Microsoft

# Modern Data Platform

## Straight Through File Processing

Azure Global Bootcamp 2019 - Neil Millington

# Capability Architecture

Casual Users
80%

Business Analysts
10%

Citizen Data Scientists
5%

Data Scientists
5%

## Delivery

| Natural Language | Data Visualisation | Data Mashup/Discovery | Data Science Tooling |

## Modern Data Platform

**Semantic Layer & Data Virtualization**

Abstract above data structures

Collate heterogenous data sources

Drill paths, metrics, KPIs

Caching & In Memory

**AI/Machine Learning**

Machine Learning

Deep Learning

Vision, Speech & Text

Deployment at Scale

**Metadata & Governance**

Catalog data objects

Define scope & provenance

Manage & implement Data Quality/Profiling (PaaS - tbc)

Data Lineage (PaaS - tbc)

**Real Time Analytics**

Analyse event streams/rolling time window

Anomaly detection

**Relational Analytics**

"Traditional" Data Warehousing

SQL on Relational data (star schemas, OLAP)

**Big Data Analytics**

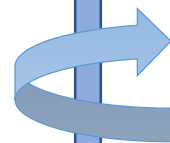Scale out processing

Semi & unstructured data

**Data Movement**

Orchestrate batch & real time data sets

Manage data transforms

**Data Lake**

| Curated Data | Land data in raw form |
| Cleansed Data | Cleanse & transform in higher level tiers |
| Raw Data | Present integrated, curated data and all layers below |

Microsoft

Data Sources

Apps

Sensors & Devices

# Modern Data Platform – wants & needs

**Casual Users** 80%  
**Business Analysts** 10%  
**Citizen Data Scientists** 5%  
**Data Scientists** 5%

## Delivery

| Chatbot Framework | Power BI Premium, Apps, Embedded | Power BI Desktop | ML Studio, DataBricks Notebooks |
|---|---|---|---|
| Natural Language | Data Visualisation | Data Mashup/Discovery | Data Science Tooling |

## Azure PaaS Data Services

### Semantic Layer & Data Virtualization
Analysis Services — Power BI

Collaborative reporting (B2B, B2C, self-service)

Unconstrained data sources - social, legacy, data lake. Flexible & rapid adoption and sharing

Reason over data, any format, any size

### AI/Machine Learning
Cognitive Services  
ML Containers & APIs  
ML Server

### Metadata & Governance
Data Catalog Gen-2

### Real Time Analytics
Stream/ Log Analytics — Storm — ADX — Azure DataBricks

### Relational Analytics
SQL DB — OSS - MySQL / PostgreSQL — SQL Data Warehouse — Snowflake — ADX — HDInsight — Azure DataBricks

### Big Data Analytics

Data Science as a Service

### Data Movement
Data Factory — Azure Databricks — Kafka — Event/IoT Hub

### Data Lake
Curated Data  
Cleansed Data  
Raw Data  
Azure Data Lake Store (V2) — Cosmos DB

Data as a Service (Data Lake, Data Pools)

Microsoft

Data Sources — Apps — Sensors & Devices

# On Demand Scalability…..

- Wide range of scale on demand services

- Sandpit capability, scale as needed

- Adapt to peaks and troughs as needed

| Service | Scale | Capability |
|---|---|---|
| Azure Analysis Services | Scale up: 400GB compressed Scale out: up to 8 instances | Pause & Resize |
| SQL DW | Scale out: 100 to 18000 DWUs | Pause & Resize |
| ADX / Kusto | Cluster creation and teardown | Pause & Resize |
| Databricks | Lightweight cluster creation and teardown | |
| Cosmos DB | Enterprise Scale NoSQL as a Service | |
| Azure Data Lake Store | Near infinite scale, no limit on object size | |
| Azure Data Factory | Data Movement as a Service | Up to 1GB/s |

Microsoft

# ...and availability



Casual Users 80%
Business Analysts 10%
Citizen Data Scientists 5%
Data Scientists 5%

**Delivery**

| Chatbot Framework | Power BI Premium, Apps, Embedded | Power BI Desktop | ML Studio, DataBricks Notebooks |
|---|---|---|---|
| Natural Language | Data Visualisation | Data Mashup/Discovery | Data Science Tooling |

**Azure PaaS Data Services**

Scalable, Available, Affordable

**Semantic Layer & Data Virtualization**

Analysis Service — 99.9% SLA

Power BI

Reduce operational burdens and constraints

**AI/Machine Learning**

Cognitive Services

ML Containers & APIs

**Metadata & Governance**

Data Catalog — 99.9% SLA

**Real Time Analytics**

Stream/ Log Analytics — 99.9% SLA
Storm — 99.9% SLA
ADX / Kusto — 99.95% SLA
Azure DataBricks

**Relational Analytics**

SQL DB — 99.9% SLA
OSS - MySQL / PostgreSQL — 99.99% SLA
SQL Data Warehouse — 99.99% SLA
Snowflake
ADX

**Big Data Analytics**

HDInsight — 99.9% SLA
Azure DataBricks — 99.9% SLA
ML Server

1st party Service, Enterprise Grade SLA

**Data Movement**

Data Factory — 99.9% SLA
Azure Databricks — 99% SLA
Kafka
Event/IoT Hub — 99.9% SLA

