# From Metadata to Lakehouse
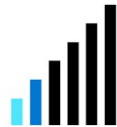
**Andreas Bergstedt**
**Data and AI Global Black Belt**
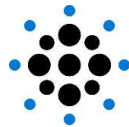
# Introduction / Agenda

**The Evolution of ETL**

**Why Lakehouse**

**Metadata Driven Engineering**

**Driving at Scale**

**Key take-aways**

# Introduction / Agenda

**The Evolution of ETL**

## Why Lakehouse

# MACHINE LEARNING, ARTIFICIAL INTELLIGENCE, AND DATA (MAD) LANDSCAPE 2021

## INFRASTRUCTURE

- STORAGE
- HADOOP
- DATA LAKES
- DATA WAREHOUSES
- STREAMING / IN-MEMORY
- RDBMS
- NoSQL DATABASES
- NewSQL DATABASES
- REAL TIME DATABASES
- GRAPH DBs
- MPP DBs
- ETL / ELT / DATA TRANSFORMATION
- REVERSE ETL
- DATA INTEGRATION
- DATA GOVERNANCE & ACCESS
- DATA CATALOG AND DISCOVERY
- METRICS STORE
- LOG ANALYTICS
- PRIVACY & SECURITY
- DATA OBSERVABILITY
- MGMT / MONITORING
- SERVER-LESS
- CLUSTER SVCS
- DATA QUALITY
- QUERY ENGINE
- SEARCH

## ANALYTICS

- BI PLATFORMS
- VISUALIZATION
- DATA ANALYST PLATFORMS
- AUGMENTED ANALYTICS

## MACHINE LEARNING & ARTIFICIAL INTELLIGENCE

- DATA SCIENCE NOTEBOOKS
- DATA SCIENCE PLATFORMS
- ML PLATFORMS
- DATA GENERATION & LABELLING
- MODEL BUILDING
- FEATURE STORE
- DEPLOYMENT & PRODUCTION
- COMPUTER VISION
- SPEECH
- NLP
- SYNTHETIC MEDIA
- HORIZONTAL AI
- GPU DBS & CLOUD
- AI HARDWARE

## APPLICATIONS – ENTERPRISE

- SALES
- MARKETING - B2B
- MARKETING - B2C
- CUSTOMER EXPERIENCE / SERVICE
- HUMAN CAPITAL
- LEGAL
- REGTECH & COMPLIANCE
- FINANCE
- AUTOMATION & RPA
- SECURITY
- PARTNERSHIPS

## APPLICATIONS – INDUSTRY

- ADVERTISING
- EDUCATION
- REAL ESTATE
- GOVT & INTELLIGENCE
- COMMERCE
- FINANCE - LENDING
- INSURANCE
- FINANCE - INVESTING
- HEALTHCARE
- LIFE SCIENCES
- TRANSPORTATION
- AGRICULTURE
- INDUSTRIAL
- OTHER

## OPEN SOURCE

- FRAMEWORKS
- FORMAT
- QUERY / DATA FLOW
- DATA ACCESS
- DATABASES
- ORCHESTRATION
- INFRA-STRUCTURE
- DATA OPS
- STREAMING & MESSAGING
- STAT TOOLS & LANGUAGES
- ML OPS & INFRA
- AI / MACHINE LEARNING / DEEP LEARNING
- SEARCH
- LOGGING & MONITORING
- VISUALIZATION
- COLLABORATION
- SECURITY

## DATA SOURCES & APIs

- DATA MARKETPLACES & DISCOVERY
- FINANCIAL & ECONOMIC DATA
- AIR / SPACE / SEA
- PEOPLE / ENTITIES
- LOCATION INTELLIGENCE
- OTHER

