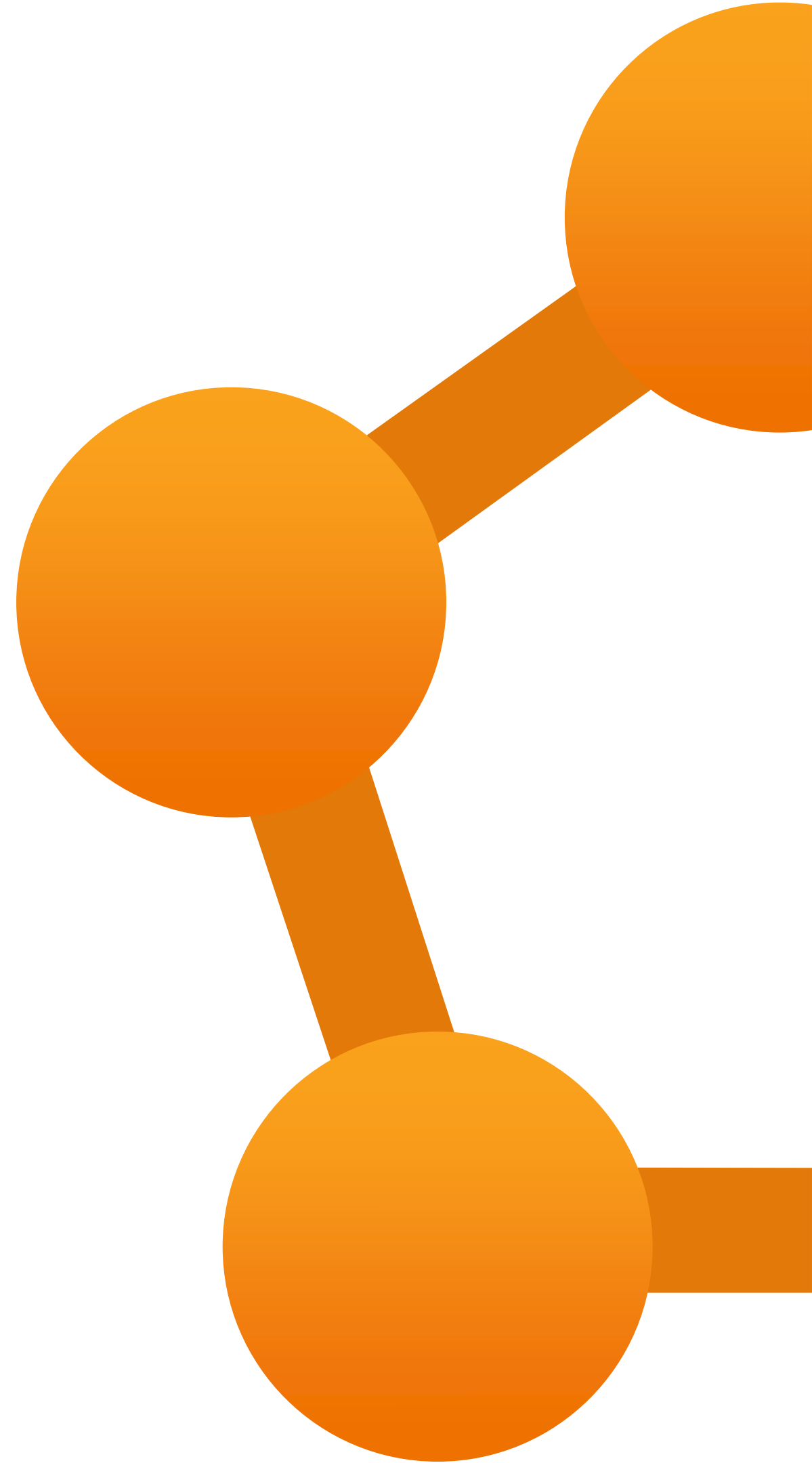
Kappa



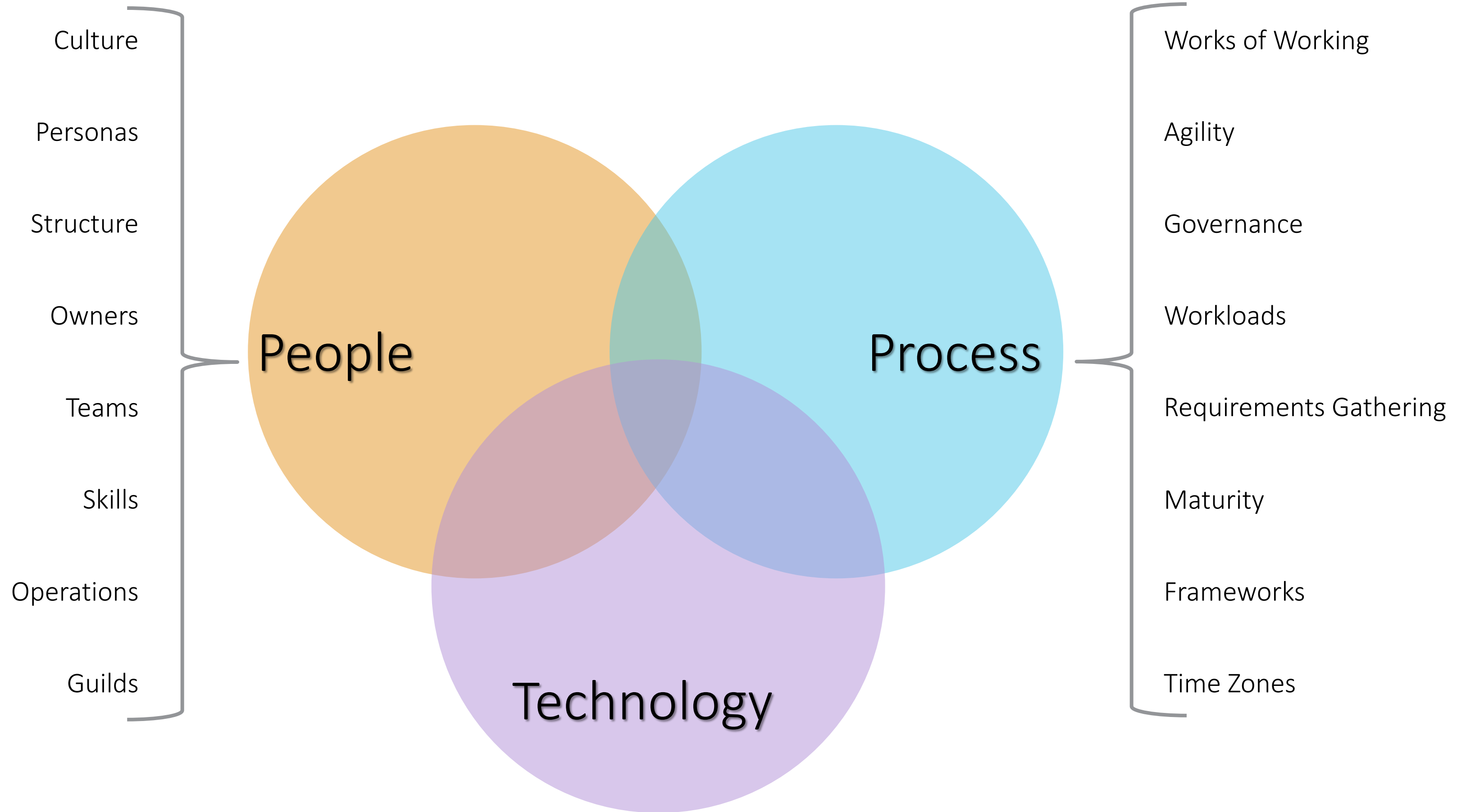
Mesh



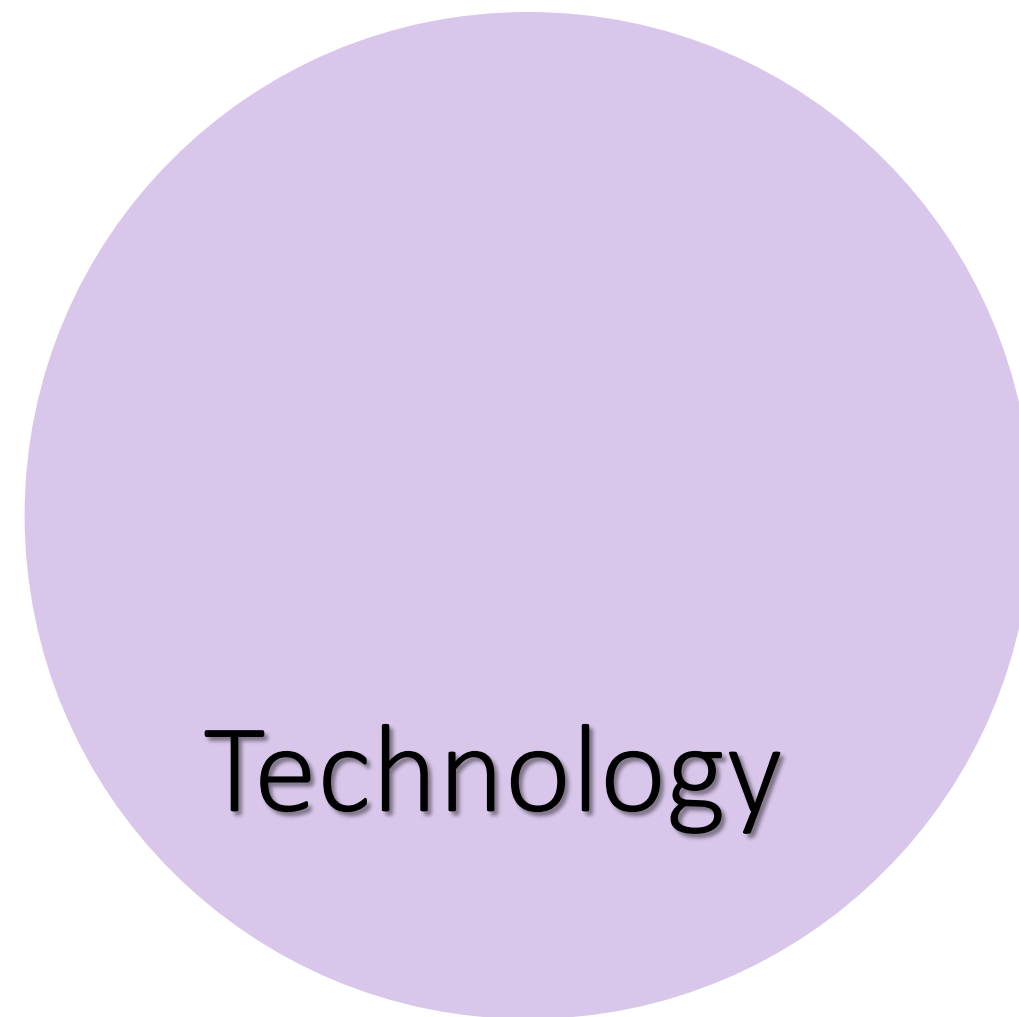
Cloud Formations



Data Mesh – What is it about?



Data Mesh – What is it about?



Data Mesh – What is it about?