**Data Lake**

Curated Data
Cleansed Data
Raw Data

Azure Data Lake Store (V2) — 99.9% SLA
Cosmos DB

99.999% SLA (Geo Rep.) for availability, throughput, latency & consistency

Data Sources

Apps

Sensors & Devices

Microsoft

# Warehousing on Azure

**Data Factory**
- Data Movement as a Service
- Host SSIS packages in pipelines
- Visual Tooling

**Power BI**
- Data discovery for everyone
- Class leading Data Shaping

Casual Users 80%

Business Analysts 10%

Data Scientists 5%

Delivery — Chatbot Framework — Natural Language

Power BI Premium, Apps, Embedded — Data Visualisation

Power BI Desktop — Data Mashup/Discovery

ML Studio, Notebooks — Data Science Tooling

**Azure Analysis Services**
- Enterprise wide semantic layer
- Stunning response times
- Scale up/Scale out

**Data Catalog**
- Catalog Enterprise
- GDPR initiatives

Azure PaaS — Data Services

Semantic Layer & Data Virtualization

Analytics Services — Power BI

**Analytics**
- Scalable
- Lambda (Hot, Warm, Cool data)

**Data Warehouse**
- DW as a Service

AI/Machine Learning — Cognitive Services — ML Containers & APIs

Metadata & Governance — Data Catalog

**Real Time Analytics**
Stream/ Log Analytics | Storm | ADX / Kusto | Azure DataBricks

**Relational Analytics**
SQL DB | OSS - MySQL / PostgreSQL | SQL Data Warehouse | Snowflake

**Big Data Analytics**
ADX | HDInsight | Azure DataBricks

**HDInsight**
- Hadoop as a Service

ML Server

**Data Movement**
Data Factory | Azure Databricks | Kafka

**Azure DataBricks**
- 1st class service
- 1 click set up
- Integration with Azure

Curated Data

Cleansed Data

Raw Data

Azure Data Lake Store (V2) | Cosmos DB

**Data Lake Store**
- No limit on object size/storage
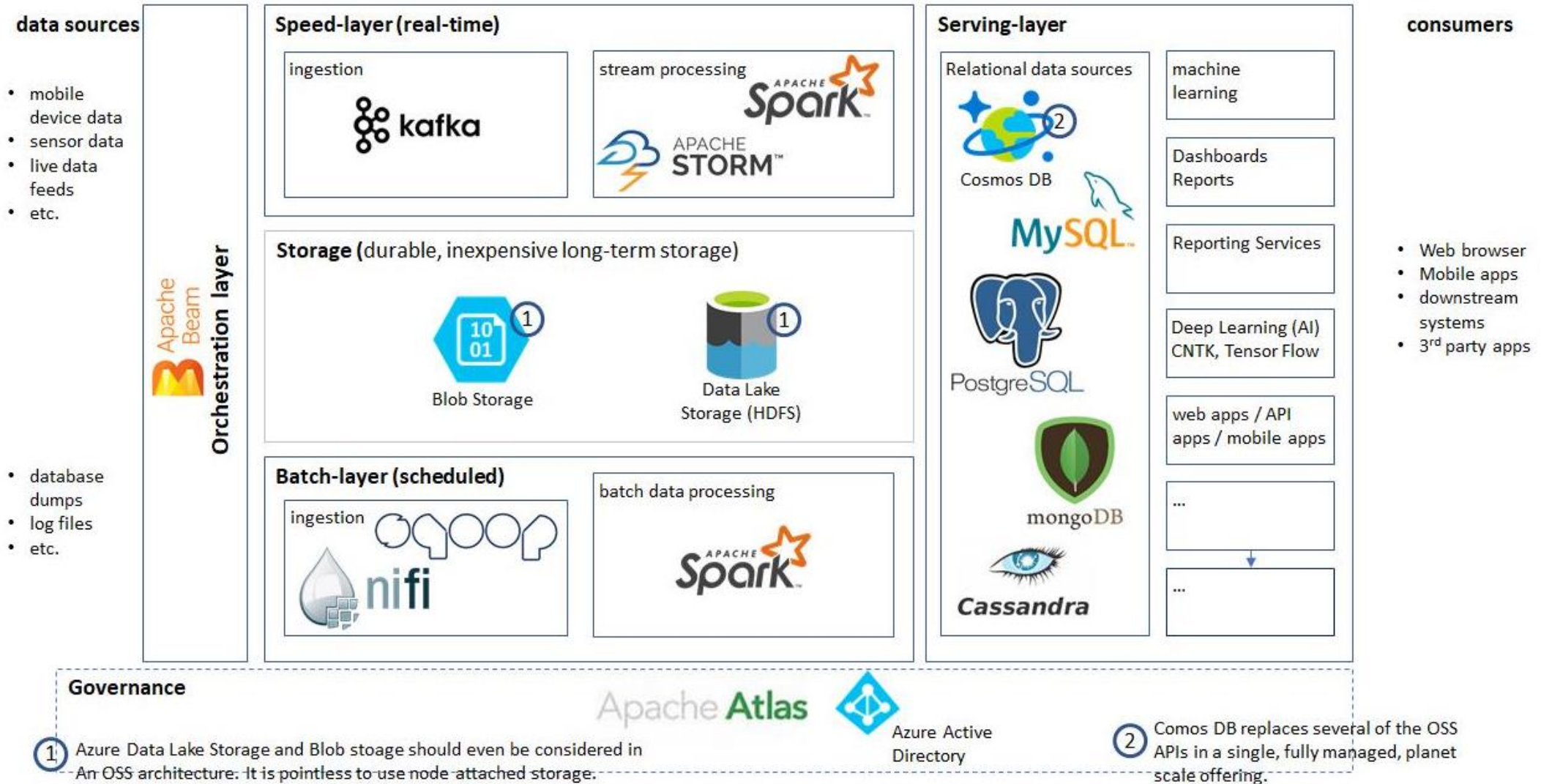- Optimised for parallel workloads at scale