## DATA RESOURCES

- DATA SERVICES
- INCUBATORS & SCHOOLS
- RESEARCH
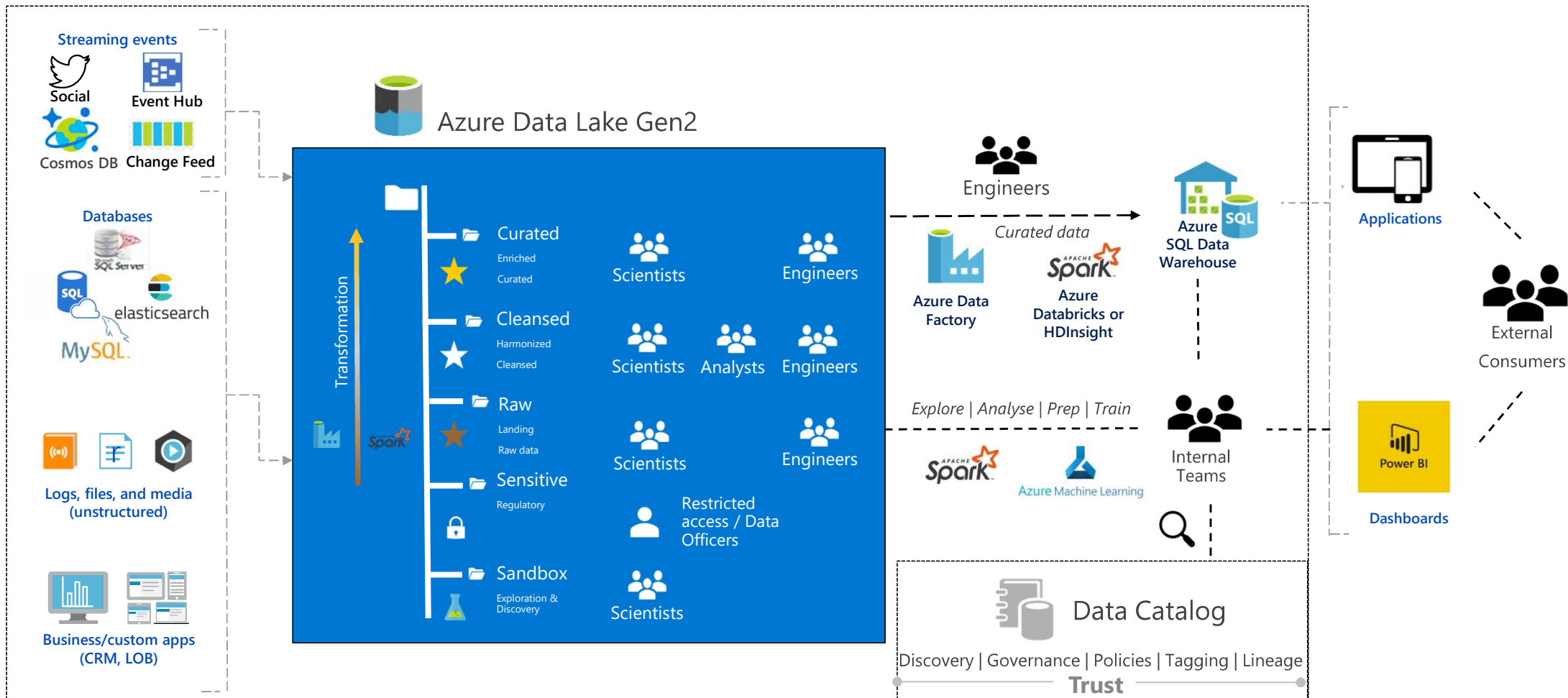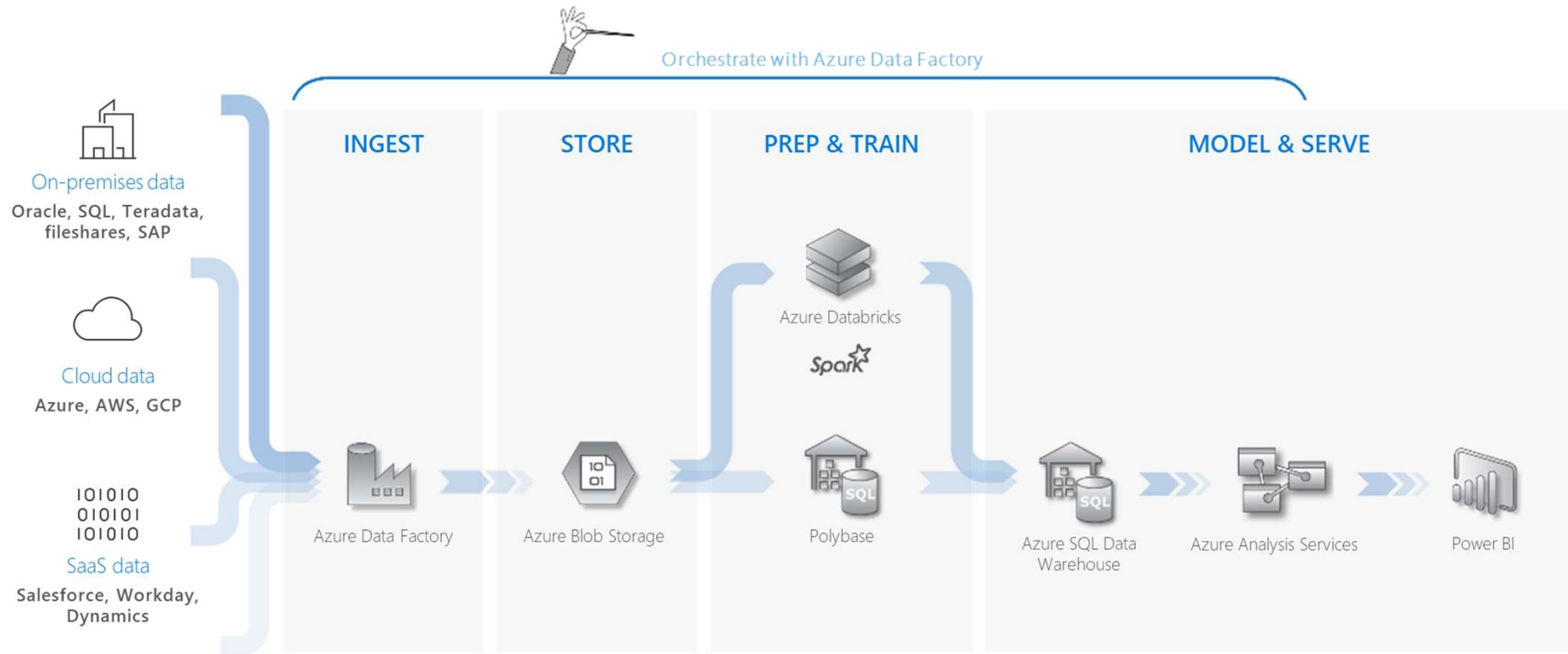
Microsoft Azure