Zhamak Dehghani

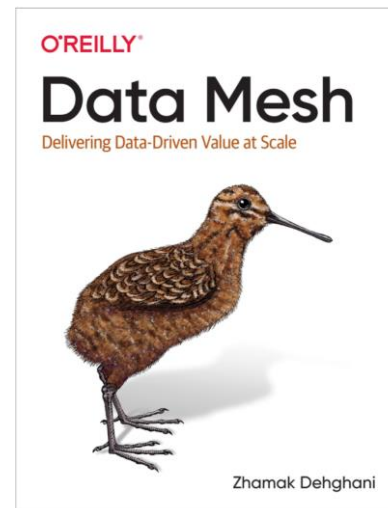
@zhamakd



<https://martinfowler.com/articles/data-mesh-principles.html>

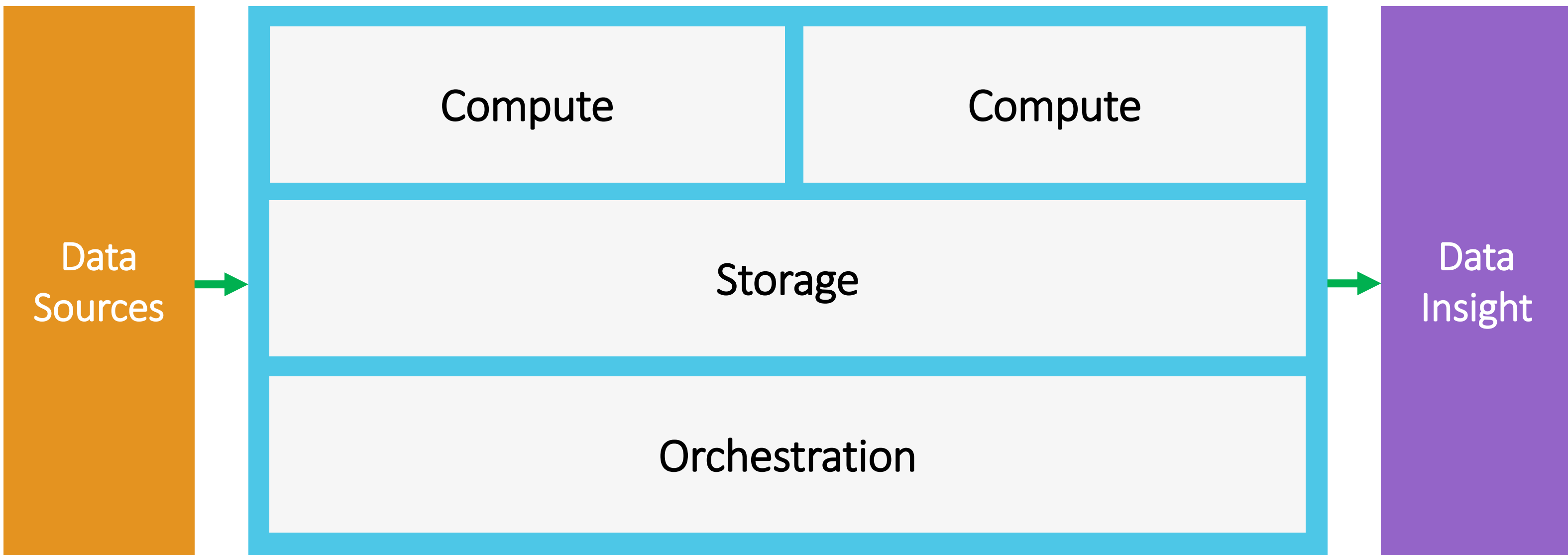
ISBN-10
1492092398

ISBN-13
978-1492092391



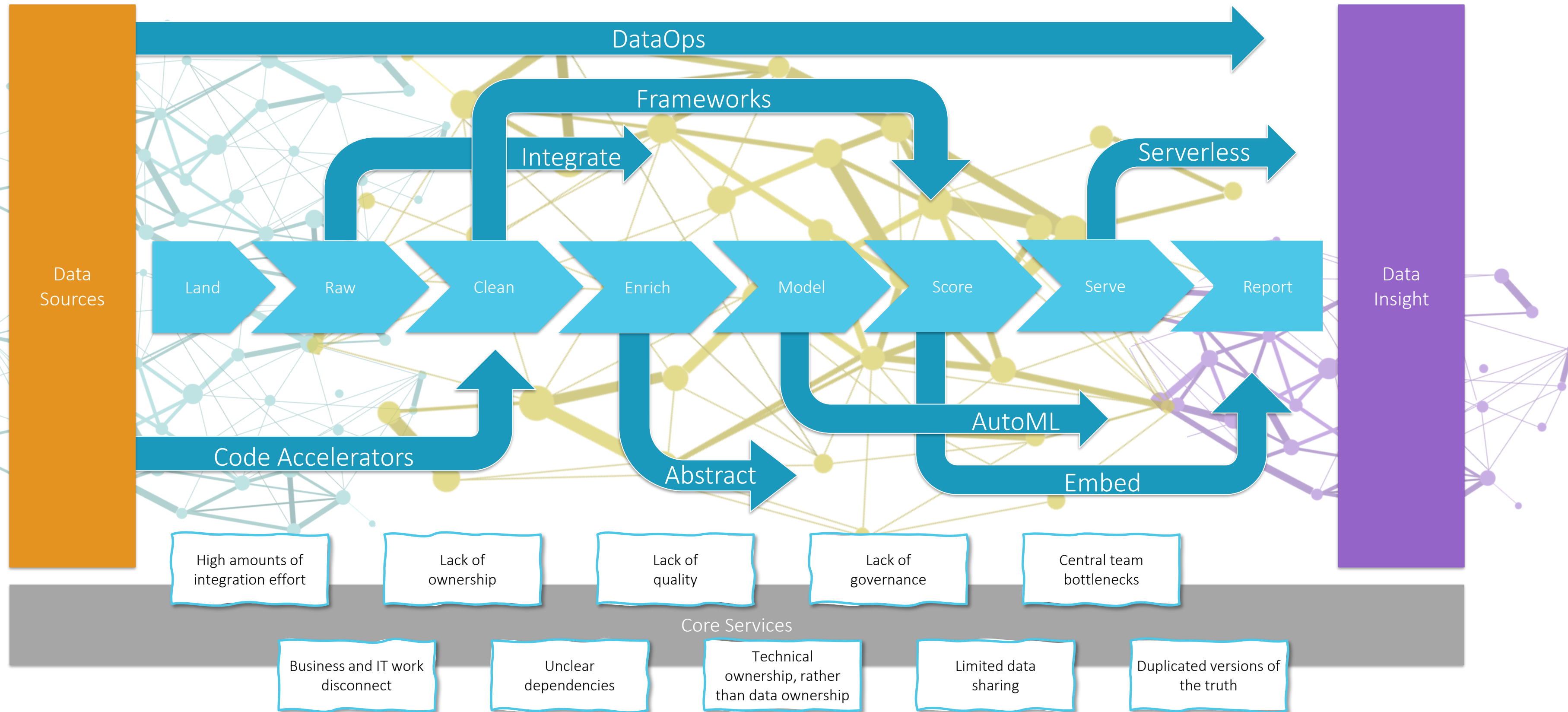
1. Domain-oriented decentralised data ownership and architecture.
2. Data as a product.
3. Self-serve data infrastructure as a platform.
4. Federated computational governance.

My First Reference Architecture

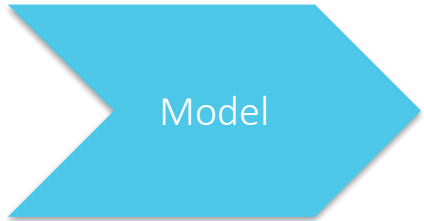


Data Mesh – Why should we build it?

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

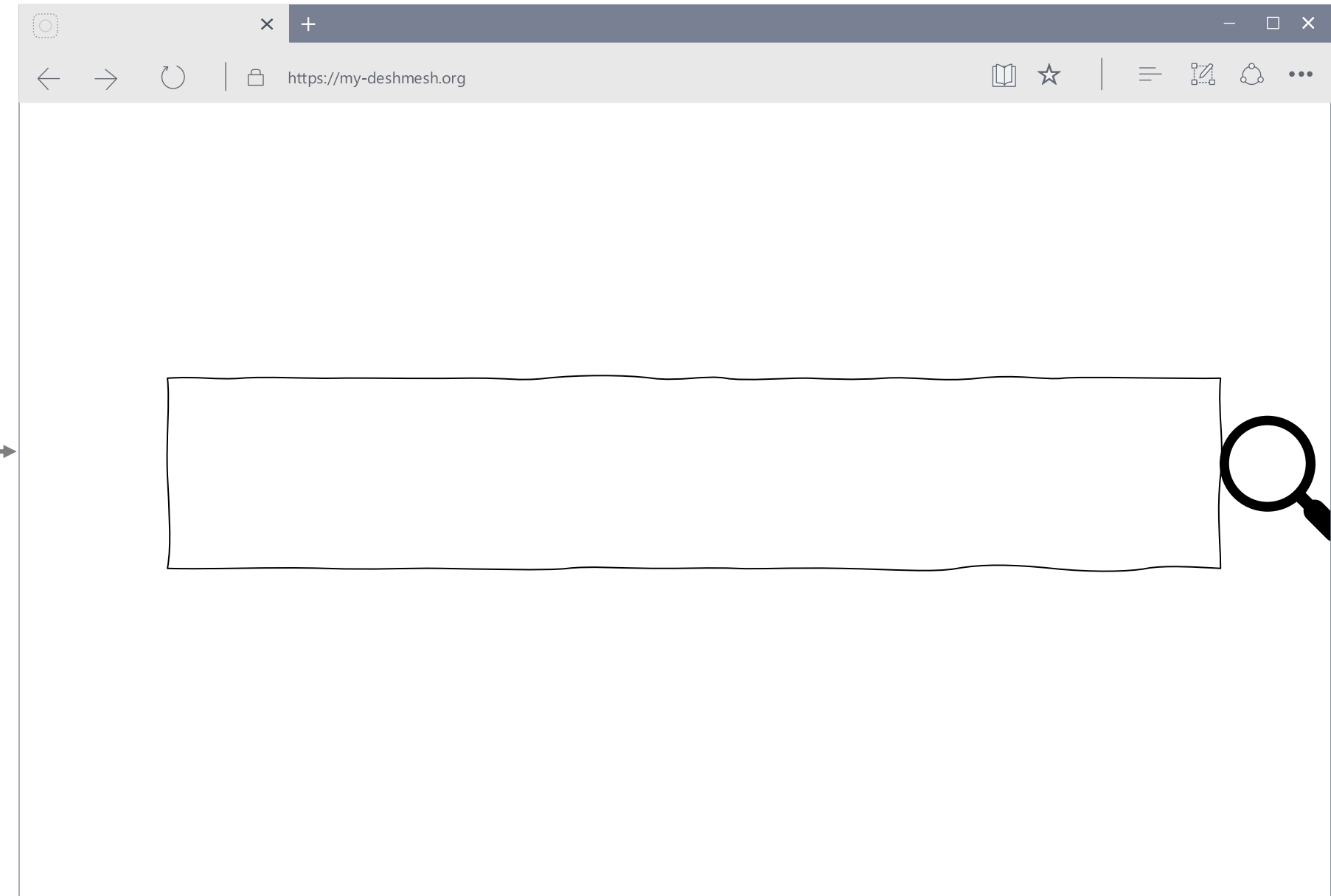
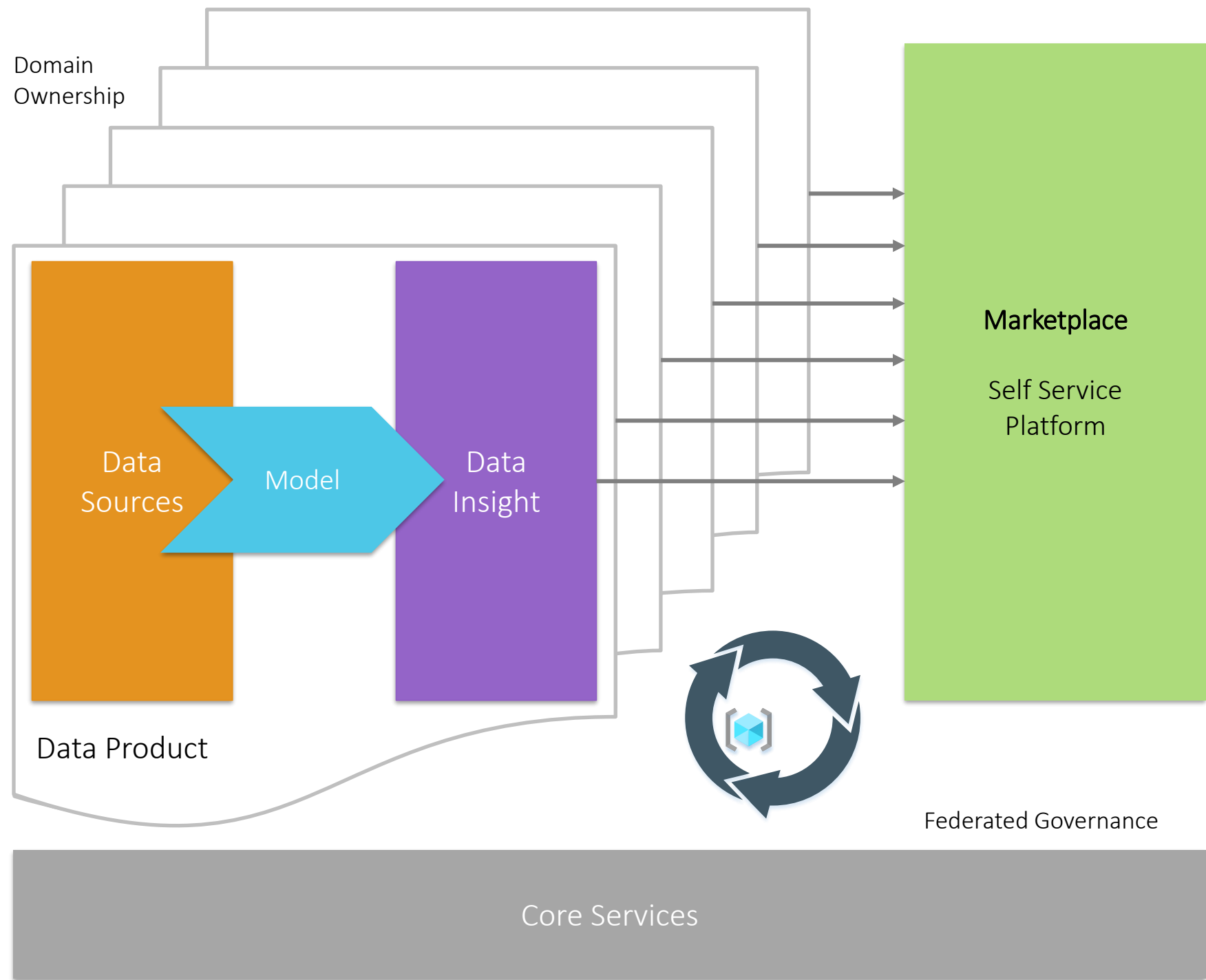


