

# Data Lakehouse, Data Mesh, and Data Fabric

(the alphabet soup of data architectures)

James Serra

Data & AI Solution Architect

Microsoft

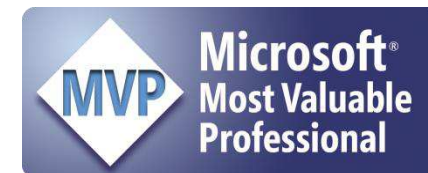
[jameserra@microsoft.com](mailto:jameserra@microsoft.com)

Blog: [JamesSerra.com](https://JamesSerra.com)



# About Me

- Microsoft, Data & AI Solution Architect in Microsoft Consulting Services (MCS), now called Industry Solutions Delivery (ISD)
- At Microsoft for most of the last eight years, with a brief stop at EY
- Was previously a Data & AI Architect at Microsoft for seven years
- In IT for 35 years, worked on many BI and DW projects
- Worked as desktop/web/database developer, DBA, BI and DW architect and developer, MDM architect, PDW/APS developer
- Been perm employee, contractor, consultant, business owner
- Presenter at PASS Summit, SQLBits, Enterprise Data World conference, Big Data Conference Europe, SQL Saturdays
- Blog at [JamesSerra.com](http://JamesSerra.com)
- Former SQL Server MVP
- Author of book "Reporting with Microsoft SQL Server 2012"



# Agenda

- Data Warehouse
- Data Lake
- Modern Data Warehouse
- Data Fabric
- Data Lakehouse
- Data Mesh

I tried to figure out all these data platform buzzwords on my own...

And it did not turn out well:



Let's prevent that from happening...

***The view in this deck are my own and not that of Microsoft!***



# What is a Data Warehouse and why use one?

(or, why do we need a copy of the source data?)

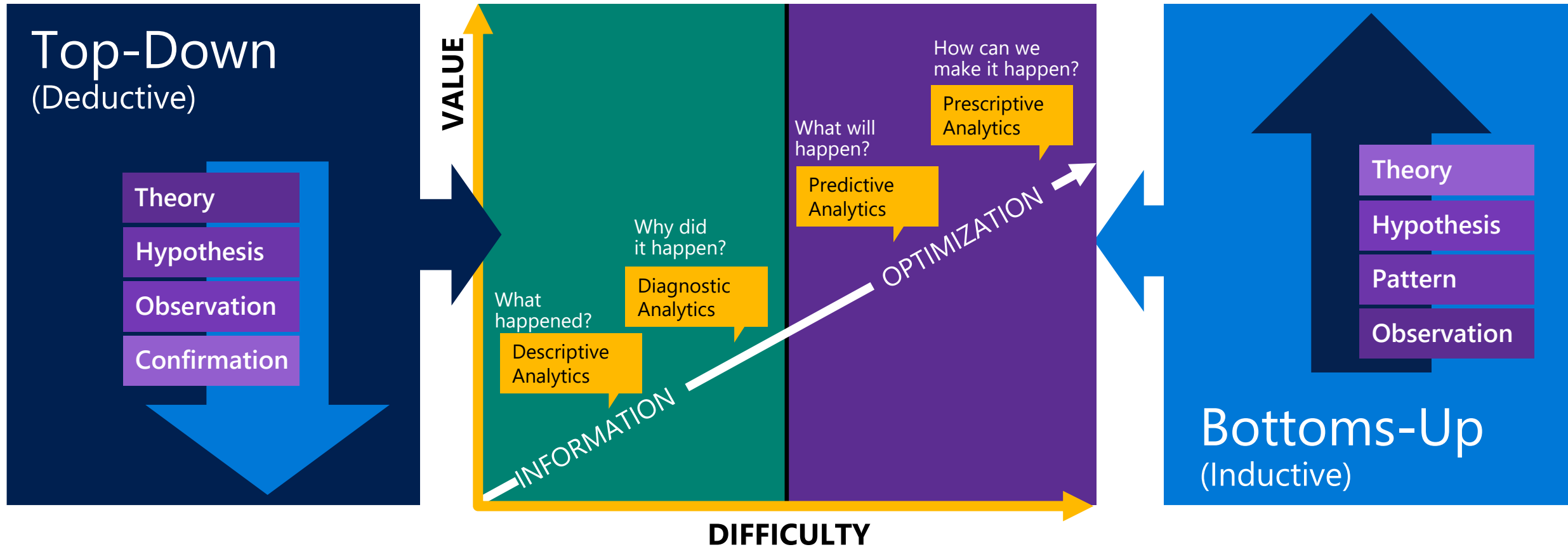
A data warehouse is where you store data from multiple data sources to be used for historical and trend analysis reporting. It acts as a central repository for many subject areas and contains the "single version of truth". It is NOT to be used for OLTP applications.

Reasons for a data warehouse:

- Reduce stress on production system
- Optimized for read access, sequential disk scans
- Integrate many sources of data
- Keep historical records (no need to save hardcopy reports)
- Restructure/rename tables and fields, model data
- Protect against source system upgrades
- Use Master Data Management, including hierarchies
- No IT involvement needed for users to create reports
- Improve data quality and plugs holes in source systems
- One version of the truth
- Easy to create BI solutions on top of it (i.e. Power BI tabular model)
- Don't need to provide security access for many users to the production systems
- Make better business decisions by getting greater insights into your company

[Why You Need a Data Warehouse](#)

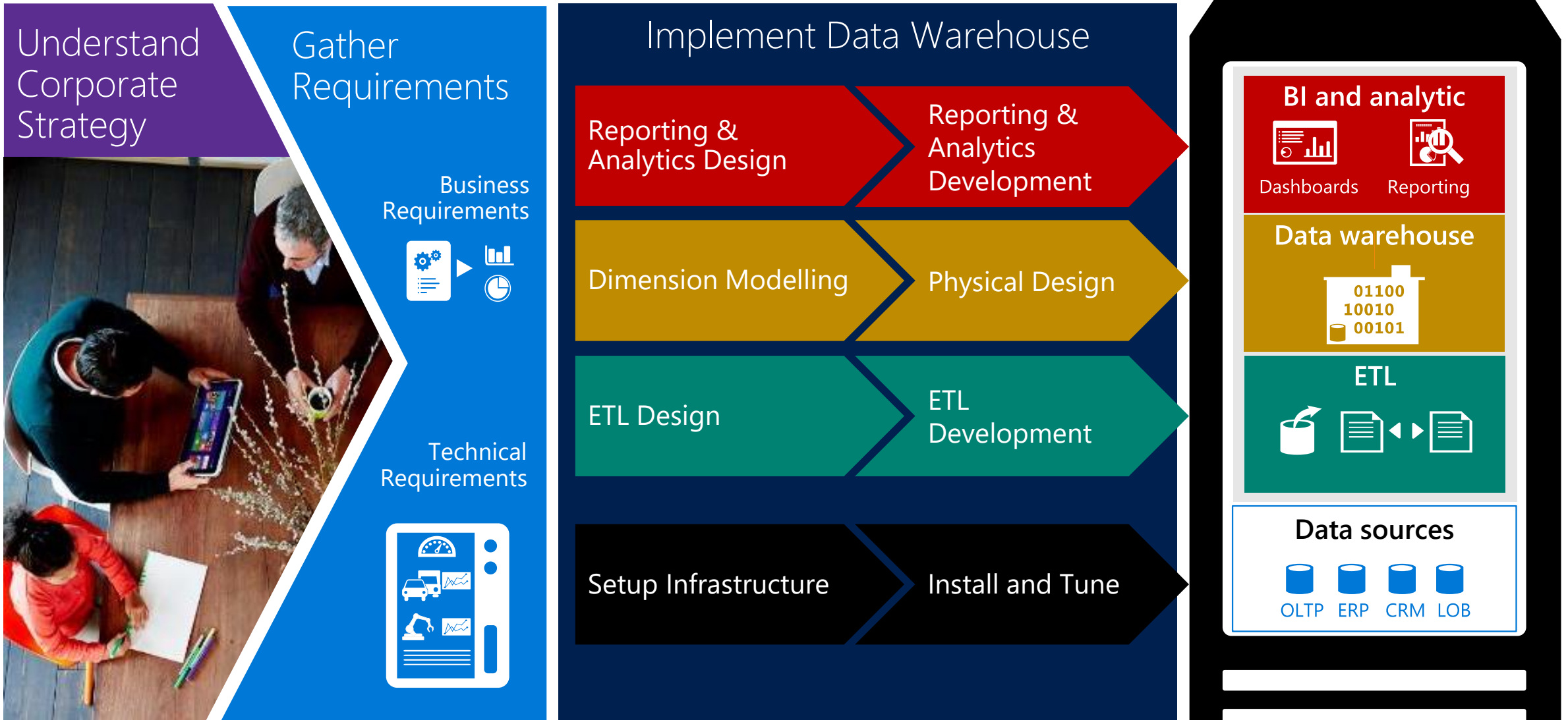
# Two Approaches to getting value out of data: Top-Down + Bottoms-Up