**Cosmos DB**
- Mongo, Cassandra, Graph, Key value
- Global scale

Data Sources

Devices

# An OSS viewpoint

# Our job

**Plug n Play Products**

Microsoft

# Modern Data Platform (**commoditized patterns**)

| | INGEST | STORE | PREP | MODEL & SERVE (& store) | |
|---|---|---|---|---|---|
| LOB | | | | | |
| CRM | | | | | → BI + Reporting |
| Graph | | | | | |
| Image | Data orchestration and monitoring | Big data store | Transform & Clean | Data warehouse | → Advanced Analytics |
| Social | | | | | |
| IoT | | | | | → AI |
| Cloud | | | | | |

| INGEST | STORE | PREP | MODEL & SERVE |
|---|---|---|---|
| Azure Data Factory | Azure Data Lake Storage | Azure Databricks | Azure SQL Data Warehouse |
| SSIS | Blob Storage | Azure HDInsight | Azure Analysis Services |
| | SQL Server | PolyBase & Stored Procedures | SQL Database (Single, MI, HyperScale, Serverless) |
| | OSS (PostgreSQL, MySQL, MariaDB) | Power BI Dataflow | SQL Server in a VM |
| | | Azure Data Lake Analytics | Cosmos DB |
| | | | Power BI Aggregations |

# Azure Database Platform – Microsoft & OSS redefined

Relational Database Services

- Azure SQL Database
- Azure Database for MySQL
- Azure Database for PostgreSQL
- Azure Database for MariaDB
- SQL Server on VM

Non-Relational Database Services

- Azure Cosmos DB
- Azure Cache for Redis

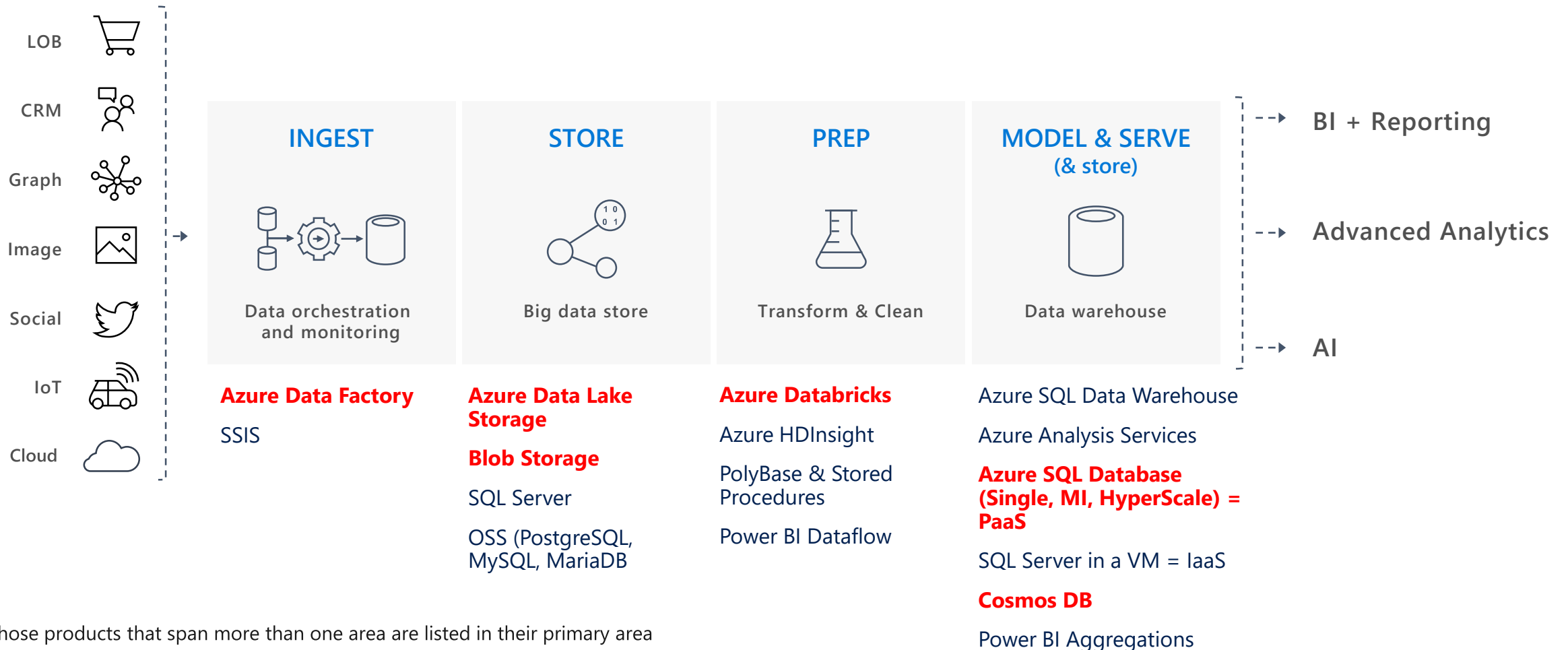| Fully-Managed | Flexible | Enterprise Scale & Performance | Security & Compliance |

Microsoft

# Products not covered

- Streaming (i.e. Iot Hub, Event Hub, Azure Stream Analytics)
- On-prem products (i.e. APS)
- Open source (i.e. Kafka, PostgreSQL, MySQL, MariaDB, Storm, Spark, Hbase, Redis)
- Machine learning/AI tools (i.e. Azure ML, Machine Learning Services, Cognitive Services)
- Reporting tools (i.e. Power BI)
- 3rd-party products (i.e. Informatica, Profisee)
- Competitor products (AWS, Google)
- Snowflake
- Azure Data Catalog (ADC) Gen2
- Azure Data Explorer
- Azure Database Migration Service
- Data Box, Data Box Disk
- Azure Search / Cognitive Search / Knowledge Mining
- Master Data Services (MDS)
- Azure Functions (for Prep) / Azure Logic Apps (for Ingest)
- OSS !!