Data Mesh – Why should we build it?



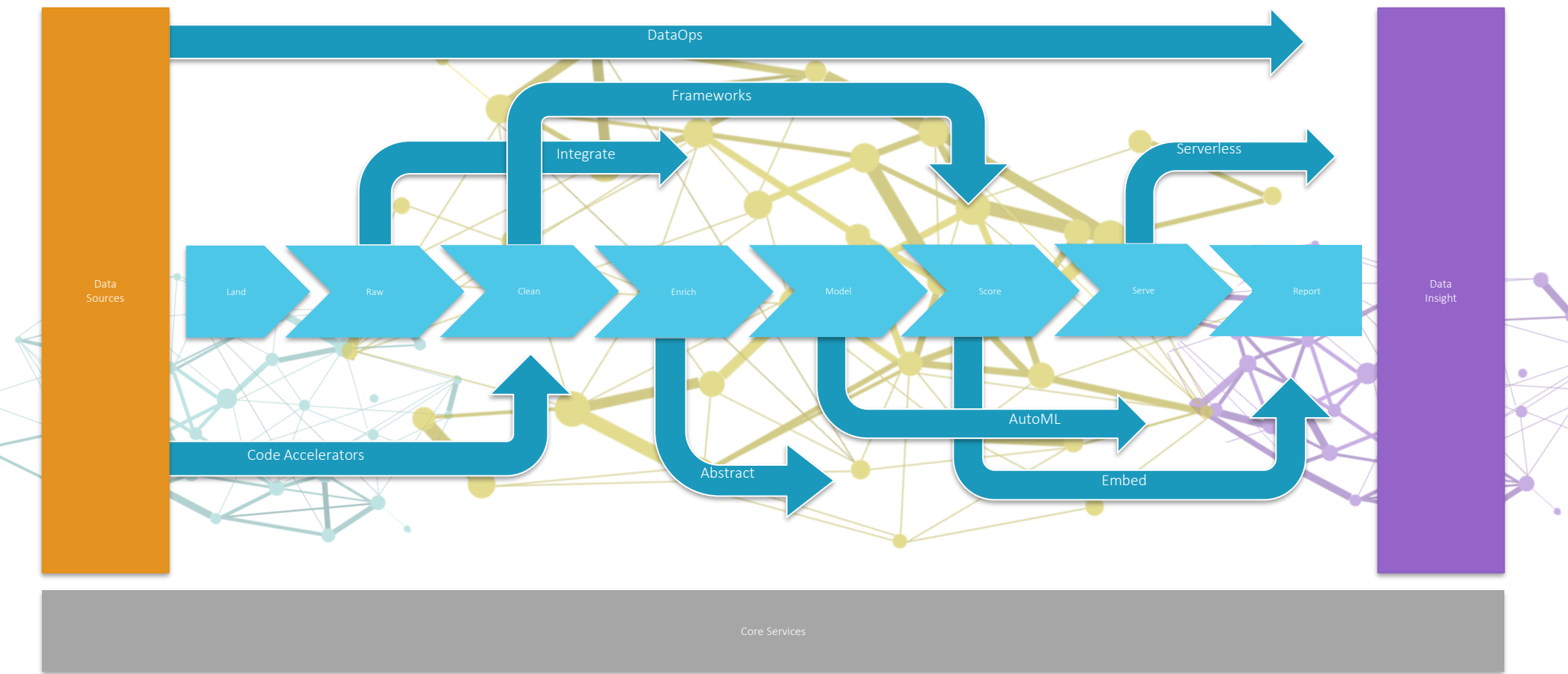
Data Mesh – Why should we build it?

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



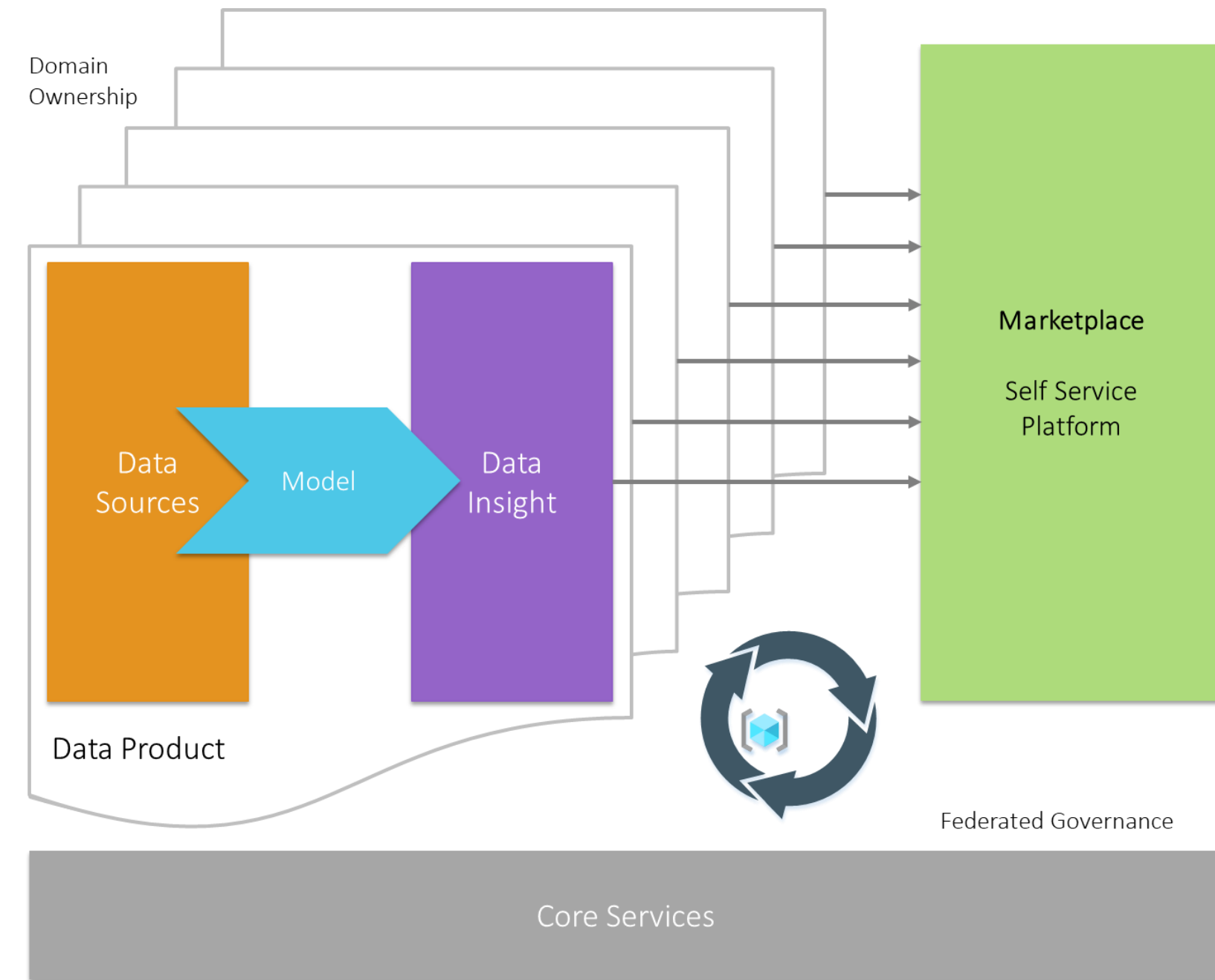
Data Mesh – Why should we build it? A: Time to Insight