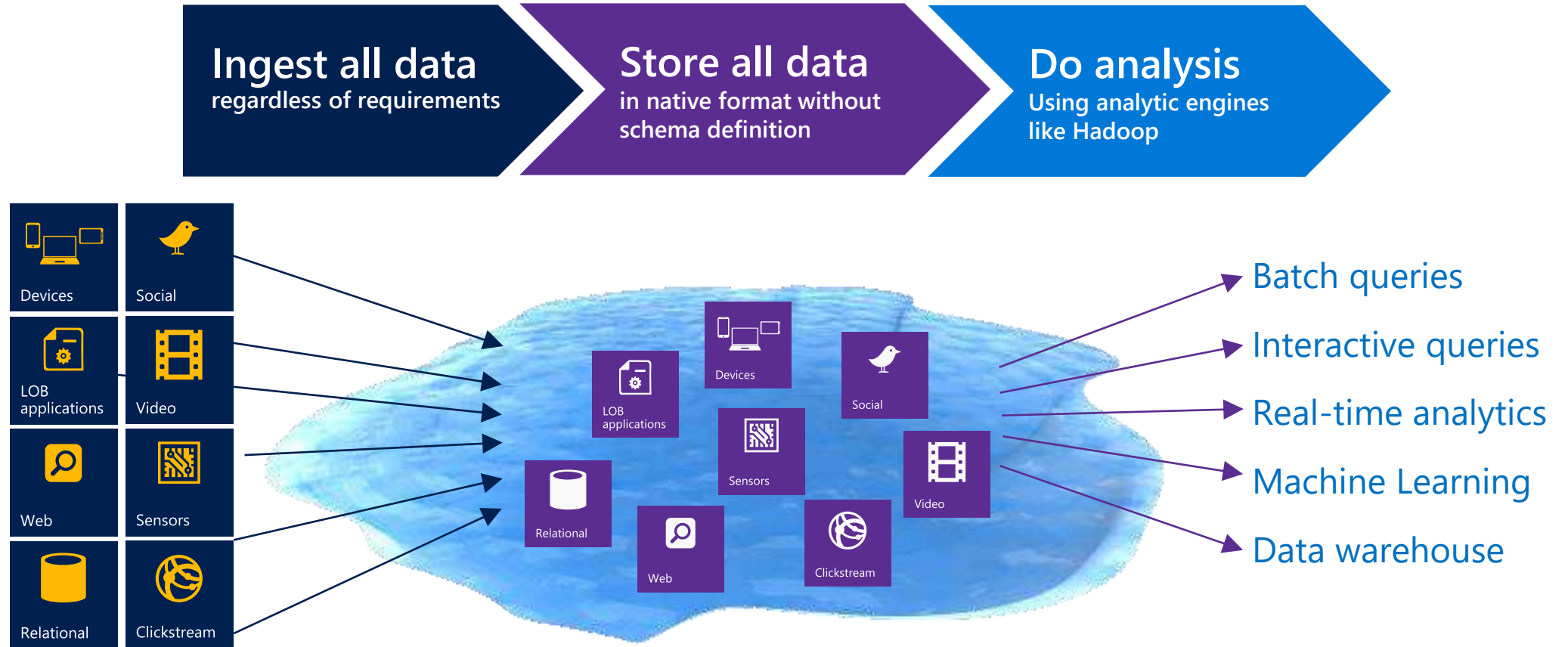
- Know the questions to ask
- Lot's of upfront work to get the data to where you can use it
- Model first (schema-on-write)

- Don't know the questions to ask
- Little upfront work needs to be done to start using data
- Model later (schema-on-read)

# Data Warehousing Uses A Top-Down Approach



# The "data lake" Uses A Bottoms-Up Approach





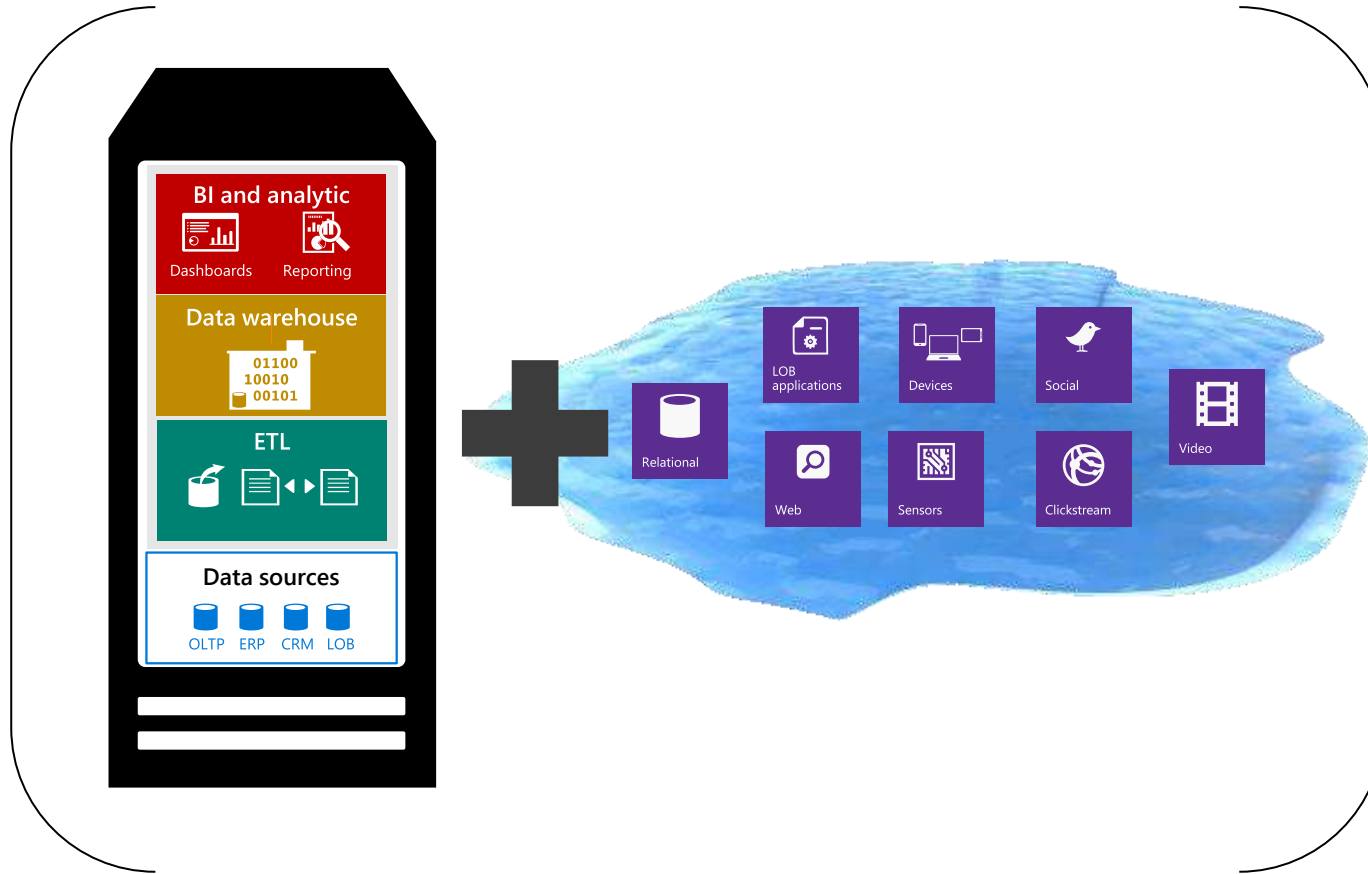
# Data Lake + Data Warehouse Better Together