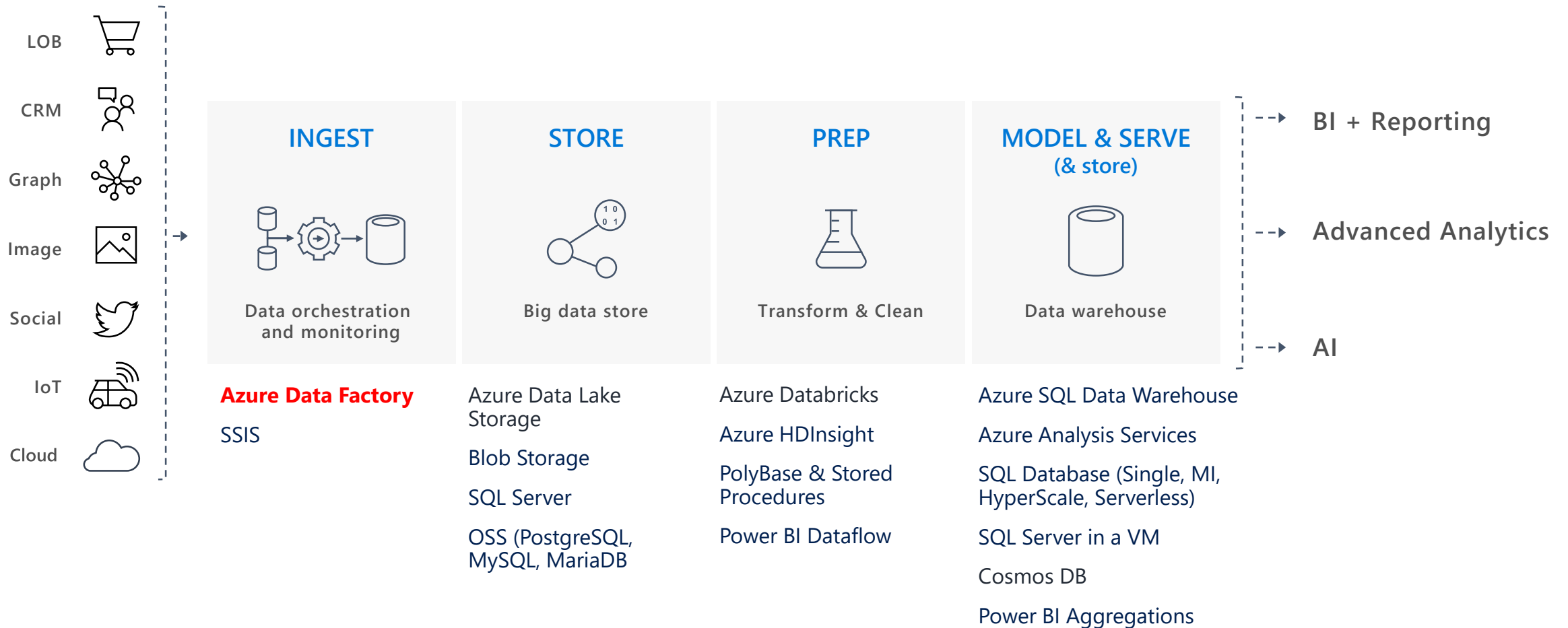Microsoft

# Modern Data Platform (possible products by four areas)

| | INGEST | STORE | PREP | MODEL & SERVE (& store) | |
|---|---|---|---|---|---|
| LOB | | | | | BI + Reporting |
| CRM | | | | | |
| Graph | | | | | Advanced Analytics |
| Image | Data orchestration and monitoring | Big data store | Transform & Clean | Data warehouse | |
| Social | | | | | AI |
| IoT | | | | | |
| Cloud | | | | | |

**INGEST**

**Azure Data Factory**

SSIS

**STORE**

**Azure Data Lake Storage**

**Blob Storage**

SQL Server

OSS (PostgreSQL, MySQL, MariaDB

**PREP**

**Azure Databricks**

Azure HDInsight

PolyBase & Stored Procedures

Power BI Dataflow

**MODEL & SERVE**

Azure SQL Data Warehouse

Azure Analysis Services

**Azure SQL Database (Single, MI, HyperScale) = PaaS**

SQL Server in a VM = IaaS

**Cosmos DB**

Power BI Aggregations

Note: Those products that span more than one area are listed in their primary area

Microsoft

# Modern Data Platform (commoditized patterns)



LOB
CRM
Graph
Image
Social
IoT
Cloud

| INGEST | STORE | PREP | MODEL & SERVE (& store) |
|---|---|---|---|
| Data orchestration and monitoring | Big data store | Transform & Clean | Data warehouse |

BI + Reporting

Advanced Analytics

AI

**Azure Data Factory**

SSIS

Azure Data Lake Storage

Blob Storage

SQL Server

OSS (PostgreSQL, MySQL, MariaDB

Azure Databricks

Azure HDInsight

PolyBase & Stored Procedures

Power BI Dataflow

Azure SQL Data Warehouse

Azure Analysis Services

SQL Database (Single, MI, HyperScale, Serverless)

SQL Server in a VM

Cosmos DB

Power BI Aggregations

Microsoft

# Ingest – Data Orchestration and Monitoring

**Product: Azure Data Factory (ADF)**

## Overview:

With Mapping Data Flow, can now transform data, so ETL/ELT tool. Copy Data tool to easily copy from source to destination.