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



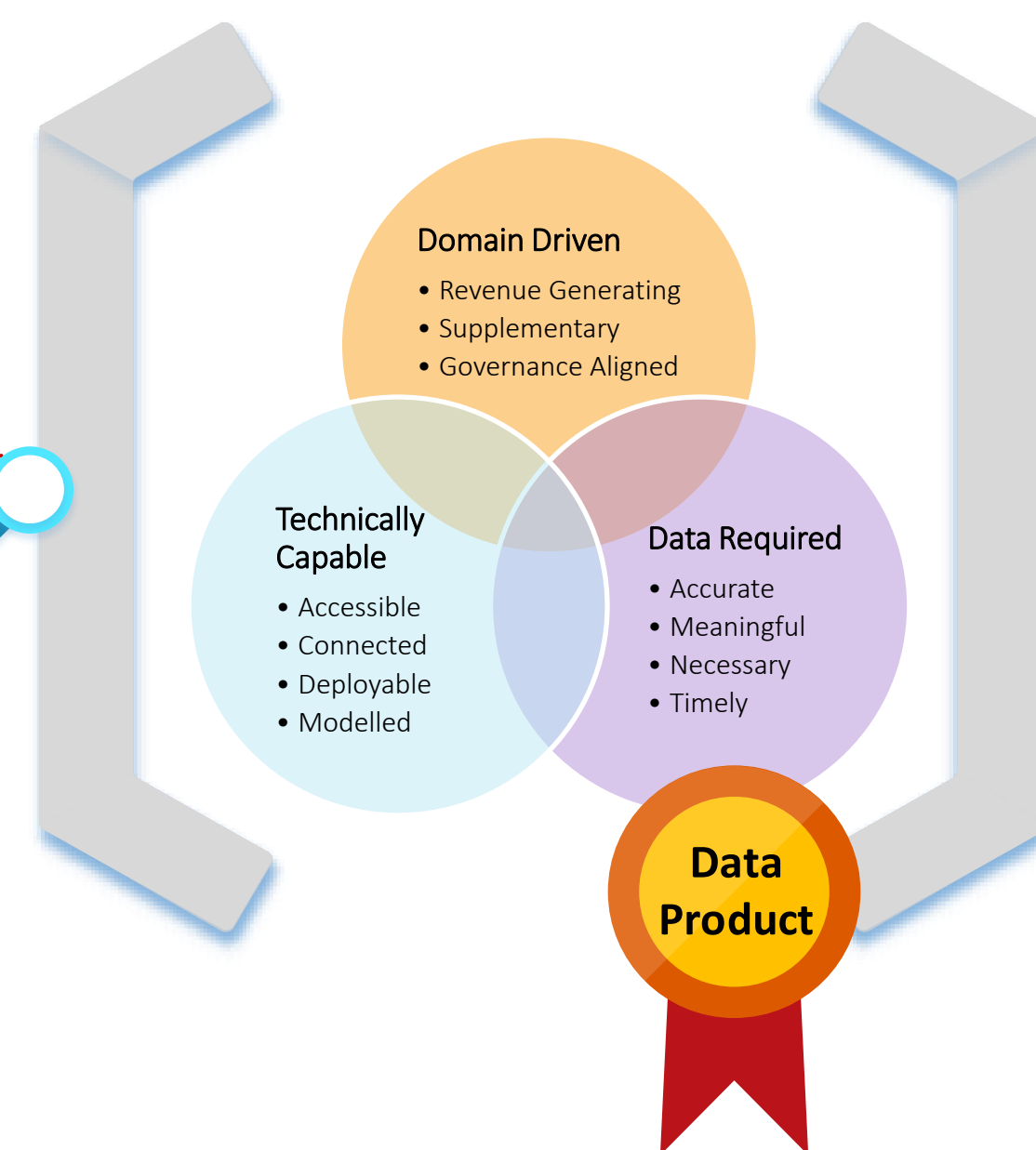
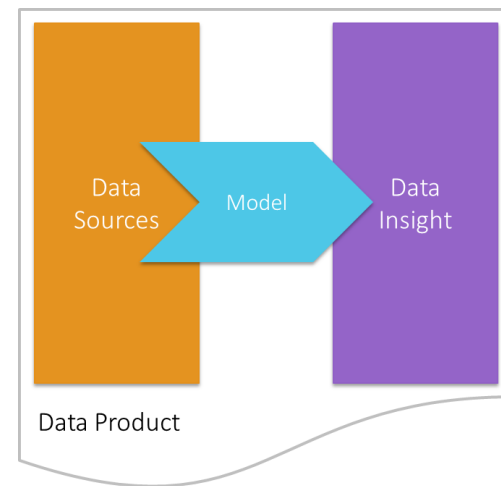
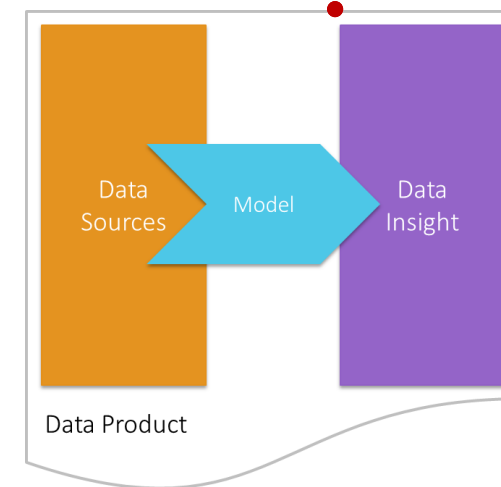
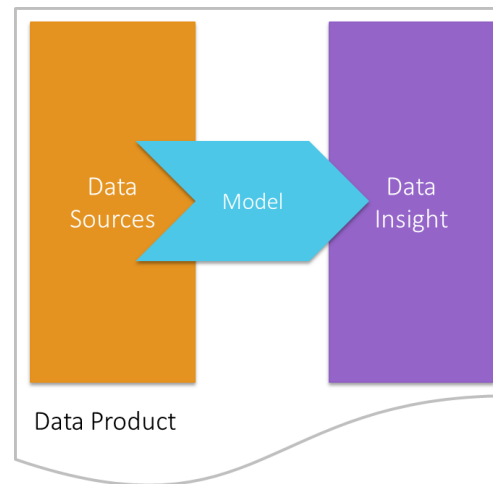
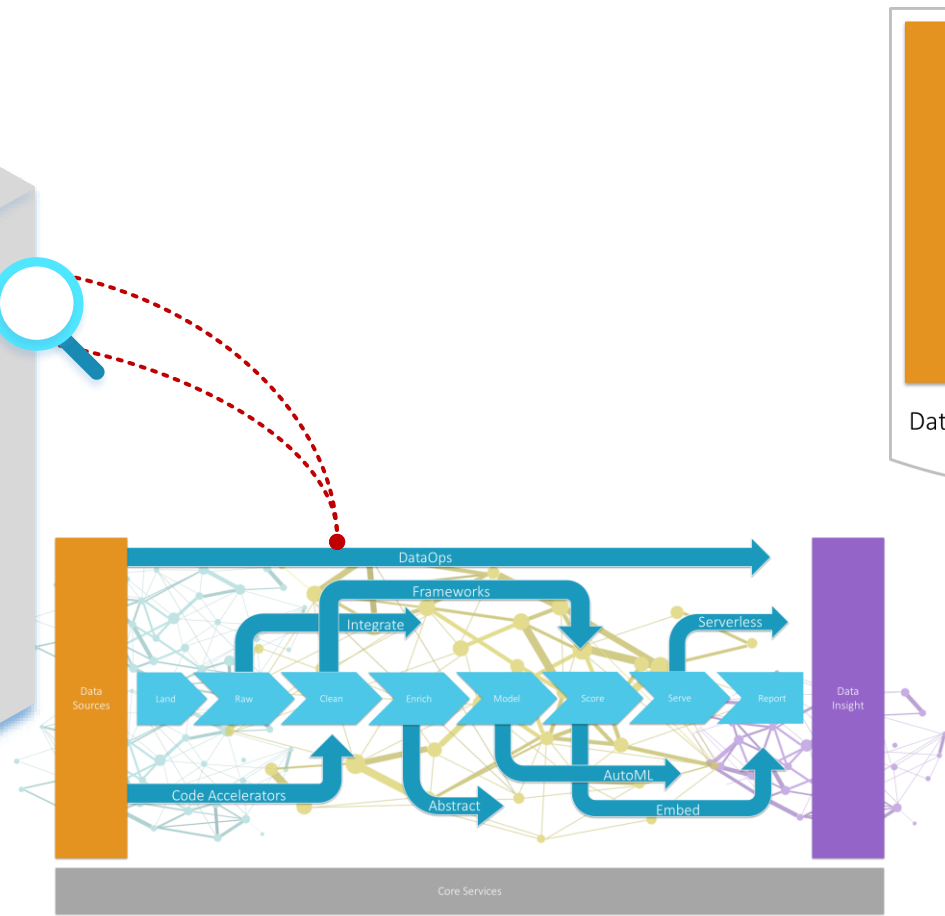
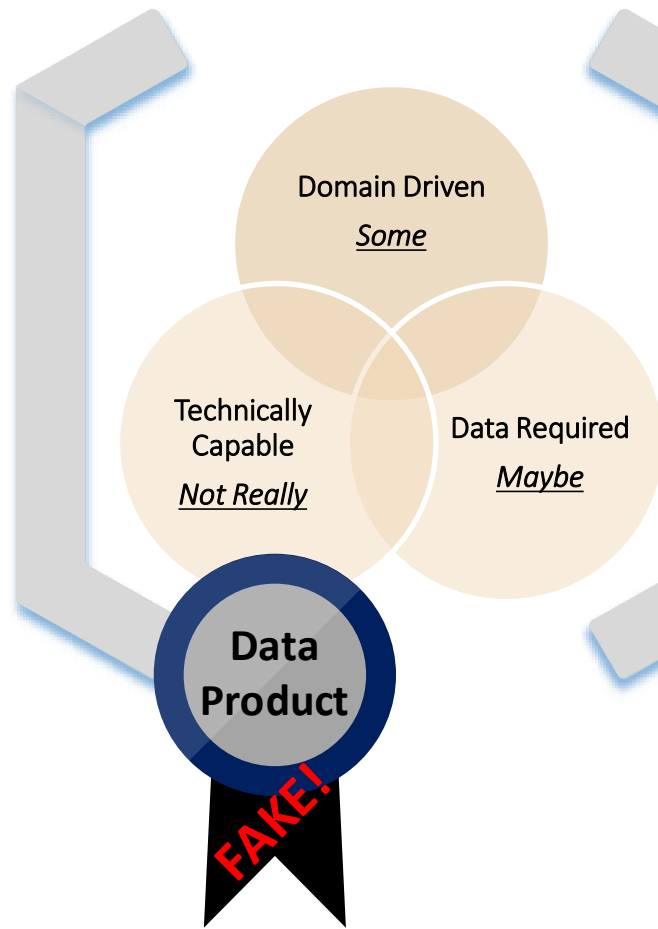
... Weeks/Months

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

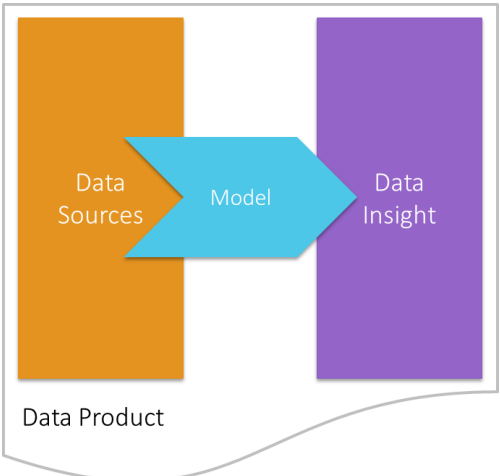
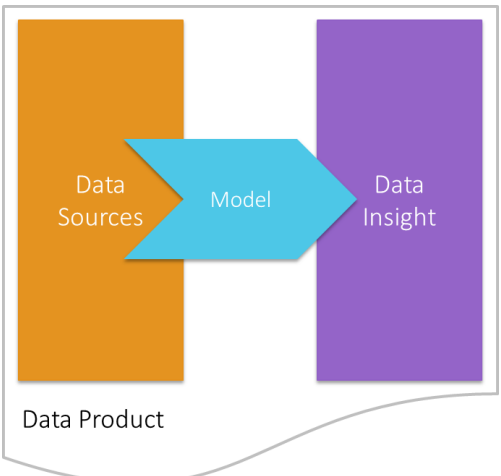
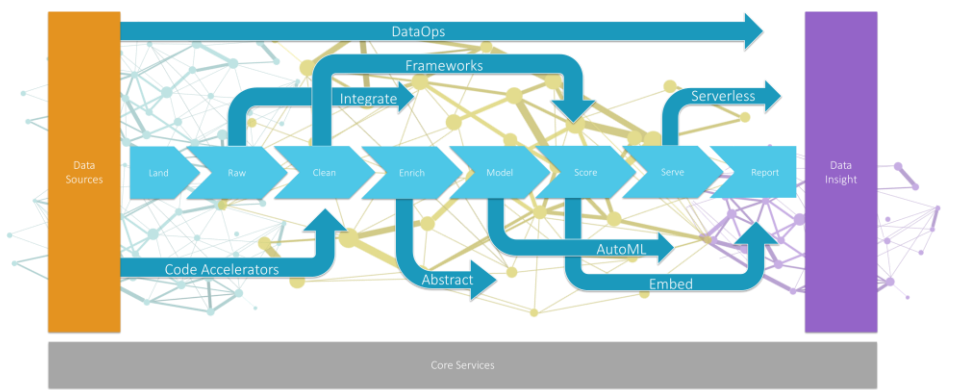
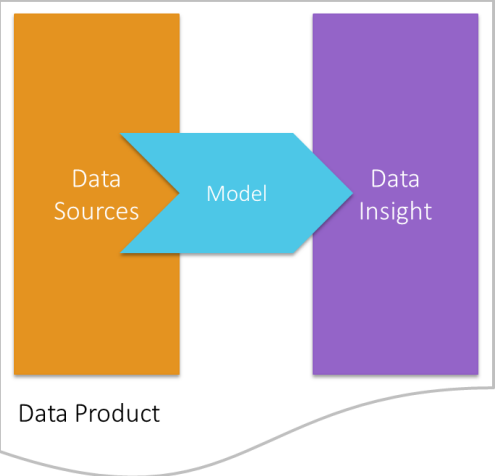