What happened?

Descriptive  
Analytics

Why did it happen?

Diagnostic  
Analytics



What will happen?

Predictive  
Analytics

How can we make it happen?

Prescriptive  
Analytics

# What is a data lake and why use one?

A schema-on-read storage repository that holds a vast amount of raw data in its native format until it is needed.

Reasons for a data lake:

- Inexpensively store unlimited data
- Centralized place for multiple subjects (single version of the truth)
- Collect all data “just in case” (data hoarding). The data lake is a good place for data that you “might” use down the road
- Easy integration of differently-structured data
- **Store data with no modeling – “Schema on read”**
- Complements enterprise data warehouse (EDW)
- **Frees up expensive EDW resources for queries instead of using EDW resources for transformations (avoiding user contention)**
- Wanting to use technologies/tools (i.e Databricks) to refine/filter data that do the refinement quicker/better than your EDW
- **Quick user access to data for power users/data scientists (allowing for faster ROI)**
- **Data exploration to see if data valuable before writing ETL and schema for relational database, or use for one-time report**
- Allows use of Hadoop tools such as ETL and extreme analytics
- Place to land IoT streaming data
- On-line archive or backup for data warehouse data (i.e. keep three years of data in DW and have older data in data lake with an external table pointing to it)
- With Hadoop/ADLS, high availability and disaster recovery built in
- It can ingest large files quickly and provide data redundancy
- ELT jobs on EDW are taking too long because of increasing data volumes and increasing rate of ingesting (velocity), so offload some of them to the Hadoop data lake
- Have a backup of the raw data in case you need to load it again due to an ETL error (and not have to go back to the source). You can keep a long history of raw data
- Allows for data to be used many times for different analytic needs and use cases
- Cost savings and faster transformations: storage tiers with lifecycle management; separation of storage and compute resources allowing multiple instances of different sizes working with the same data simultaneously vs scaling data warehouse; low-cost storage for raw data saving space on the EDW
- **Extreme performance for transformations by having multiple compute options each accessing different folders containing data**
- The ability for an end-user or product to easily access the data from any location

# Data Lake with DW use cases

## Data Lake

### Staging & preparation

- Data scientists/Power users
- Batch processing
- Data refinement/cleaning
- ETL workloads
- Store older/backup data
- Sandbox for data exploration
- One-time reports
- Quick access to data
- Don't know questions

## Data Warehouse

### Serving, Security & Compliance

- Business people
- Low latency
- Complex joins
- Interactive ad-hoc query
- High number of users
- Additional security
- Large support for tools
- Dashboards
- Easily create reports (Self-service BI)
- Know questions

# Enterprise Data Maturity Stages

Digital transformation accelerates along this journey

## STAGE 1: Reactive

Structured data is transacted and locally managed. Data used reactively

## STAGE 2: Informative

Structured data is managed and analyzed centrally and informs the business

## STAGE 3: Predictive

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

## STAGE 4: Transformative

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



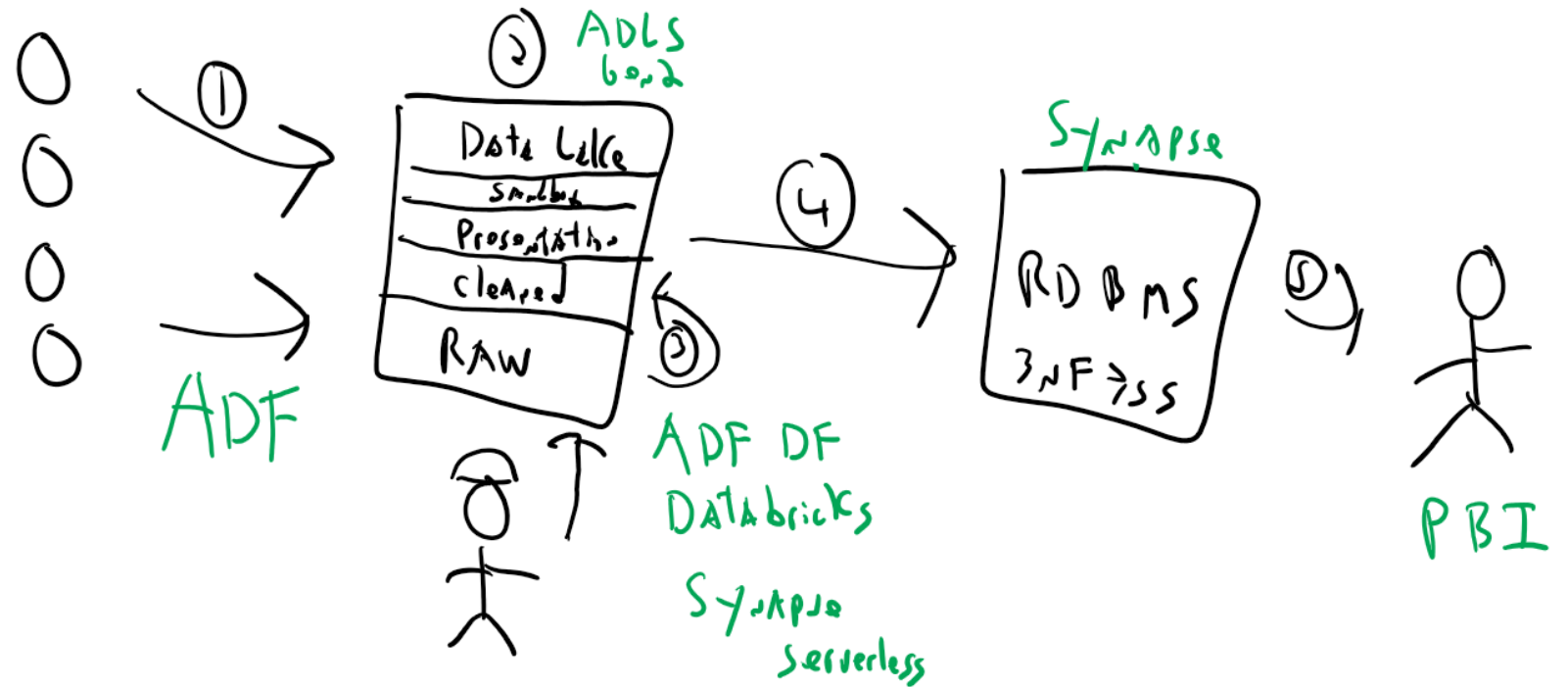
← Rear-view  
mirror →

← Real-time  
intelligence →

# Modern Data Warehouse

Modern Data Warehouse (MDW)