## Use cases:

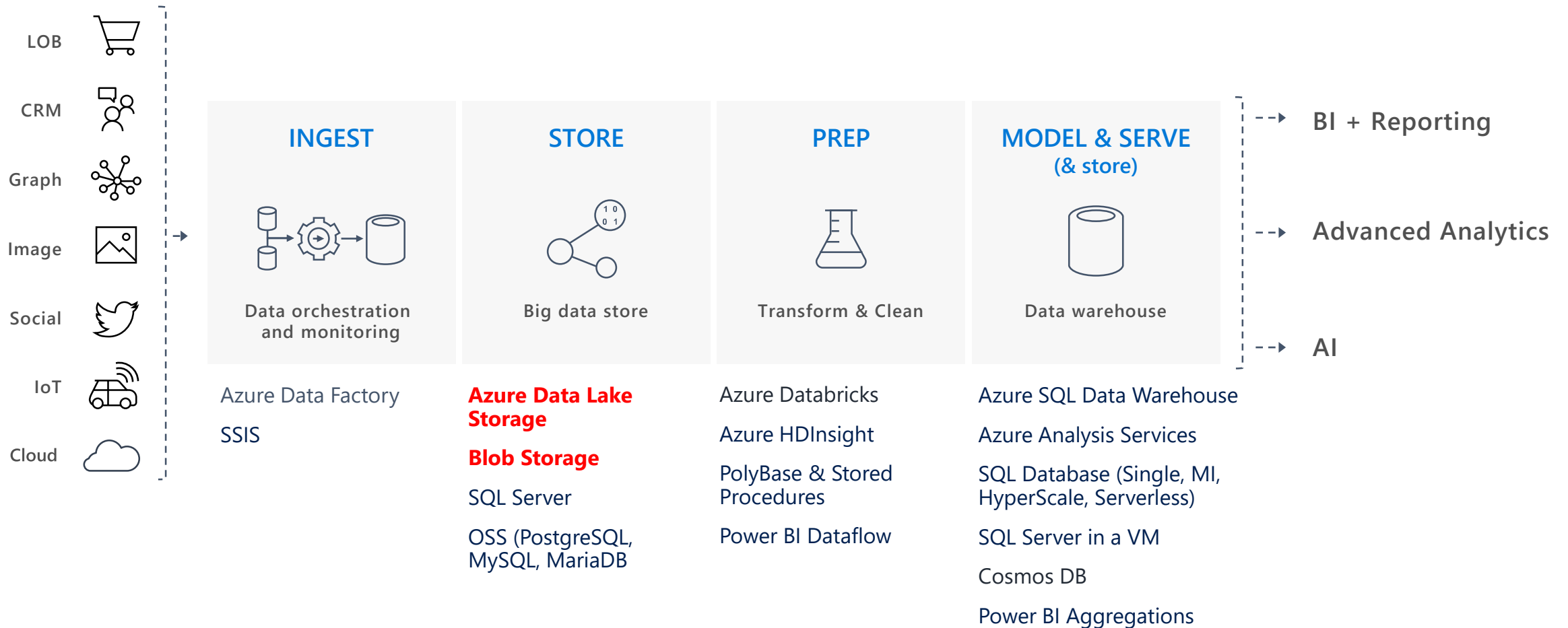Any new usecase for data movement and curation. SSIS packages migration.

## How to use:

PaaS based service – *Hybrid configuration, Hybrid operation*

## Area also used for:

Preparation of data (as part of a pipeline)

Microsoft

# Modern Data Platform (commoditized patterns)

| | INGEST | STORE | PREP | MODEL & SERVE (& store) | |
|---|---|---|---|---|---|
| LOB | | | | | → BI + Reporting |
| CRM | Data orchestration and monitoring | Big data store | Transform & Clean | Data warehouse | |
| Graph | | | | | → Advanced Analytics |
| Image | | | | | |
| Social | | | | | → AI |
| IoT | | | | | |
| Cloud | Azure Data Factory<br>SSIS | **Azure Data Lake Storage**<br>**Blob Storage**<br>SQL Server<br>OSS (PostgreSQL, MySQL, MariaDB | Azure Databricks<br>Azure HDInsight<br>PolyBase & Stored Procedures<br>Power BI Dataflow | Azure SQL Data Warehouse<br>Azure Analysis Services<br>SQL Database (Single, MI, HyperScale, Serverless)<br>SQL Server in a VM<br>Cosmos DB<br>Power BI Aggregations | |

Microsoft

# Store – Data Lake / Big Data

**Product:** Azure Data Lake Storage Gen2 (ADLS Gen2)

**Overview:**

Combines features of blob storage and ADLS Gen1. ADLS Gen2 adds a high performance HDFS Endpoint to Azure Blob Storage and inherits the rich feature set of Azure Blob Storage

**Use cases:**

Any new project.  Convert Blob and Gen1 over time

**How to use:**

PaaS

**Area also used for:**

n/a

# Store – Blob / Big Data

**Product: Blob Storage**

## Overview:

Original storage, foundational to Azure

## Use cases:

non-analytical use cases that only need object storage rather than hierarchical storage (i.e. video, images, backup files).