... Hours/Days

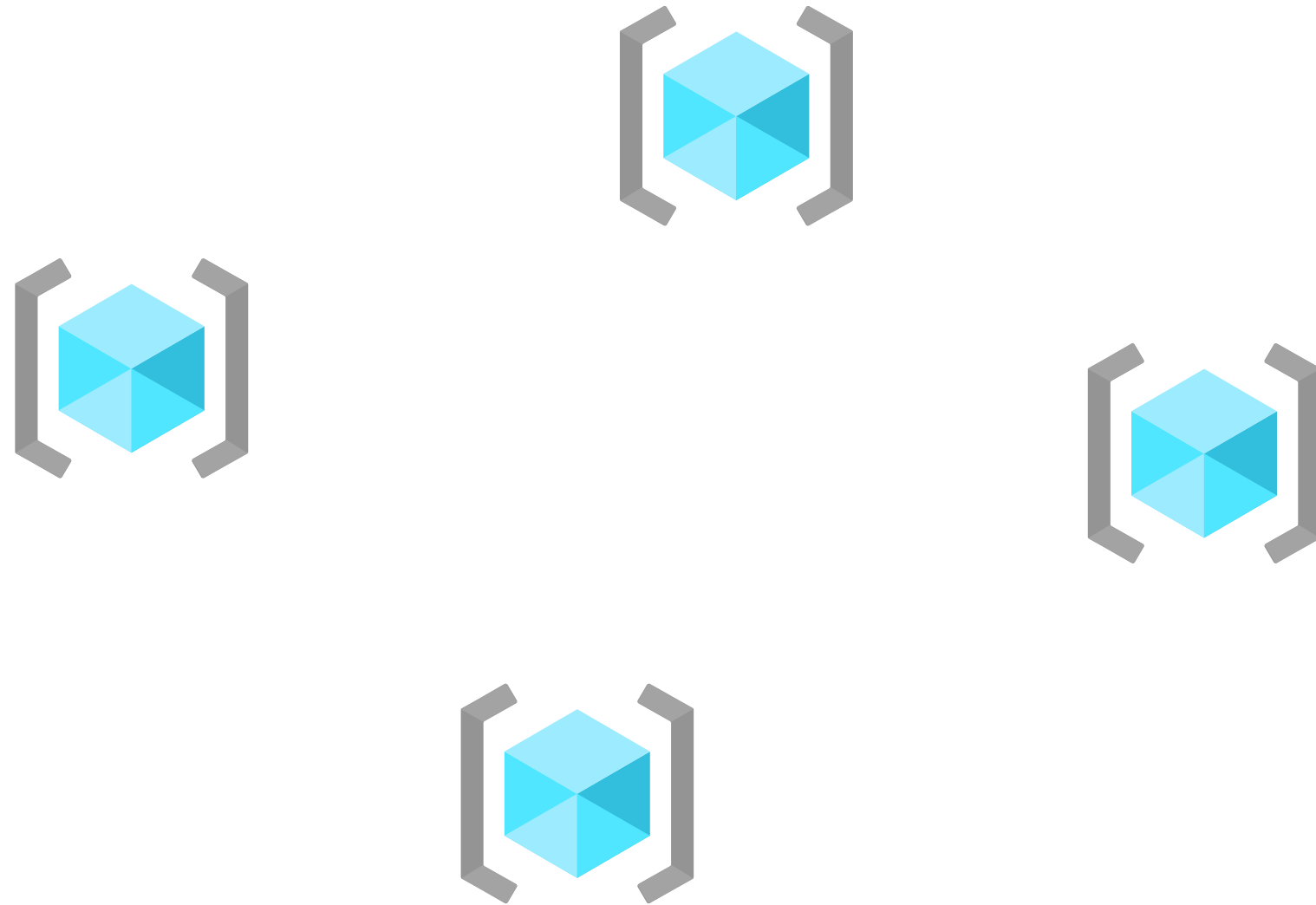
Data Mesh – Data Products



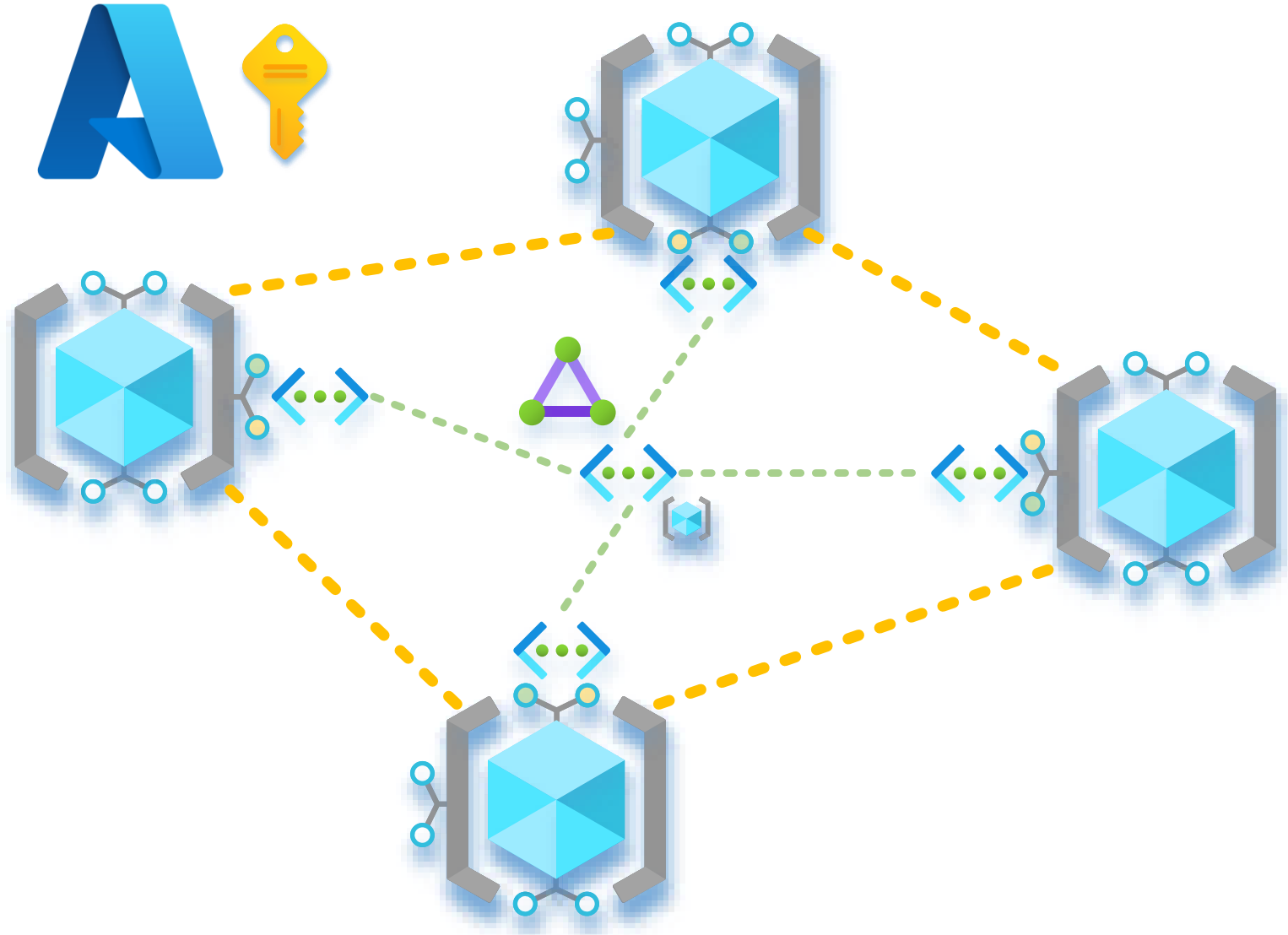
Data Mesh – Data Products in Azure



Data Mesh – Data Products in Azure



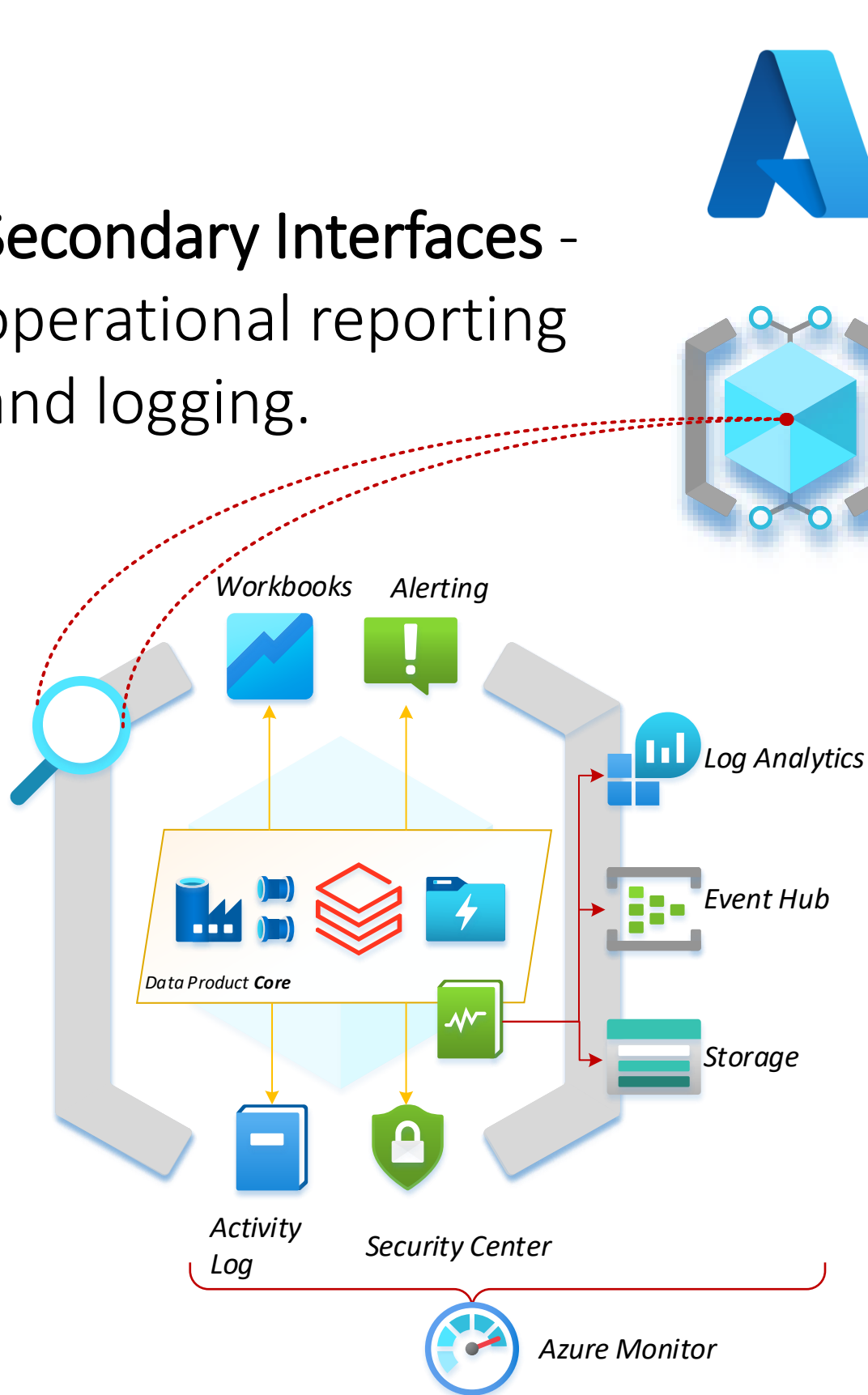
Data Mesh – Data Products in Azure with Interfaces