- 1) Ingest
- 2) Store
- 3) Transform
- 4) Model
- 5) Visualize/ML





# Data Fabric

Data Fabric adds to a modern data warehouse:

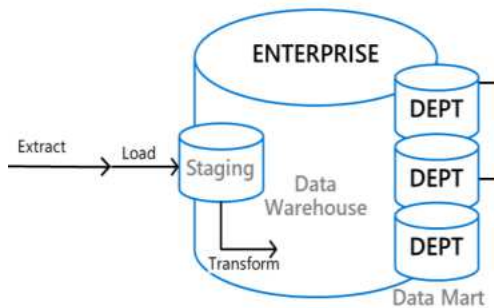
- Data access
- Data policies
- Metadata catalog/Lineage
- Master Data Management (MDM)
- Data virtualization
- Real-time processing
- Data scientist tools
- APIs
- Building blocks/Services
- Products

Bottom line: Additional technology to source more data, secure it, and make it available

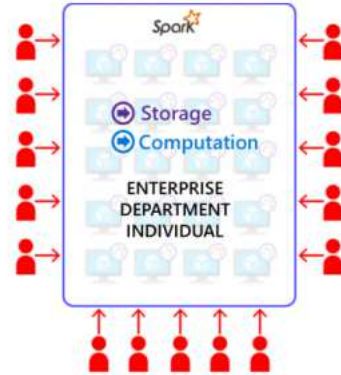
[Data Fabric defined](#)

# Data Lakehouse

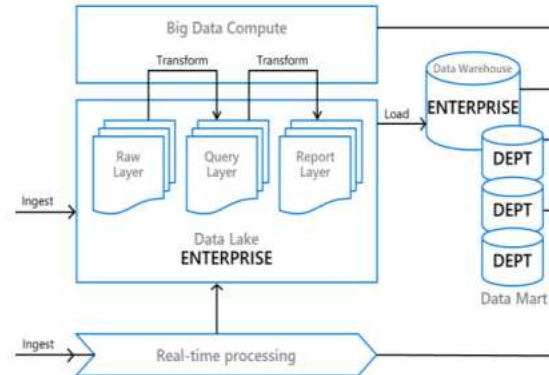
**Late 1980s**  
Data Warehouse



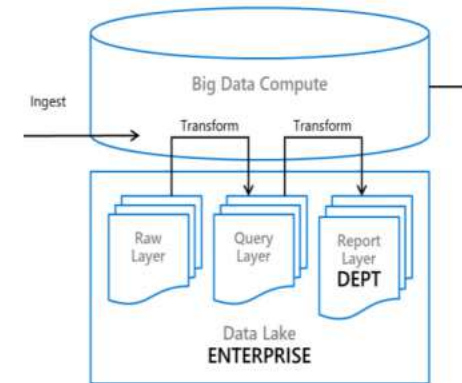
**Late 2000s**  
Data Lake



**Mid 2010s**  
Cloud Data Platform



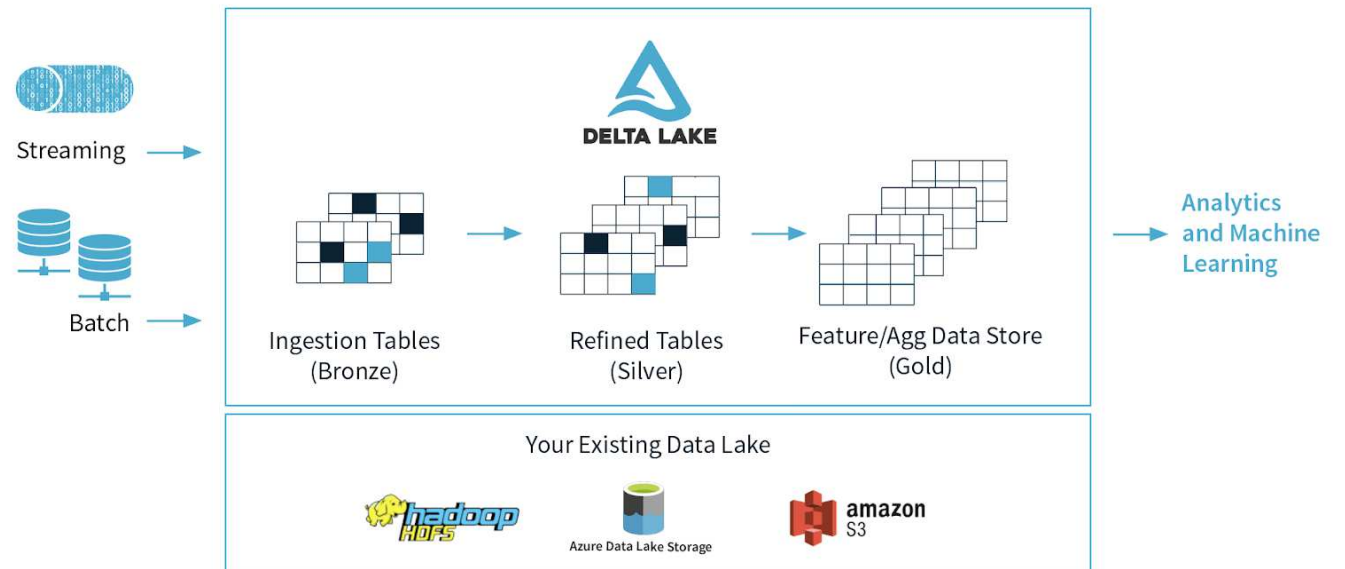
**2020**  
Data Lakehouse



# Delta Lake

Top features:

- ACID transactions
- Time travel (data versioning enables rollbacks, audit trail)
- Streaming and batch unification
- Schema enforcement
- Supports commands DELETE, UPDATE, and MERGE
- Performance improvement



# Use cases for Data Lakehouse

Today's data architectures commonly suffer from four problems:

- Reliability: Keeping the data lake and warehouse consistent
- Data staleness: Data in warehouse is older
- Limited support for advanced analytics: Top ML systems don't work well on warehouses
- Total cost of ownership: Extra cost for data copied to warehouse

# Concerns skipping relational database

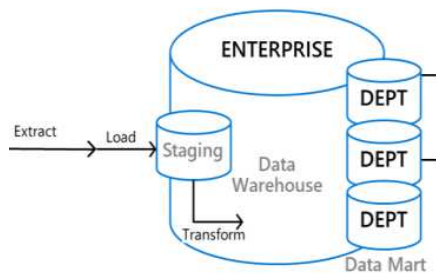
- Speed: Relational databases faster, especially Massively Parallel Processing (MPP)
- Security: No RLS, column-level, dynamic data masking
- Complexity: Metadata separate from data, file-based world
- Concurrency: Multiple reads of a file at the same time
- Missing features: Referential integrity, TDE, workload management, MDM; other features lock you into Spark
- Having to use Spark SQL instead of T-SQL
- People are used to using a relational database

Azure Synapse: starting to see data lake only solutions because can use T-SQL, Power BI (speed, RLS), cost savings with Serverless