ADLS Gen2 – no need to use if current data does not need features of ADLS Gen2
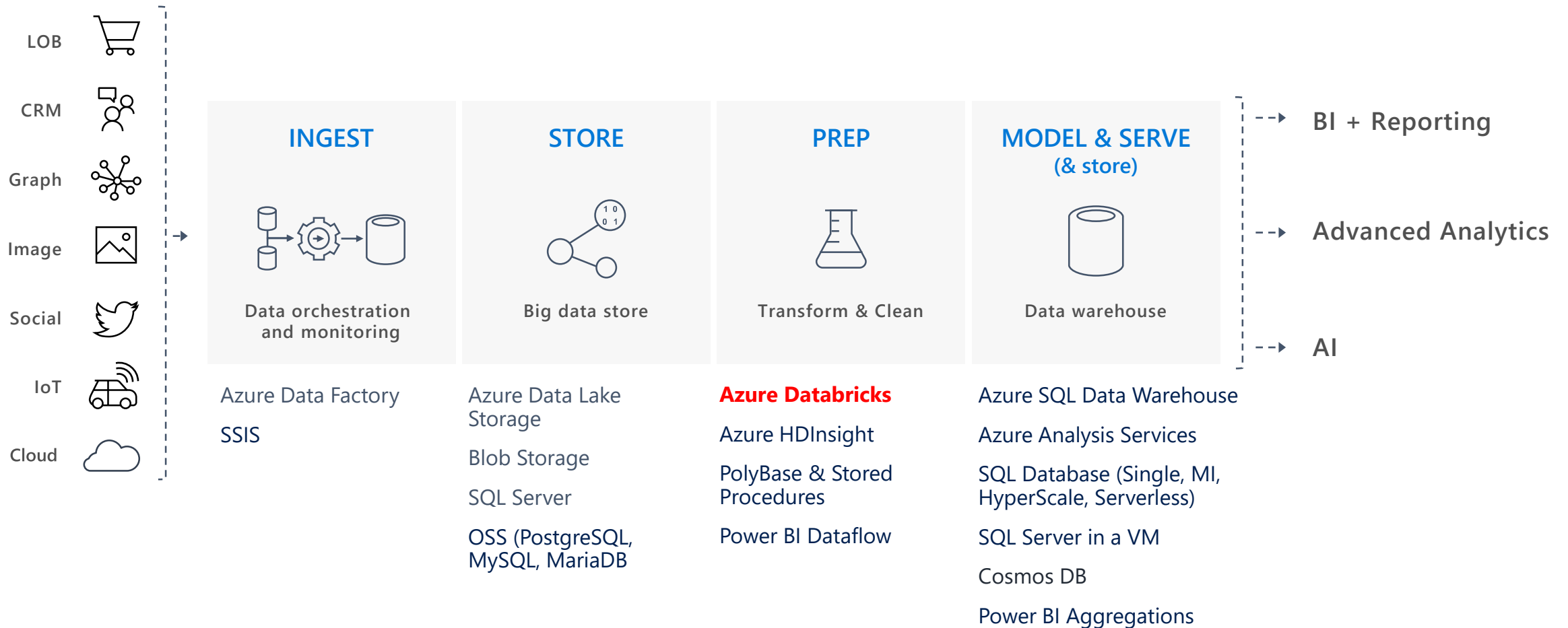
## How to use:

PaaS

## Area also used for:

n/a

Microsoft

# Modern Data Platform (commoditized patterns)

| | INGEST | STORE | PREP | MODEL & SERVE (& store) |
|---|---|---|---|---|
| LOB | | | | |
| CRM | | | | |
| Graph | | | | |
| Image | Data orchestration and monitoring | Big data store | Transform & Clean | Data warehouse |
| Social | | | | |
| IoT | | | | |
| Cloud | | | | |

BI + Reporting

Advanced Analytics

AI

**INGEST**
Azure Data Factory

SSIS

**STORE**
Azure Data Lake Storage

Blob Storage

SQL Server

OSS (PostgreSQL, MySQL, MariaDB

**PREP**
**Azure Databricks**

Azure HDInsight

PolyBase & Stored Procedures

Power BI Dataflow

**MODEL & SERVE**
Azure SQL Data Warehouse

Azure Analysis Services

SQL Database (Single, MI, HyperScale, Serverless)

SQL Server in a VM

Cosmos DB

Power BI Aggregations

Microsoft

# Prep – Big Data

**Product: Azure Databricks**

**Overview:**

Tool for curating and processing massive amounts of data and developing, training and deploying models on that data, and managing the whole workflow process throughout the project. SPARK.

**Use cases:**

Comfortable with Spark and notebooks, integration with ADLS, SQL DW, PBI, etc, need auto-scaling and auto-termination, need fast Spark
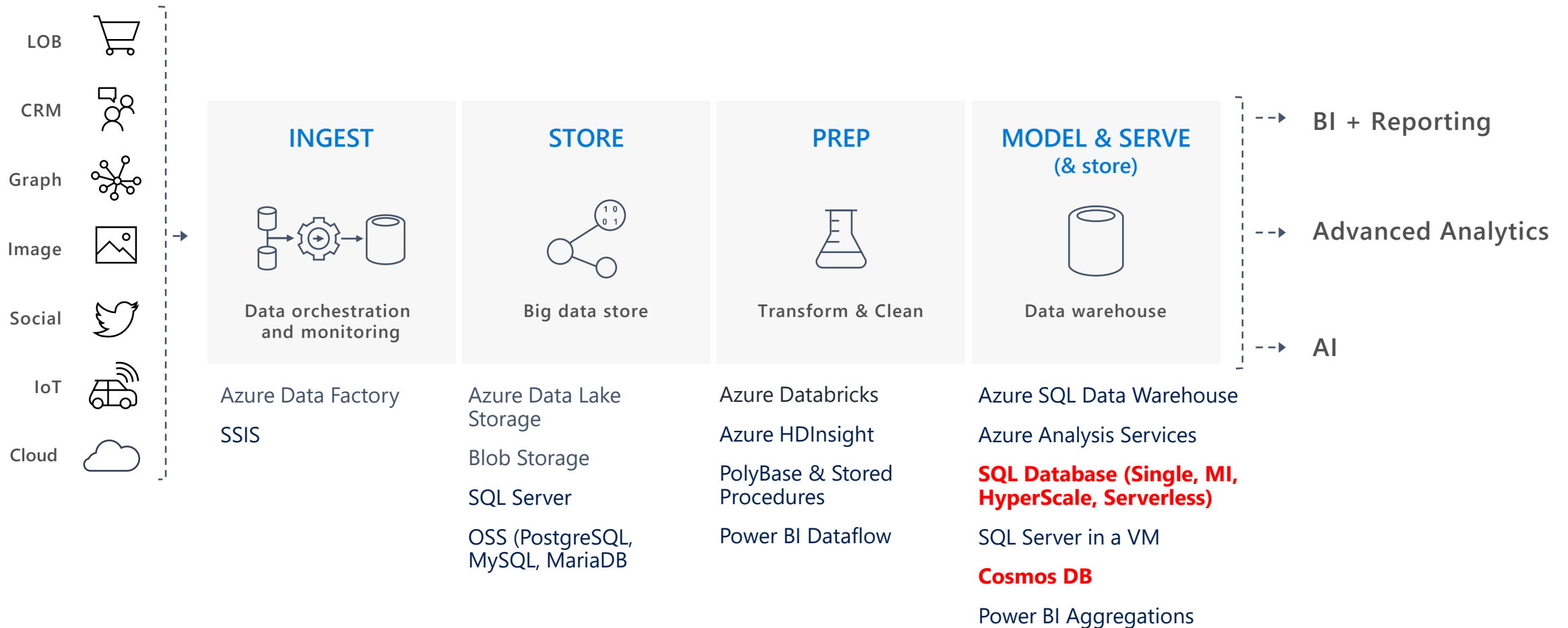
**How to use:**