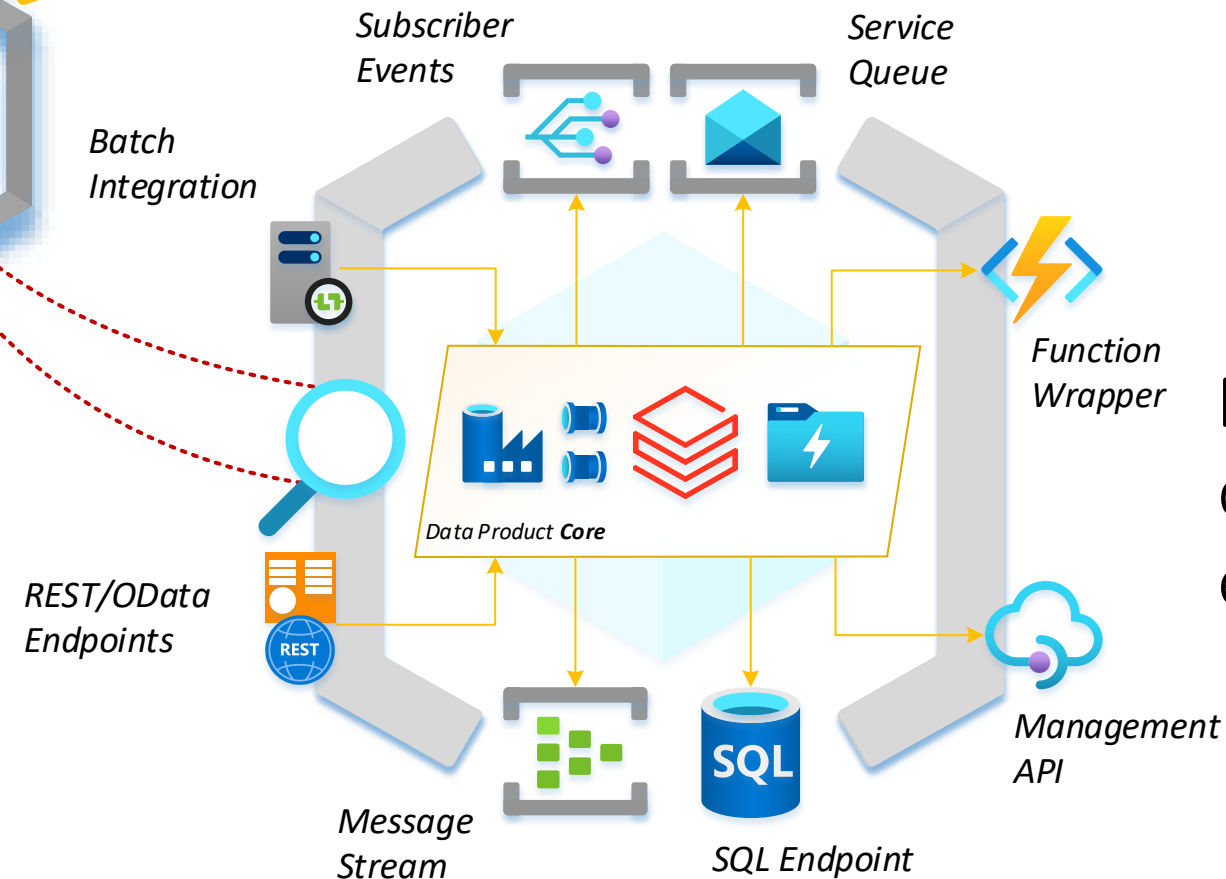
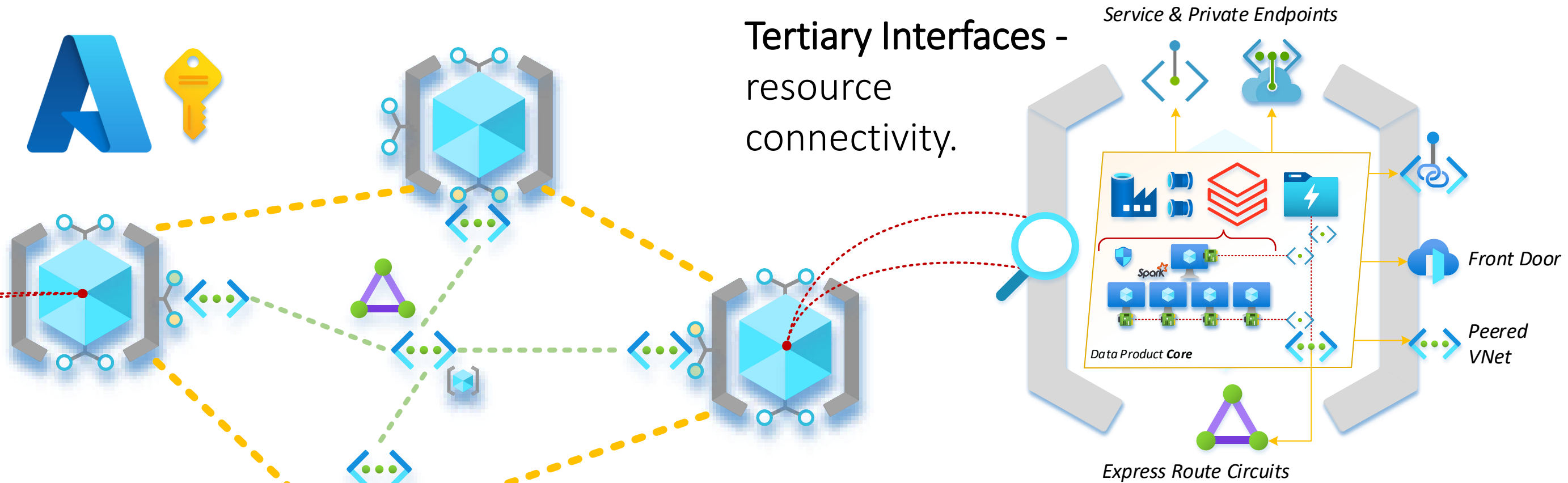
Data Mesh – Data Products in Azure with Interfaces



Secondary Interfaces -
operational reporting
and logging.

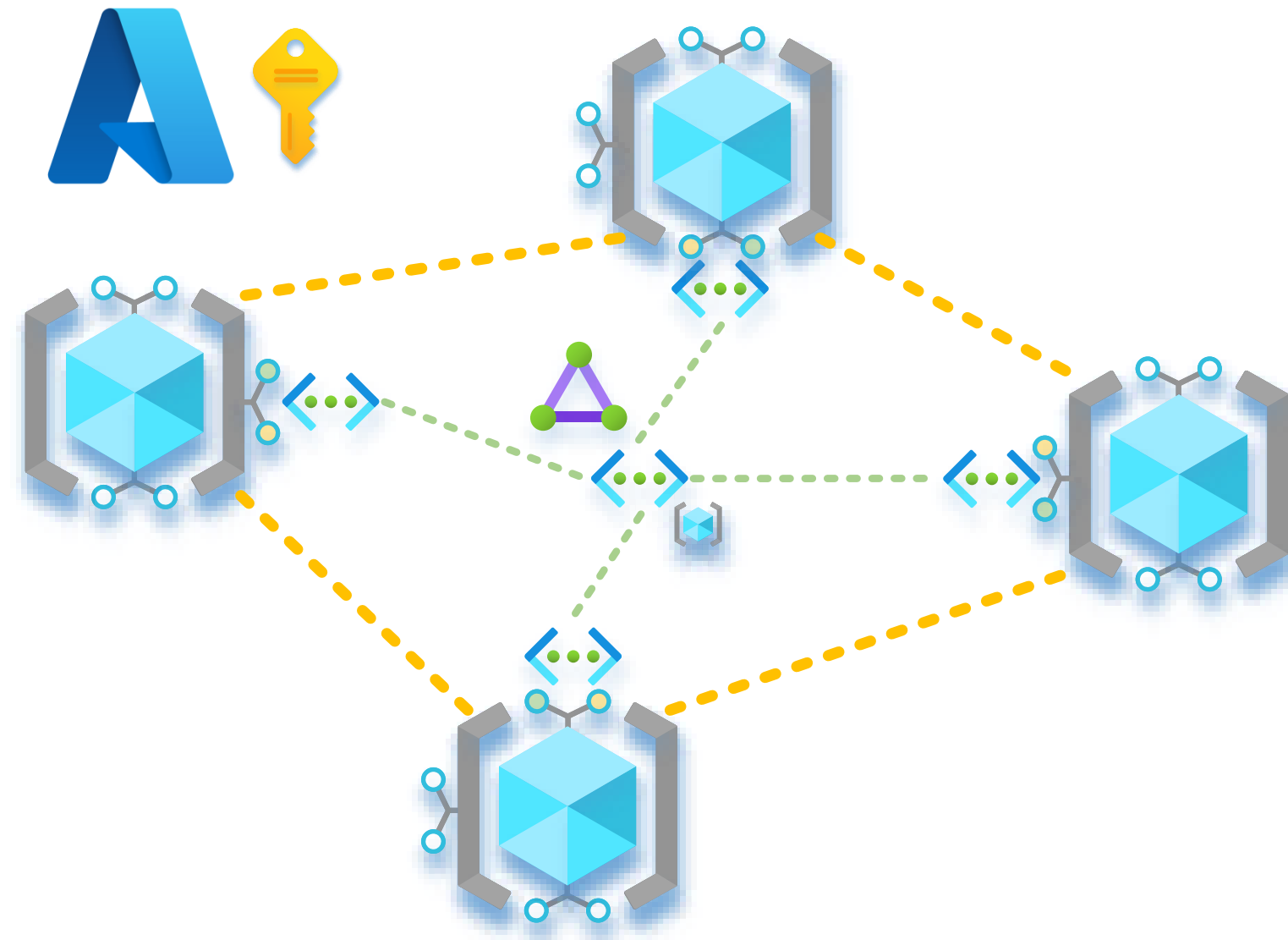


Tertiary Interfaces -
resource
connectivity.

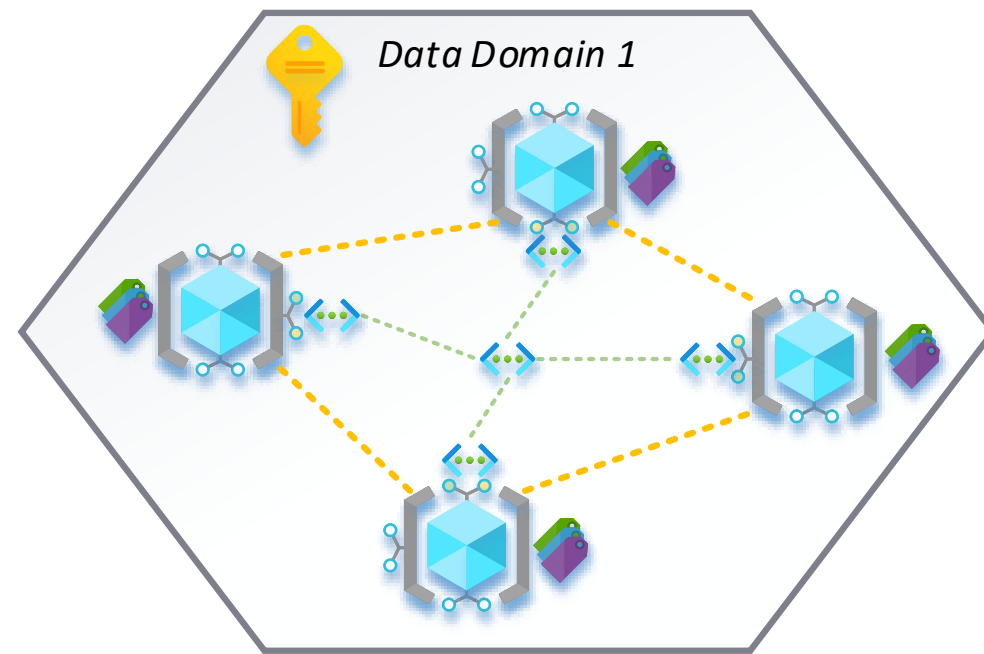


Primary Interfaces –
data integration and
exchange.