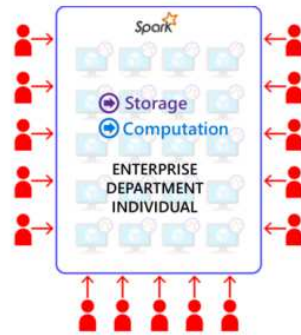


# Data Mesh

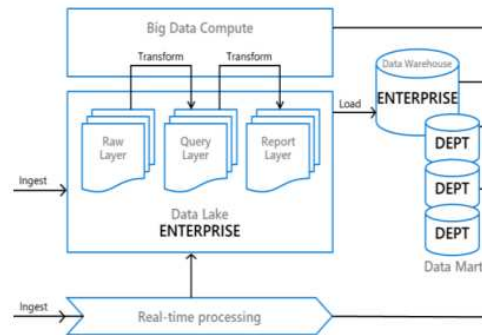
**Late 1980s**  
Data Warehouse



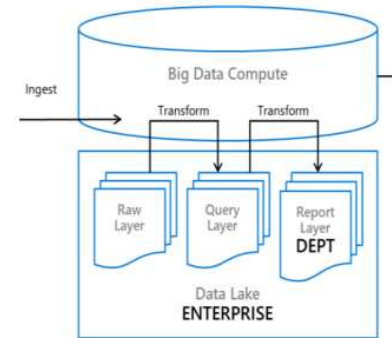
**Late 2000s**  
Data Lake



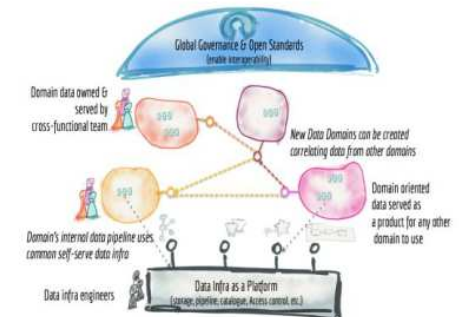
**Mid 2010s**  
Cloud Data Platform



**2020**  
Data Lakehouse



**2021**  
**Data Mesh??**

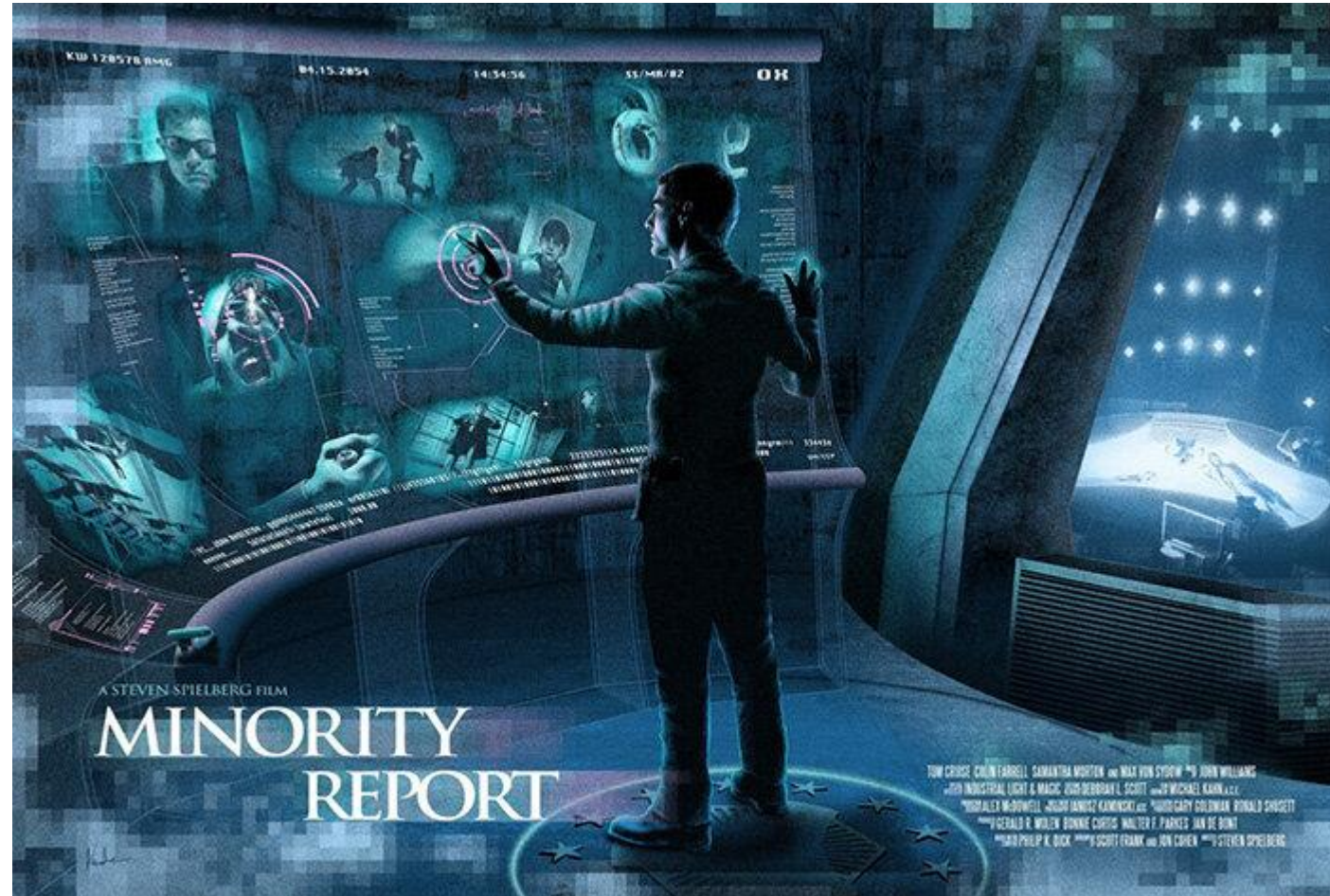


Centralization

Decentralization

# Data Mesh in theory

*Lots of things sound great in theory...*



*Data Mesh is a concept, not technology*

# Data Mesh - Overview

Data mesh is an intentionally designed distributed data architecture, under centralized governance and standardization for interoperability, enabled by a shared and harmonized self-serve data infrastructure

## Data Mesh Principles

### #1) Domain Ownership

Decentralize and distribute responsibility to people who are closest to the data in order to support continuous change and scalability (i.e. manufacturing, sales, supplier)

### #2) Data as a product

Analytical data provided by the domains are treated as a product and the consumers of that data are treated as customers (domain teams, API code, data and metadata, infrastructure)

### #3) Self-serve data infrastructure as a platform

High-level abstraction of diverse infrastructure that removes complexity and friction of provisioning and managing the lifecycle of data products (i.e. storage, compute, data pipeline, access control)