PaaS

**Area also used for:**

ETL/ELT - Ingest, Model & Serve

Microsoft

# Modern Data Platform (commoditized patterns)

| | INGEST | STORE | PREP | MODEL & SERVE (& store) | |
|---|---|---|---|---|---|
| LOB | | | | | BI + Reporting |
| CRM | | | | | |
| Graph | | | | | |
| Image | Data orchestration and monitoring | Big data store | Transform & Clean | Data warehouse | Advanced Analytics |
| Social | | | | | |
| IoT | | | | | AI |
| Cloud | | | | | |

| INGEST | STORE | PREP | MODEL & SERVE |
|---|---|---|---|
| Azure Data Factory | Azure Data Lake Storage | Azure Databricks | Azure SQL Data Warehouse |
| SSIS | Blob Storage | Azure HDInsight | Azure Analysis Services |
| | SQL Server | PolyBase & Stored Procedures | **SQL Database (Single, MI, HyperScale, Serverless)** |
| | OSS (PostgreSQL, MySQL, MariaDB | Power BI Dataflow | SQL Server in a VM |
| | | | **Cosmos DB** |
| | | | Power BI Aggregations |

Microsoft

# Serve - Big Data

**Product: Cosmos DB**

## Overview:

A globally distributed, multi-model (key-value, graph, and document) database service.  It fits into the NoSQL camp by having a non-relational model (supporting schema-on-read and JSON documents)

## Use cases:

Works really well for large-scale OLTP solutions.  Spark to Cosmos DB connector for DW aggregations.  Use for data lake to have one datastore for both operational and analytical queries
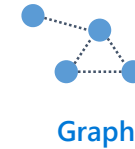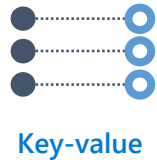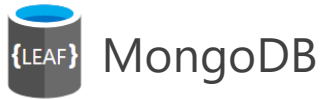
## How to use:

PaaS

## Area also used for:

Store, Prep

Microsoft

# Azure Cosmos DB

**A globally distributed, fully managed, massively scalable, multi-model database service**

SQL

MongoDB

Table API

Gremlin
$G = (V, E)$

cassandra

**Key-value**

**Column-family**

**Document**

**Graph**

Guaranteed low latency at the 99th percentile

Elastic scale out
of storage & throughput

Five well-defined consistency models

Turnkey global distribution

Comprehensive SLAs

Microsoft

# Ok, but why ?

# The Data Estate



**DATASOURCE**
- Structured/ Curated
- Logs (unstructured)
- Web/Social/Media (unstructured)
- Video / Voice / Image Files (unstructured)

**INGEST**
- AZURE DATA FACTORY
- 10 01 BLOB Landing-zone
- EVENT GRID

Generic Azure Services
- Azure Active Directory
- Azure Key Vault
- Azure Data Catalog

**STORE**

Azure Data Lake

structured - csv, txt
semi - Json, xml
semi – PDF, Office
unstructured – image, voice, video
Web/social – stream, Twitter

Enterprise / Divisional Lake

**PREP / TRAIN / MODEL / ANALYSE**

Azure Databricks
kubernetes
Analytics

Snowflake
SQL DW
Warehouse

AutomatedML Deep Learning
Keras
Flask
R
ONNX
Scala
Spark ML
python
Machine Learning / Model Execution

**SERVE**
- Business / custom apps
- Power BI
  Q & A (English Query)
  ML

Microsoft

# Data as a Service - Theoretical

| DATASOURCE | INGEST | STORE (Lake) | PREP / TRAIN / MODEL / ANALYSE | SERVE |
|---|---|---|---|---|

**DATASOURCE**

Structured/ Curated

Logs (unstructured)

Web/Social/Media (unstructured)

Video / Voice / Image Files (unstructured)

**INGEST**

AZURE DATA FACTORY

**10 01**
BLOB Landing-zone

EVENT GRID

Azure Active Directory

Azure Key Vault

Azure Data Catalog

Generic Azure Services

**STORE (Lake)**

Azure Data Lake

structured - csv, txt
semi - Json, xml
semi – PDF, Office
unstructured – image, voice, video
Web/social – stream, Twitter

Enterprise / Divisional Lake

**PREP / TRAIN / MODEL / ANALYSE**

Azure Databricks

kubernetes

Analytics

Snowflake

SQL DW

Warehouse

Keras

Flask

R

AutomatedML
Deep Learning

ONNX

Scala

Spark ML

python

Machine Learning / Model Execution

**SERVE**

Business / custom apps

Power BI

Q & A (English Query)

ML

Microsoft

# Data as a Service - Reality

| DATASOURCE | INGEST | STORE (Pools) | PREP / TRAIN / MODEL / ANALYSE | SERVE |
|---|---|---|---|---|

**DATASOURCE**

Structured/ Curated

Logs (unstructured)

Web/Social/Media (unstructured)

Video / Voice / Image Files (unstructured)

**INGEST**

AZURE DATA FACTORY

BLOB Landing-zone

EVENT GRID

Azure Active Directory

Azure Key Vault

Azure Data Catalog

Generic Azure Services

**STORE (Pools)**

SQL

Relational/Structured (SQL)
SQL Server, PostgreSQL, MySQL

Semi-relational (NoSQL)
COSMOS DB (MongoDB, Cassandra)

unstructured
(Data Lake)
PDF, Voice, Image, etc.