APACHE Spark


Azure Synapse Analytics


DELTA LAKE


databricks


ICEBERG

# Data Lake Architecture – Concepts, Tools & Process

# Modernize your enterprise data warehouse at scale

Orchestrate with Azure Data Factory

| INGEST | STORE | PREP & TRAIN | MODEL & SERVE |
|--------|-------|--------------|---------------|

**On-premises data**
Oracle, SQL, Teradata, fileshares, SAP

**Cloud data**
Azure, AWS, GCP

101010
010101
101010

**SaaS data**
Salesforce, Workday, Dynamics

Azure Databricks

Spark

Azure Data Factory

Azure Blob Storage

Polybase

Azure SQL Data Warehouse

Azure Analysis Services

Power BI

Microsoft Azure also supports other **Big Data** services like **Azure HDInsight**, **Azure SQL Database** and **Azure Data Lake** to allow customers to tailor the above architecture to meet their unique needs.

# Medallion architecture

Although the 3-layered design is common and well-known, there are many discussions on the scope, purpose, and best practices on each of these layers.

## Bronze layer
Typically raw, "as-is"

- Maintains the raw state in the structure "as-is"
- Data is immutable (read-only)
- Delivery-based partitioned tables, i.e., YYYYMMDD
- Mostly Parquet. Sometimes other formats
- Can be any combination of streaming and batch transactions
- May include extra metadata (schema)
- May be fed from a "mediation layer"
- Used for debugging, testing

## Silver layer
Matched and conformed

- Uses data quality rules for validation
- Usually only functional data
- Historization is merged (SCD2)
- Efficient storage format; Delta
- Versioning for rolling back
- Handles missing or incorrect data
- Usually enriched with reference data
- Source-oriented, although queryable and cluttered around subject areas
- Usually used by operational analytical teams

## Gold layer
Refined business-level

- What enterprises call data products: consumer-ready / user-friendly data
- Data is highly governed and well-documented
- Historization is applied only for the set of use cases or consumers
- Contains complex business rules, such as calculations and enrichments
- Efficient storage format; Delta
- Versioning for rolling back
- Might contain additional sub layers for sharing or distributing data

# Basic Lakehouse Architecture



**Operational systems**

**Blob landing zone**

**Process**

**Synapse Analytics**

**Process**

**Synapse Analytics**

**Ingest**
Copy data using Azure Data Factory

**Bronze (ADLS)**
Typically, raw and different file formats

**Silver (ADLS)**
Typically, cleansed and standardized file formats

**Gold (ADLS)**
Typically, data that is ready for consumption

**Serve**
Direct access to files

**Orchestrate** using Azure Data Factory

**Metadata Driven
Engineering**

# Available Frameworks

# The Lakehouse Pattern in Azure

# Metadata Driven Engineering Examples

# The Lakehouse Pattern in Microsoft Fabric



Data management landing zone

SQL Log

Governance using Microsoft Purview

Publish lineage

Data landing zone

Operational systems

Blob landing zone

Ingest

Process

Streaming Analytics    Spark

Process

Realtime Analytics    Spark

One Lake

SQL

PowerBI

Bronze (ADLS)
Typically, raw and layered per datetime

Silver (ADLS)
Typically, cleansed, historized and enriched

Gold (ADLS)
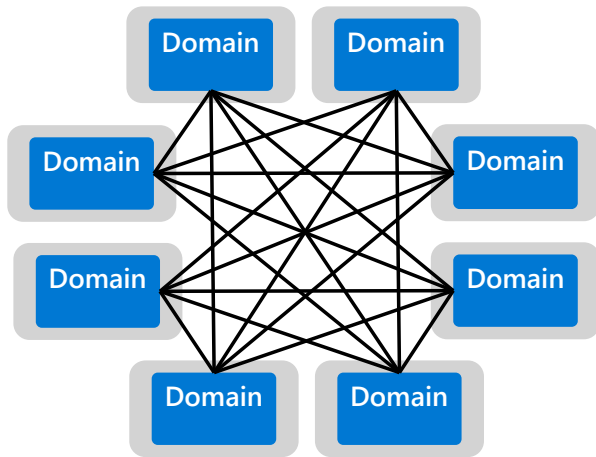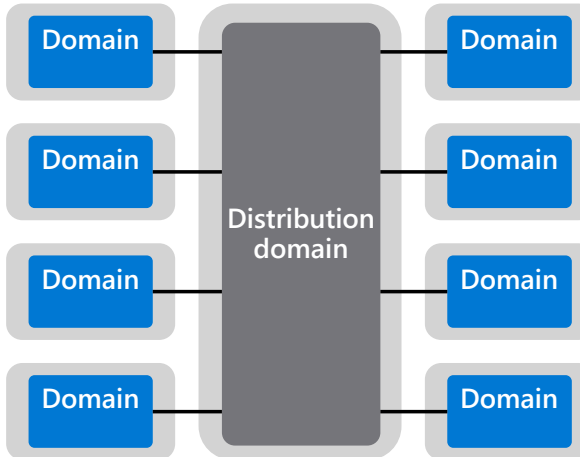Consumption -ready data