### #4) Federated computational governance

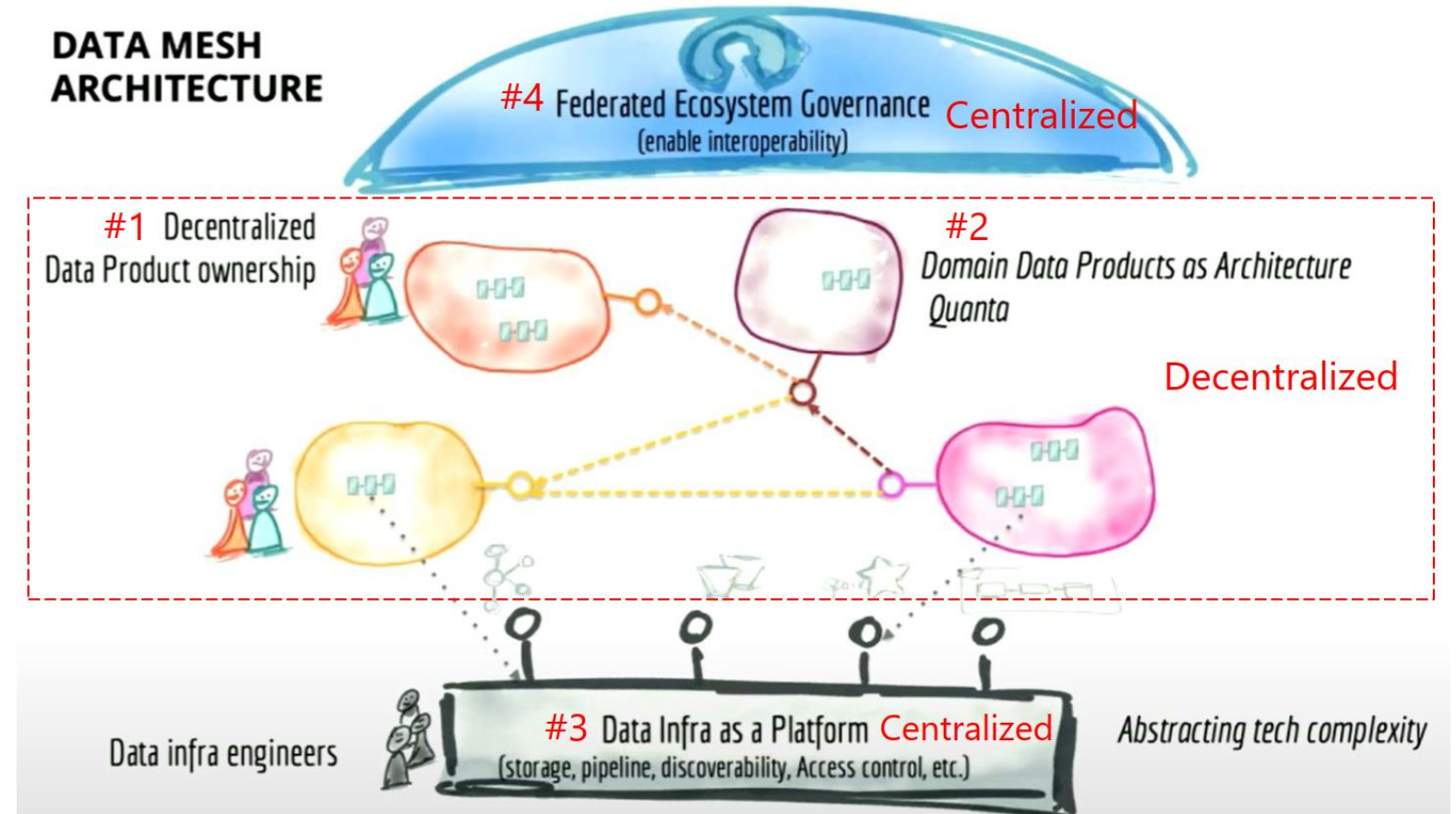
Architect global decisions and standards for interoperability, while respecting autonomy of local domains, and implement global policies effectively (i.e. data quality, data security, regulations, data modeling)

# Data Mesh

***Data Mesh is a concept, not a product***

It's a mindset shift where you go from:

- Centralized ownership to decentralized ownership
- Pipelines as first-class concern to domain data as first-class concern
- Data as a by-product to data as a product
- A siloed data engineering team to cross-functional domain-data teams
- A centralized data lake/warehouse to an ecosystem of data products





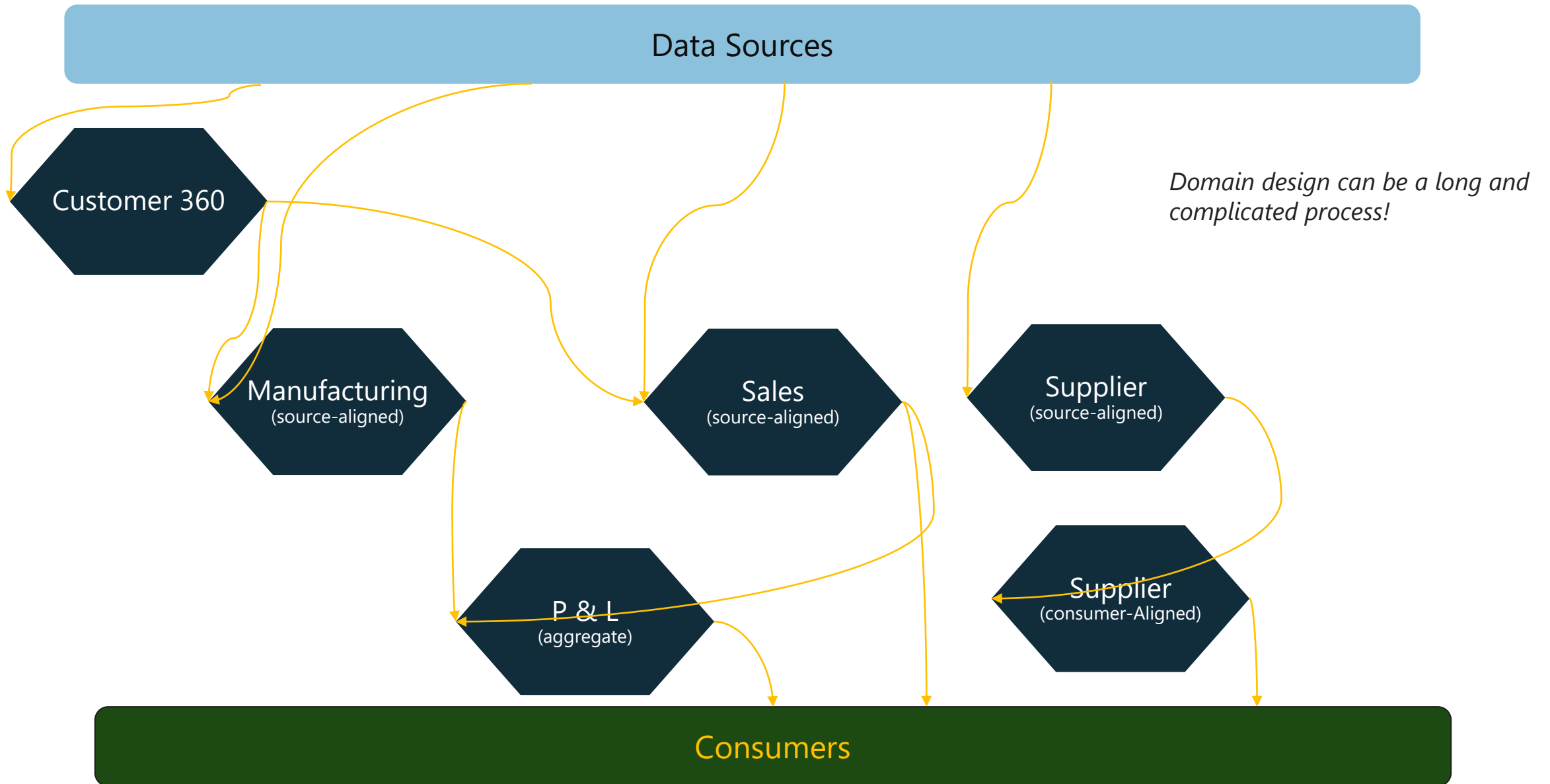
# Use cases for Data Mesh

Data mesh tries to solve four challenges with a centralized data lake/warehouse:

- Lack of ownership: who owns the data – the data source team or the infrastructure team?
- Lack of quality: the infrastructure team is responsible for quality but does not know the data well
- Organizational scaling: the central team becomes the bottleneck, such as with an enterprise data lake/warehouse
- Technical scaling: current big data solutions can't keep up with additional data requirements