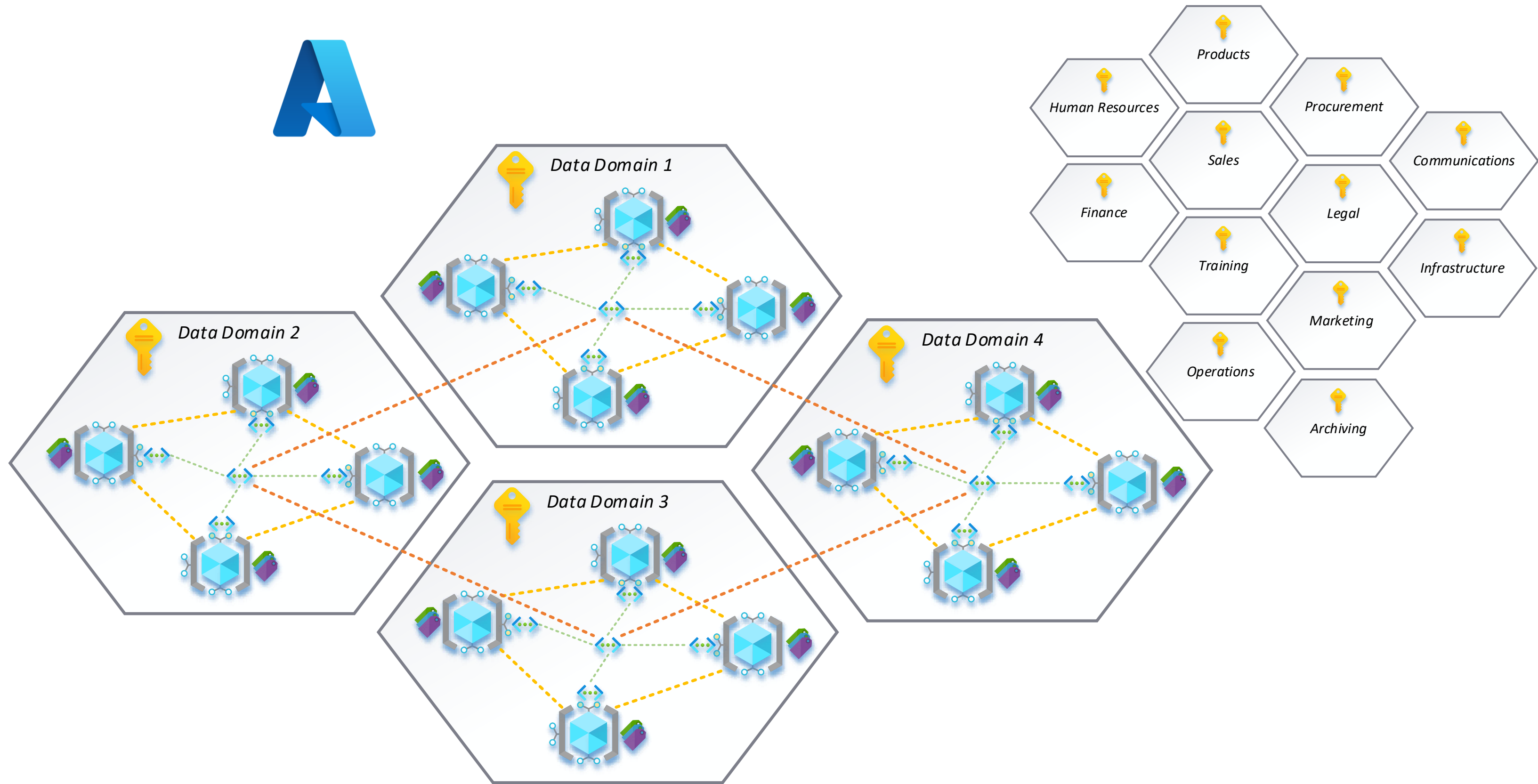
Data Mesh – Data Domains in Azure



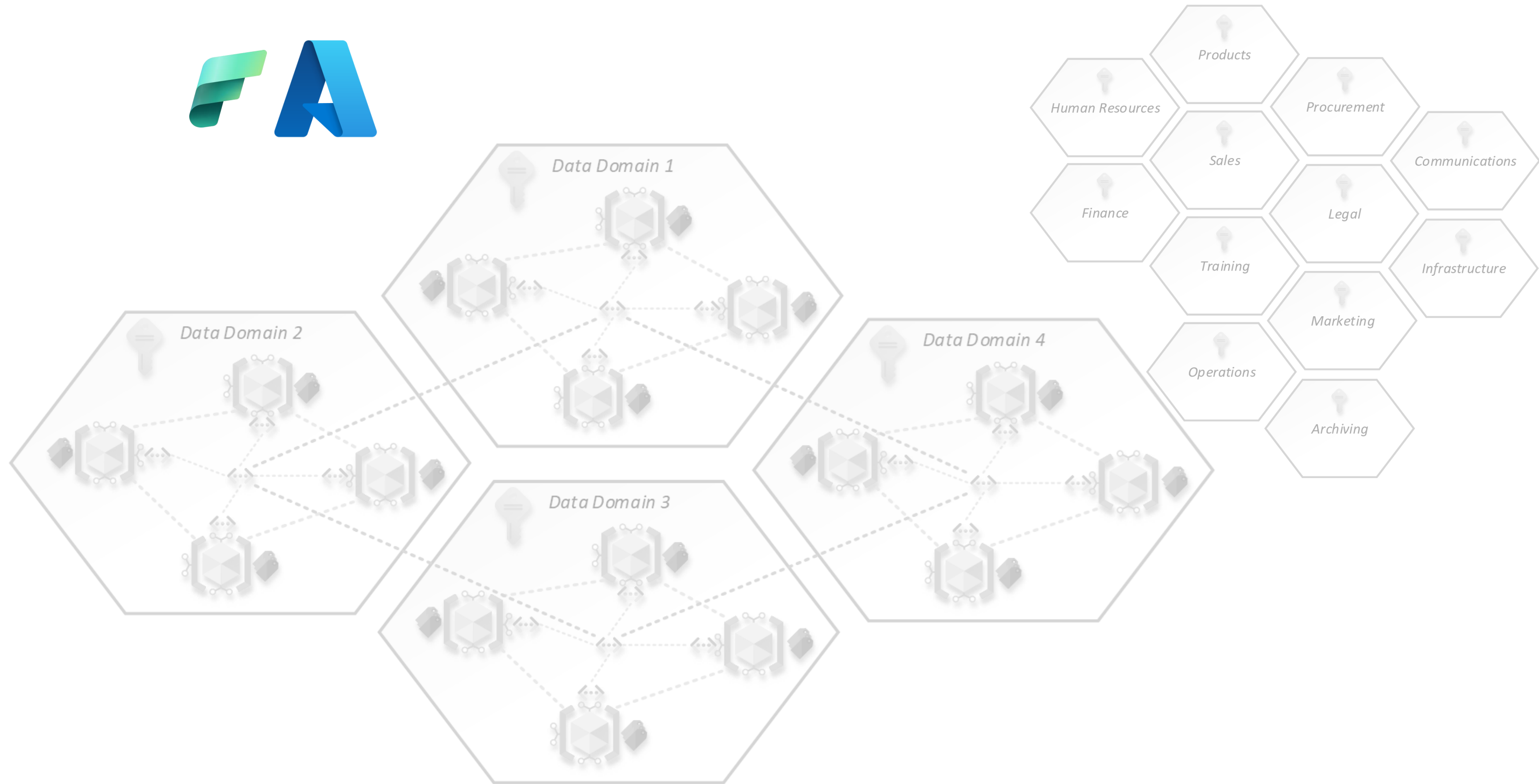
Data Mesh – Data Domains in Azure



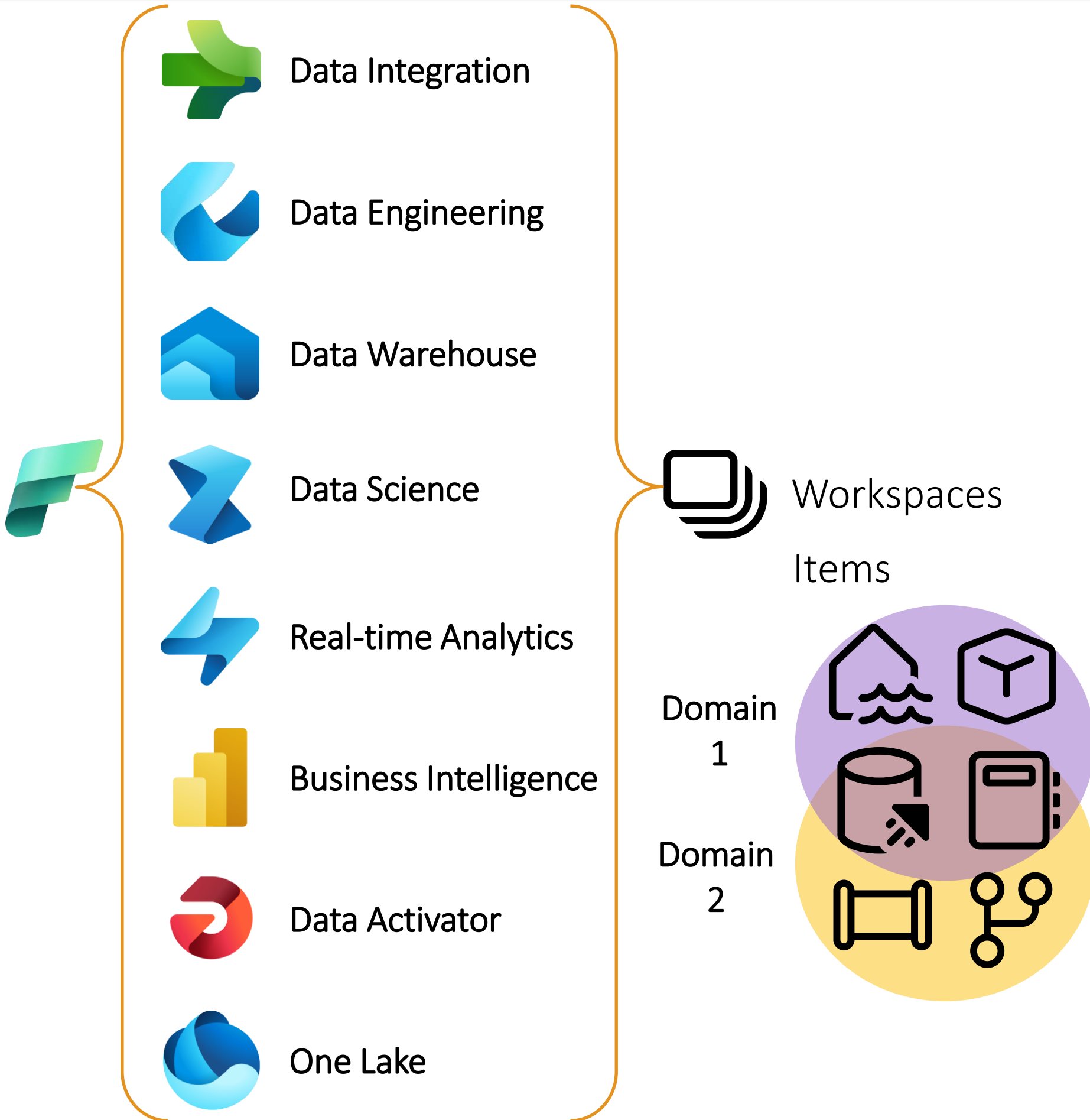
Data Mesh – Data Domains in Azure



Data Mesh – Data Domains in Fabric



Data Mesh – Data Domains in Fabric

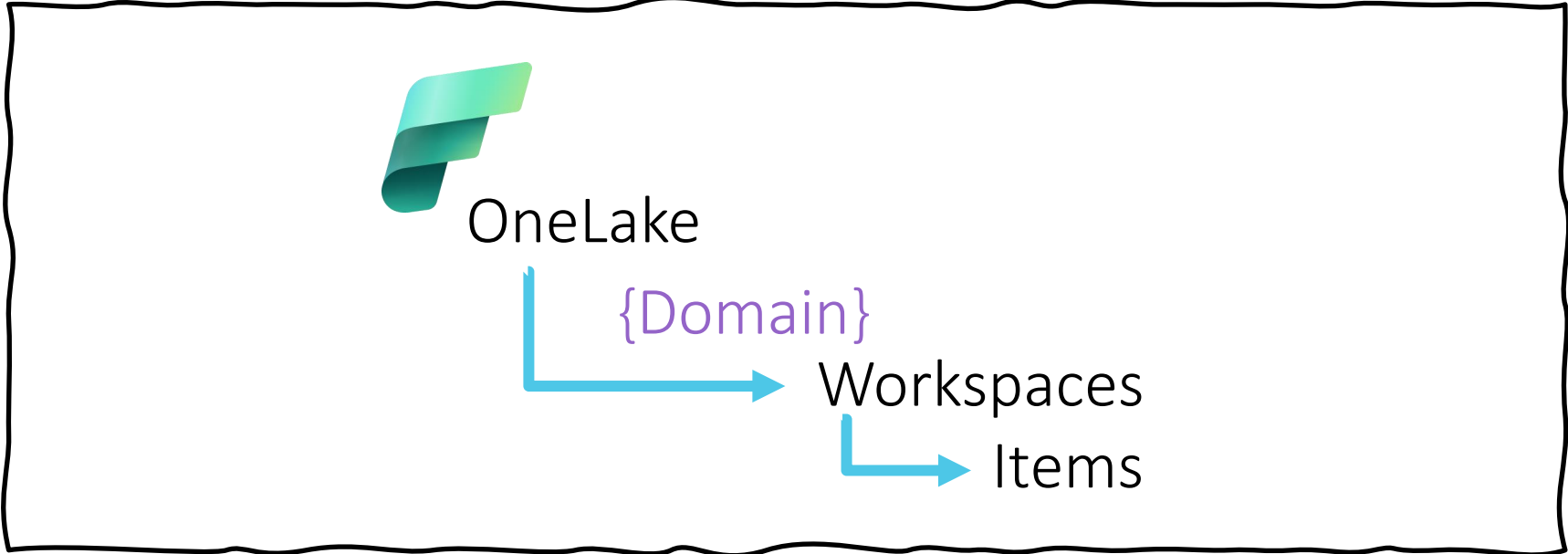


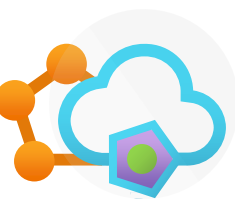
.... To meet this challenge, organizations are shifting from traditional IT centric data architectures, where the data is governed and managed centrally, to more federated models organized according to business needs. This federated data architecture is called data mesh. A data mesh is a decentralized data architecture that organizes data by specific business domains, such as marketing, sales, human resources, etc.

What are Fabric domains?

In Fabric, a domain is a way of logically grouping together all the data in an organization that is relevant to a particular area or field.

Reference: <https://learn.microsoft.com/en-us/fabric/governance/domains>





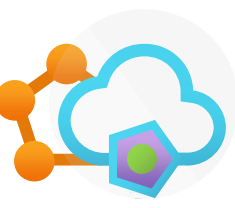
Architecture Agenda:

δ

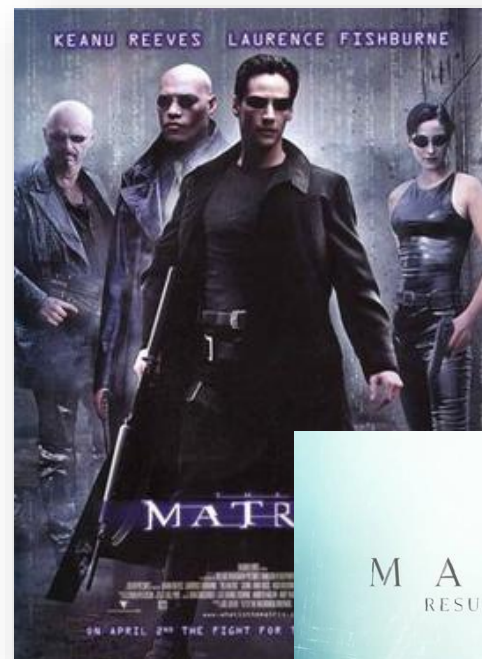
λ

K



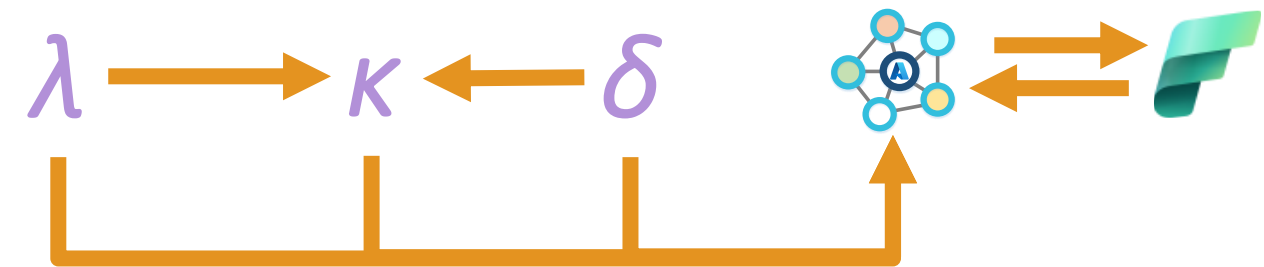


Final thoughts from me...



An Evolution of Data Platform Architectures

Lambda, Kappa, Delta, Mesh & Fabric