Orchestrate using Azure Data Factory

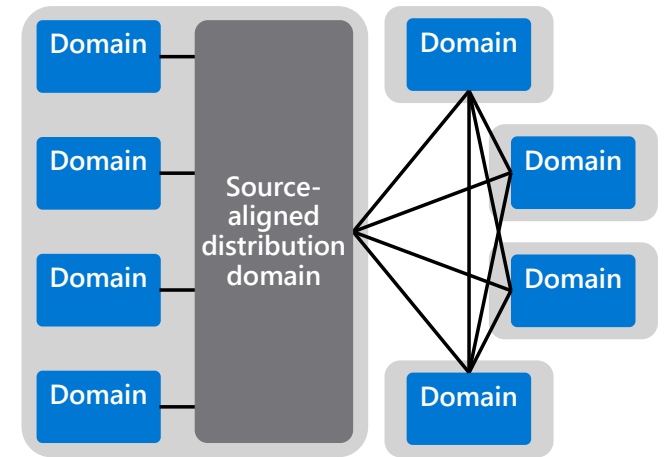Driving at Scale

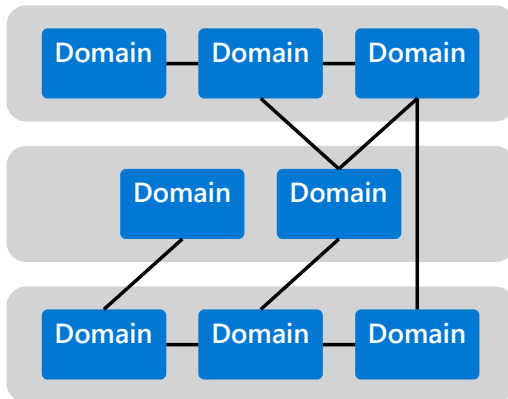= team independency

**Fine-grained fully federated mesh**
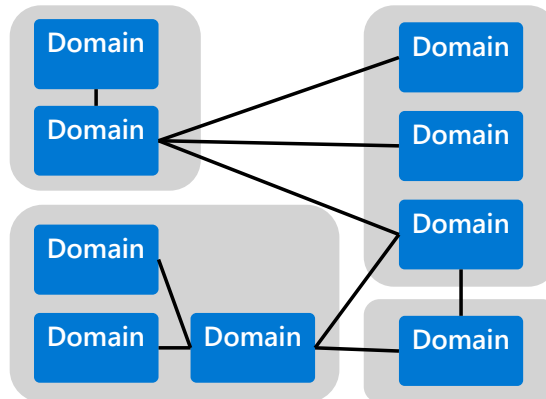
**Fine-grained and fully governed mesh**
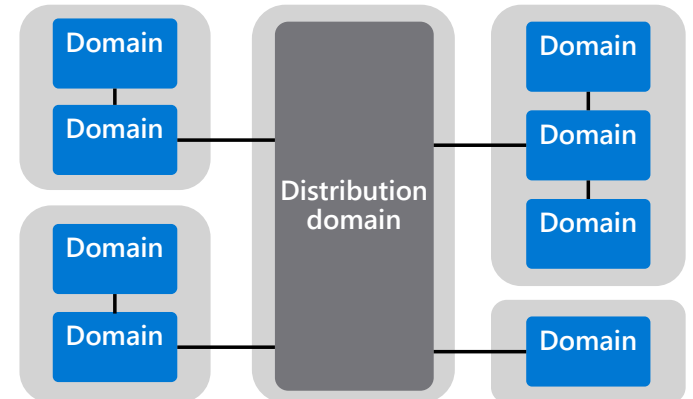
**Hybrid federated mesh**