# Data Mesh – Logical Architecture



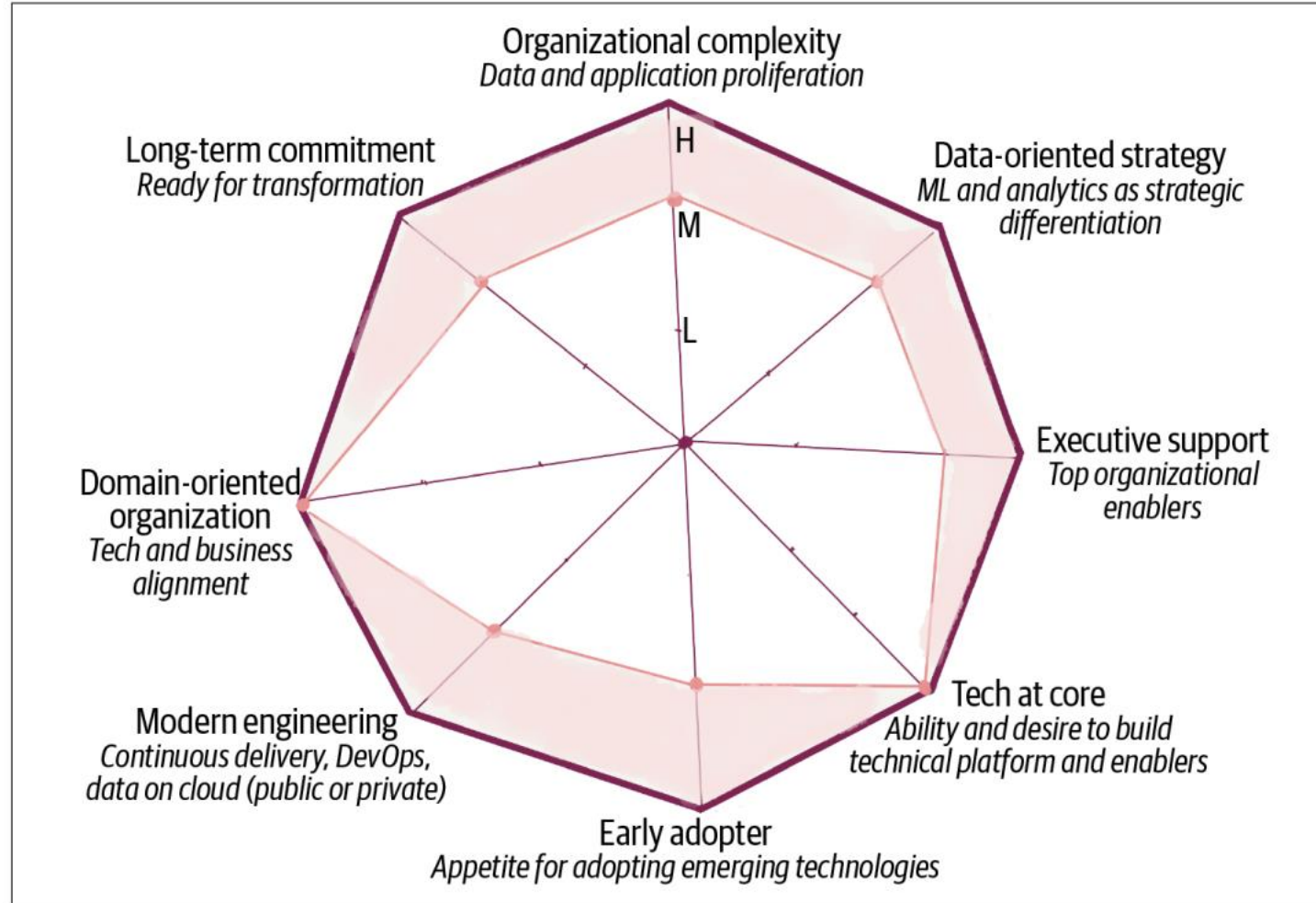
# Concerns with Data Mesh

- No standard definition of a data mesh
- Huge investment in organizational change and technical implementation
- Performance of combining data from multiple domains
- Duplication of data for performance reasons
- Getting quality engineering people for each domain
- Inconsistent technical implementations for the domains
- Domains don't want to wait for a data mesh
- Need incentives for each domain to counter extra work
- Self-serve approach of data requests could be challenging
- Duplication of data and ingestion platform
- Creation of data silos for domains not able to join data mesh
- Not seeing the big picture for combining data

[Data Mesh: Centralized vs decentralized data architecture](#)

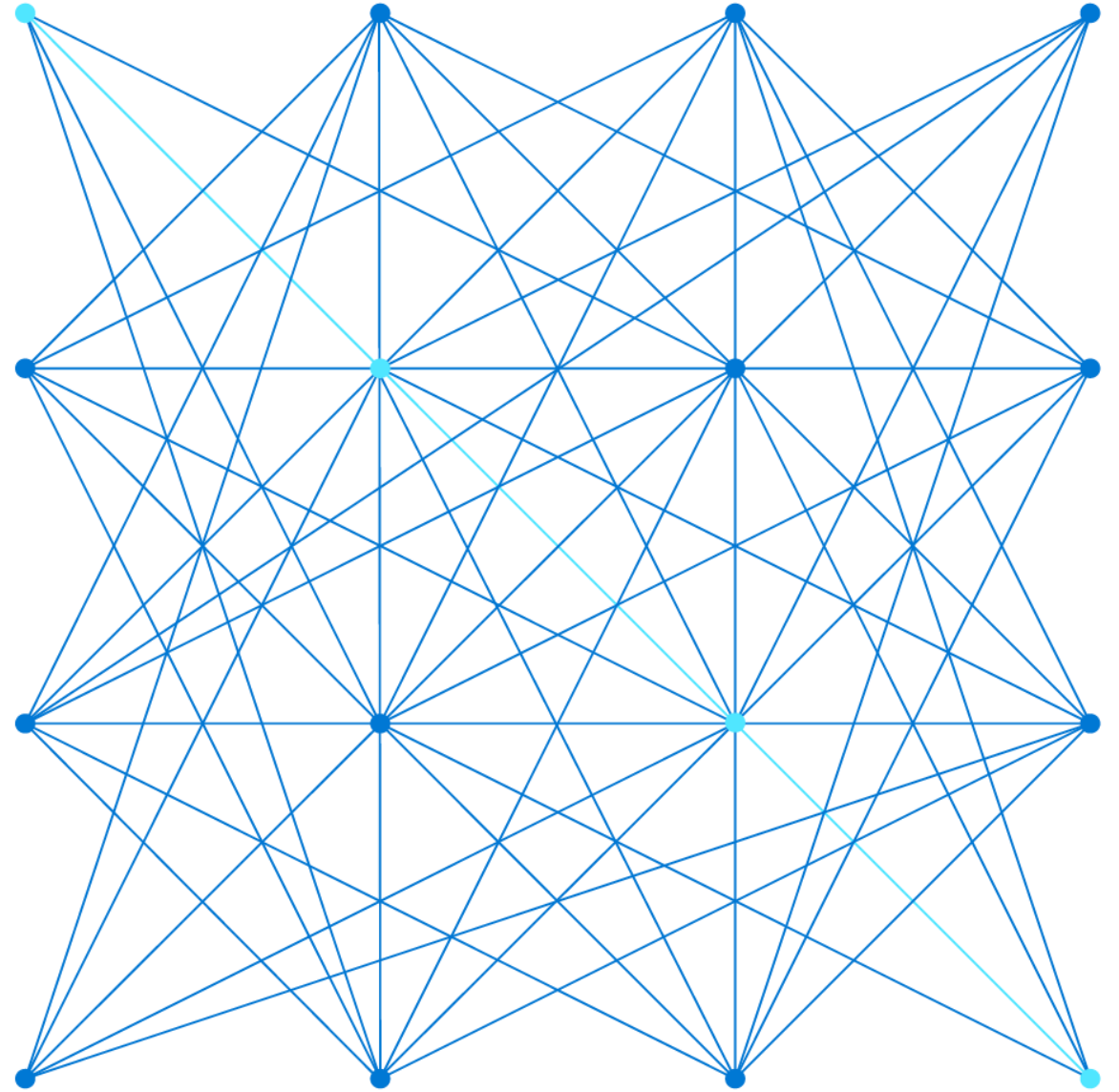
[Data Mesh: Centralized ownership vs decentralized ownership](#)

# Should you adopt data mesh today?



Need to score medium or high in ALL categories

# Data Mesh on Azure



# Cloud Scale Analytics

Cloud Scale Analytics is an architecture approach and reference implementation that enables effective construction and operationalization of landing zones on Azure, at scale and aligned with Azure Roadmap and Cloud Adoption Framework.