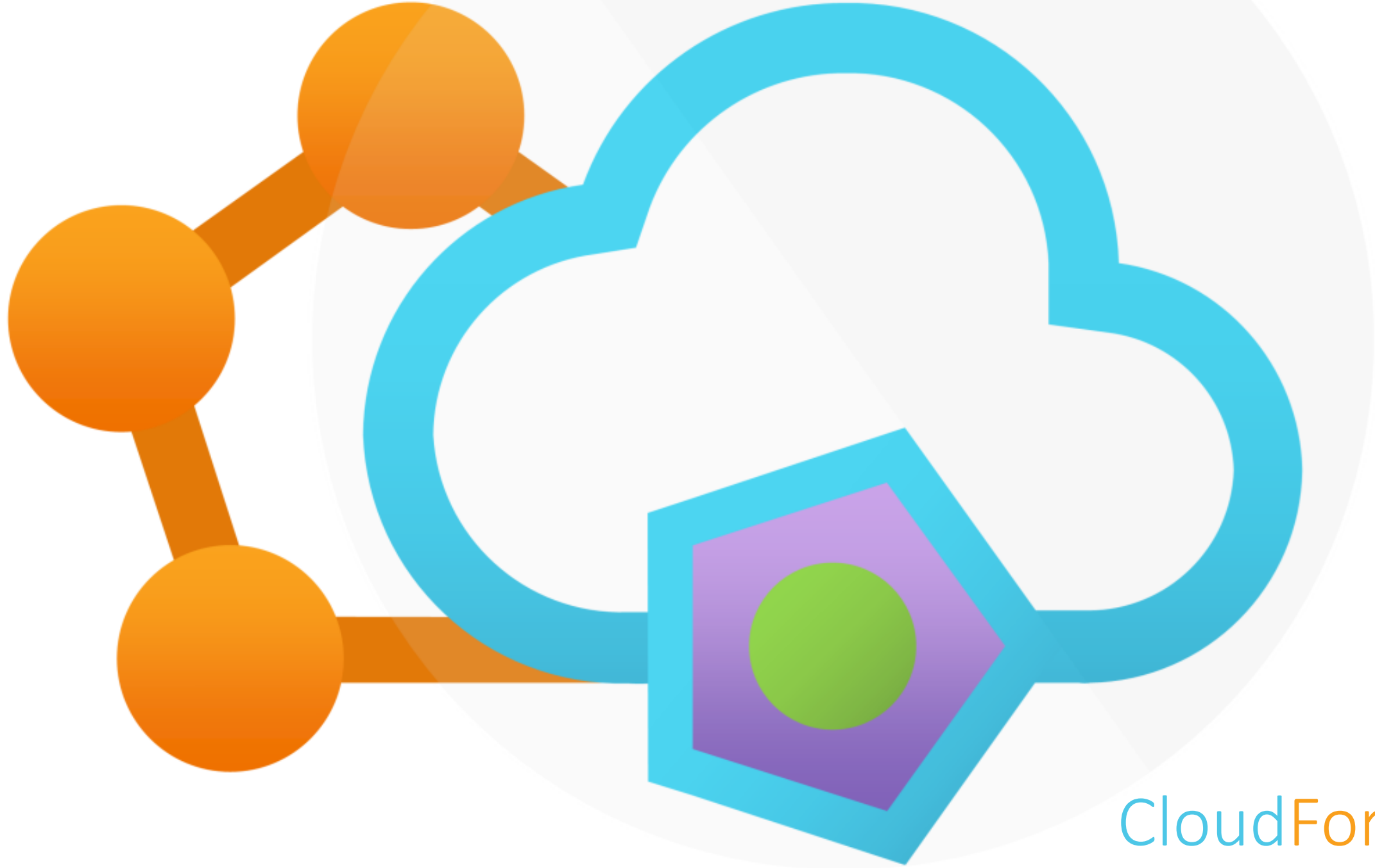
Q: What about a medallion architecture?



Source Infrastructure	Microsoft Azure	Microsoft Fabric	
Sources	Bronze	Silver	Gold
<ul style="list-style-type: none"> Any Data Structure. Any Technology. Operational Data Stores. Normalised Datasets. 	<ul style="list-style-type: none"> Simple Data Structures. Using Source Schemas. Change Data Capture. Audit Columns Applied. Limited Data Retention. 	<ul style="list-style-type: none"> Resilient Data Entities. Controlled Schemas. Cleansing, mapping, labelling & tokenisation. Merged Records. Complete History. 	<ul style="list-style-type: none"> Resilient Data Entities. Modelled Schemas. Output Aligned. Domain Orientated. De-Normalised Datasets.

Q: Should we be considering a solution/technology stack that offers all these capabilities?

A: Yes! ✓



Thank You



- mrpaulandrew.com
- paul@mrpaulandrew.com
- [In/mrpaulandrew](https://www.linkedin.com/in/mrpaulandrew)
- [@mrpaulandrew](https://twitter.com/mrpaulandrew)



- <https://cloudformations.org>
- contactus@cloudformations.org
- [In/CloudFormations](https://www.linkedin.com/company/cloudformations)
- [@CloudFormsLtd](https://twitter.com/CloudFormsLtd)
- [CloudFormationsLtd](https://www.facebook.com/CloudFormationsLtd)

[CloudFormations.org/Community-Content](https://cloudformations.org/Community-Content)