**Value chain-aligned mesh**

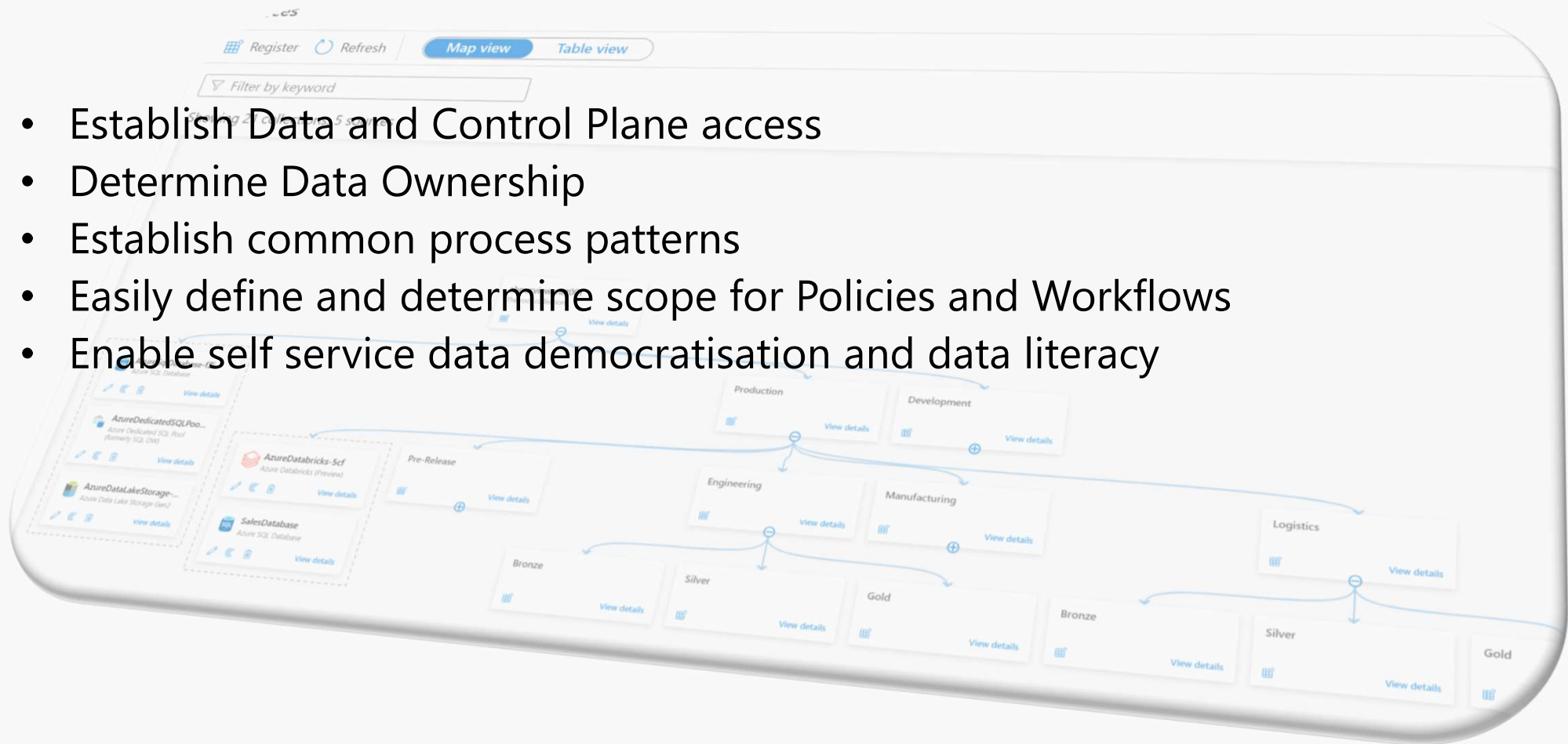**Coarse grained aligned mesh**
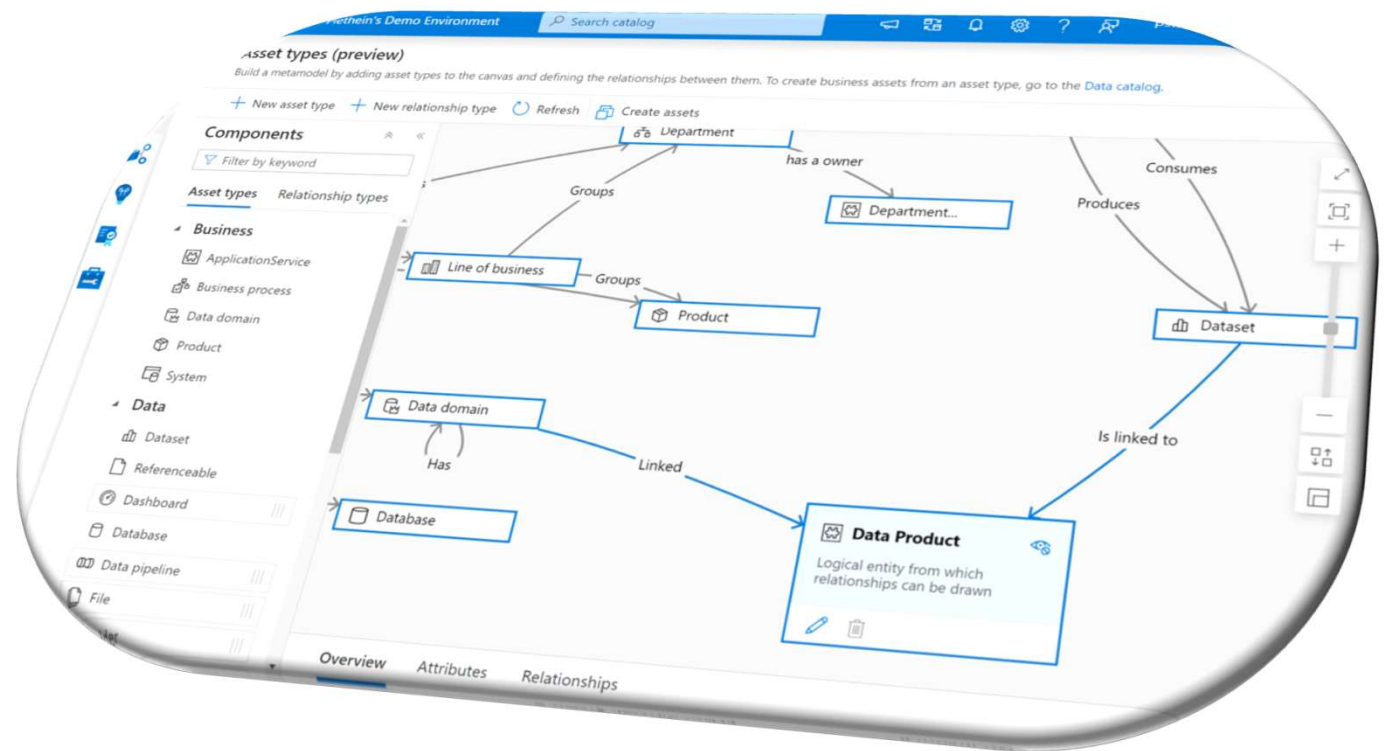
**Coarse grained and governed mesh**
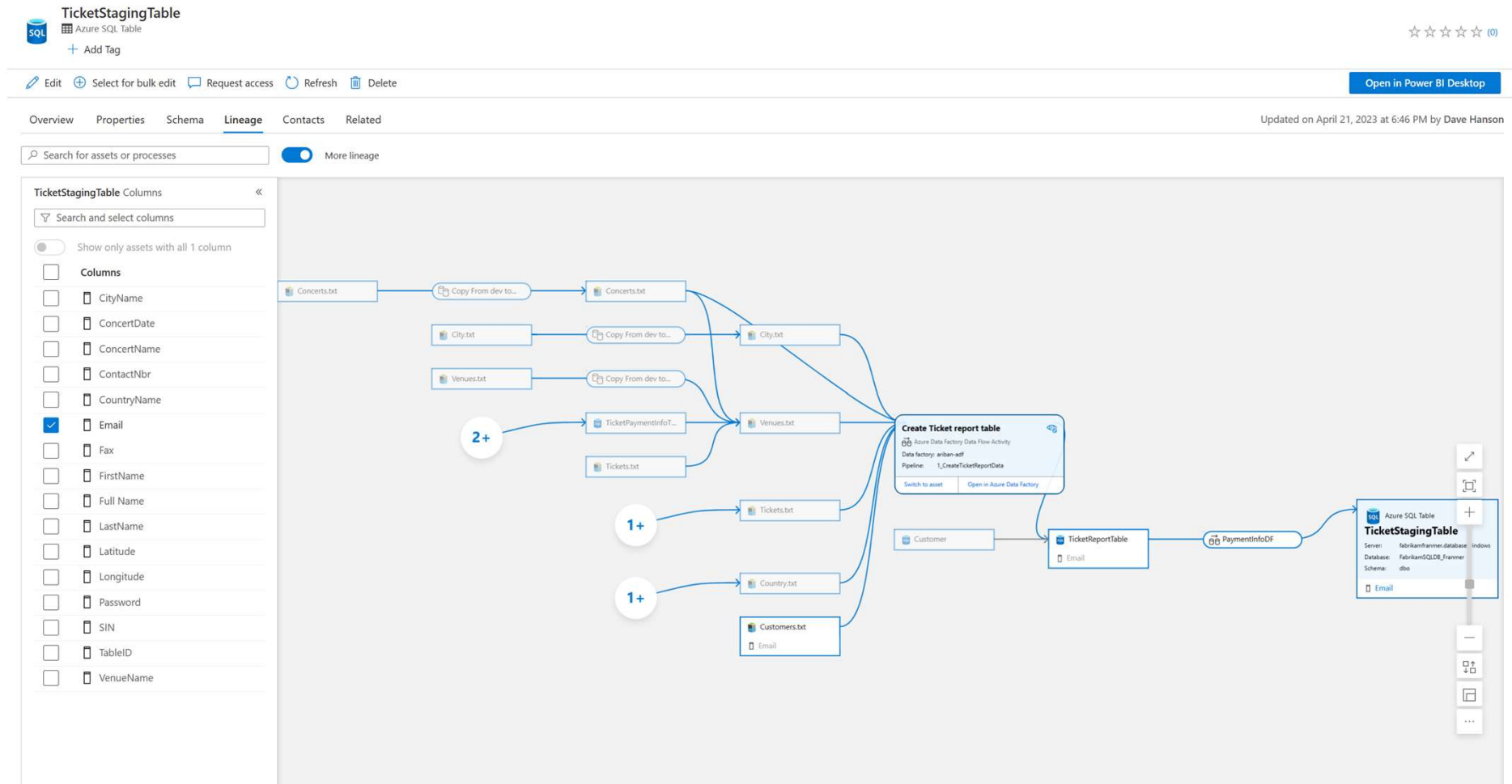
# Define the Lakehouse in Purview

- Establish Data and Control Plane access
- Determine Data Ownership
- Establish common process patterns
- Easily define and determine scope for Policies and Workflows
- Enable self service data democratisation and data literacy

# Encourage and practice good stewardship

# Use native capabilities to enable visibility

- The Lakehouse pattern is a scalable long term supportable concept

- Metadata driven data engineering provides clear data boundaries at scale

- Data engineering goes hand in hand with data governance

https://github.com/Andreas-bersgtedt/Metdata2Lakehouse

Key take-aways