What is Cloud Scale Analytics?

**A scalable analytics framework designed to enable customers building a data platform.**

- *Supports multiple topologies ranging across Data Centric, Lakehouse, Data Fabric and Data Mesh*
- Based on inputs from PG and a diverse international group of specialists working with a range of customers
- Separate guidance tailored to Small-Medium and Large enterprises
- ~80% prescribed viewpoint with 20% client customization

Helps you to create:

Data Landing Zones

Data Management Zone





ALL AROUND AZURE

# Unlocked: Cloud Scale Analytics

A **Practical** introduction. Learn **how** to construct and promote a **data mesh implementation** strategy



LIVE on LearnTV in your timezone

<https://aka.ms/aaa-unlocked-data>

**April 21, 2022**

7:00am - 9:00am PT | 14:00 - 16:00 UTC



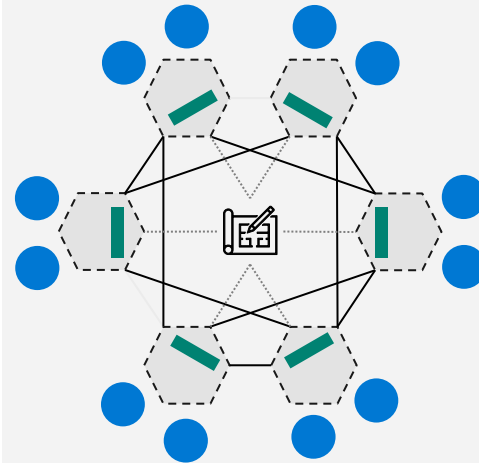
# Data Mesh on Azure Resources

- Piethein Strengholt: [Blog - Implementing Data Mesh on Azure](#), [Blog – Data Mesh topologies](#), [Book - Data Management at Scale: Best Practices for Enterprise Architecture](#)
- Cloud Adoption Framework: [Azure data management and analytics scenario](#)
- Data Management & Analytics Scenario - Data Management Zone: [Github](#)
- Data Management & Analytics Scenario - Data Landing Zone: [Github](#)
- Enterprise-Scale - Reference Implementation: [Github](#)
- Microsoft doc: [A financial institution scenario for data mesh](#)
- [Provision three Azure Data Lake Storage Gen2 accounts for each data landing zone](#)
- [Overview of Azure Data Lake Storage for the data management and analytics scenario](#)
- [The best practices for organizing Synapse workspaces and lakehouses](#)

# Governance Topologies : Different Approaches

**Private**  
(do not share publicly)

## Mesh Type 2

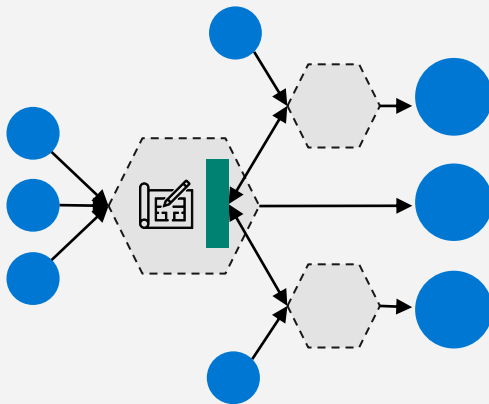


- Domains use the same technology
- Each domain has its own storage that is the same technology

Centralised  
(Control)

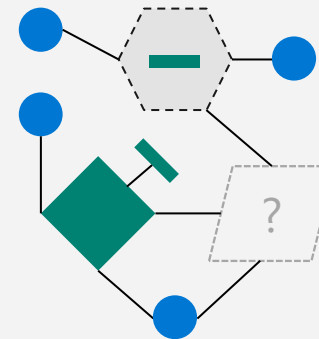
Distributed  
(Agility)

## Mesh Type 1



- Domains use the same technology
- Data is kept in one enterprise data lake with each domain getting its own container/folder

## Mesh Type 3



- Domains can use any technology they want
- Each domain has its own storage that can be any technology

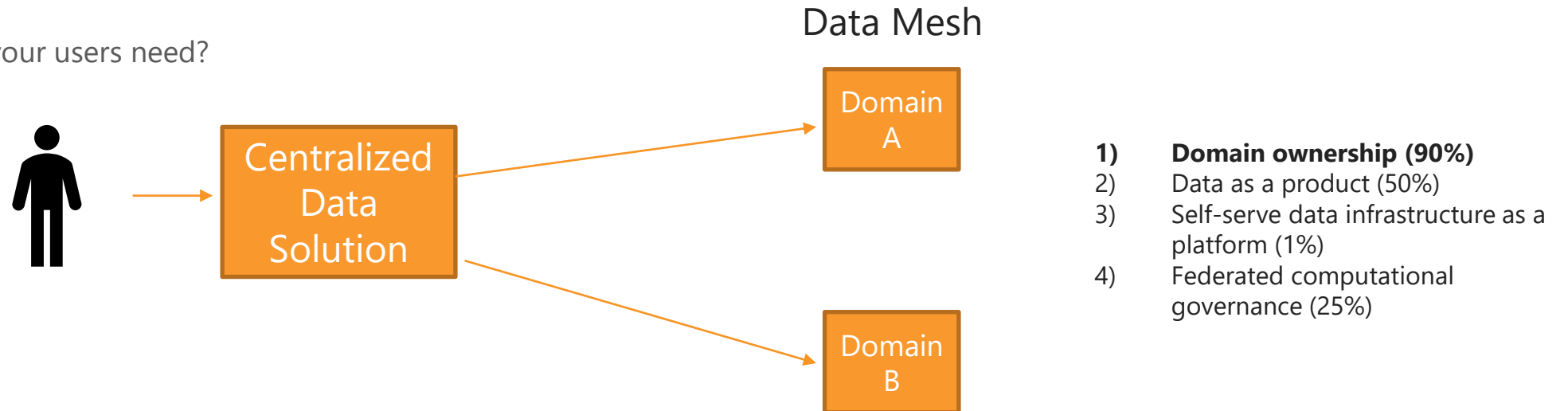
# Future

***This view is my own and not that of Microsoft!***

In the end, I predict data mesh will become an extension to a centralized data solution for a small percentage of solutions.

There will be a very small percentage of solutions that are 100% true to the pure data mesh concept (assuming mesh type 1 and 2 are true to the data mesh concept). *Ask ten people what a data mesh is and you will get eleven answers!* Some of the concepts of a data mesh will be used in a larger percentage of solutions.

Always ask: What do your users need?



[Rethinking the Data Mesh Architecture: Apply it Piecemeal \(eckerson.com\)](https://eckerson.com)

***Data Mesh concepts help with a better way of thinking how to get value out of data***

# Q & A



James Serra, Microsoft, Data & AI Solution Architect

Email me at: [jamesserra3@gmail.com](mailto:jamesserra3@gmail.com)

Follow me at: @JamesSerra

Link to me at: [www.linkedin.com/in/JamesSerra](https://www.linkedin.com/in/JamesSerra)

Visit my blog at: [JamesSerra.com](https://JamesSerra.com)