Enterprise / Divisional Lake

**PREP / TRAIN / MODEL / ANALYSE**

Azure Databricks

kubernetes

Analytics

Snowflake

SQL DW

Warehouse

K Keras

Flask

R

AutomatedML
Deep Learning

ONNX

Scala

Spark ML

python

Machine Learning / Model Execution

**SERVE**

Business / custom apps

Power BI

Q & A (English Query)

ML

Microsoft

# Data Science / Analytics as a Service

| DATASOURCE | INGEST | STORE (Lake) | PREP / TRAIN / MODEL / ANALYSE | SERVE |
|---|---|---|---|---|

**DATASOURCE**

Structured/ Curated

Logs (unstructured)

Web/Social/Media (unstructured)

Video / Voice / Image Files (unstructured)

**INGEST**

AZURE DATA FACTORY

10 01

BLOB Landing-zone

EVENT GRID

Azure Active Directory

Azure Key Vault

Azure Data Catalog

Generic Azure Services

**STORE (Lake)**

Azure Data Lake

structured - csv, txt
semi - Json, xml
semi – PDF, Office
unstructured – image, voice, video
Web/social – stream, Twitter

Enterprise / Divisional Lake

**PREP / TRAIN / MODEL / ANALYSE**

Azure Databricks

kubernetes

Analytics

Snowflake

SQL DW

Warehouse

AutomatedML
Deep Learning

Keras

Flask

R

ONNX

Scala

Spark ML

python

Machine Learning / Model Execution

**SERVE**

Business / custom apps

Power BI

Q & A (English Query)

ML

Microsoft

# Demos: Straight Through Processing
# (the real world)

Microsoft

# Real World – Pattern 1

# Real World – Pattern 2



**SOURCE**

Logs (unstructured)

Structured/ Curated

Files (unstructured)

**INGEST**

AZURE DATA FACTORY

JSON

BLOB Landing-zone

**STORE**

SQL

Relational/Structured (SQL)
SQL Server, PostgreSQL, MySQL

JSON

AZURE DATA LAKE STORE V2

**PREP & TRAIN**

SQL

JSON

Azure Databricks

JSON

**MODEL**

**SERVE**

Custom apps / Business Objects

Azure Active Directory

Azure Key Vault

Azure Data Catalog

Generic Azure Services

Microsoft

# Real World – Pattern 3

# Appendices

**Use-case: data**

Microsoft

**Ability to query across multiple entity types with a single network request.**

**For example, we have two types of documents: cat and person.**

```
{
    "id": "Andrew",
    "type": "Person",
    "familyId": "Liu",
    "worksOn": "Azure Cosmos DB"

}
```

```
{
    "id": "Ralph",
    "type": "Cat",
    "familyId": "Liu",
    "fur": {
        "length": "short",
        "color": "brown"
    }
}
```
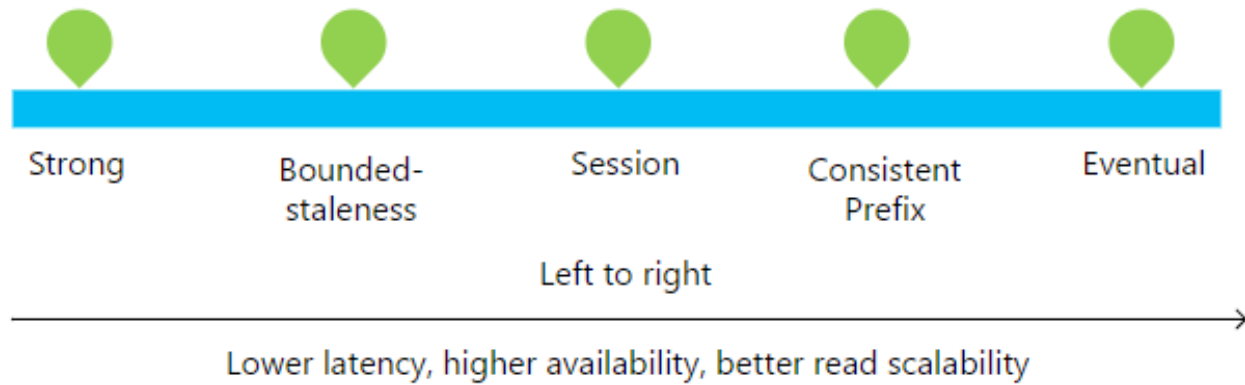
Microsoft

# Ability to query across multiple entity types with a single network request.

## For example, we have two types of documents: cat and person.

```
{
    "id": "Andrew",
    "type": "Person",
    "familyId": "Liu",
    "worksOn": "Azure Cosmos DB"
}
```

```
{
    "id": "Ralph",
    "type": "Cat",
    "familyId": "Liu",
    "fur": {
        "length": "short",
        "color": "brown"
    }
}
```

## We can query both types of documents without needing a JOIN simply by running a query without a filter on type:

SELECT * FROM c WHERE c.familyId = "Liu"

Microsoft

**Ability to query across multiple entity types with a single network request.**

**For example, we have two types of documents: cat and person.**

```
{
    "id": "Andrew",
    "type": "Person",
    "familyId": "Liu",
    "worksOn": "Azure Cosmos DB"
}
```

```
{
    "id": "Ralph",
    "type": "Cat",
    "familyId": "Liu",
    "fur": {
            "length": "short",
            "color": "brown"
    }
}
```

**If we wanted to filter on type = "Person", we can simply add a filter on type to our query:**

SELECT * FROM c WHERE c.familyId = "Liu" **AND c.type = "Person"**

Microsoft

**Use-case: application & services**

Left to right

Lower latency, higher availability, better read scalability

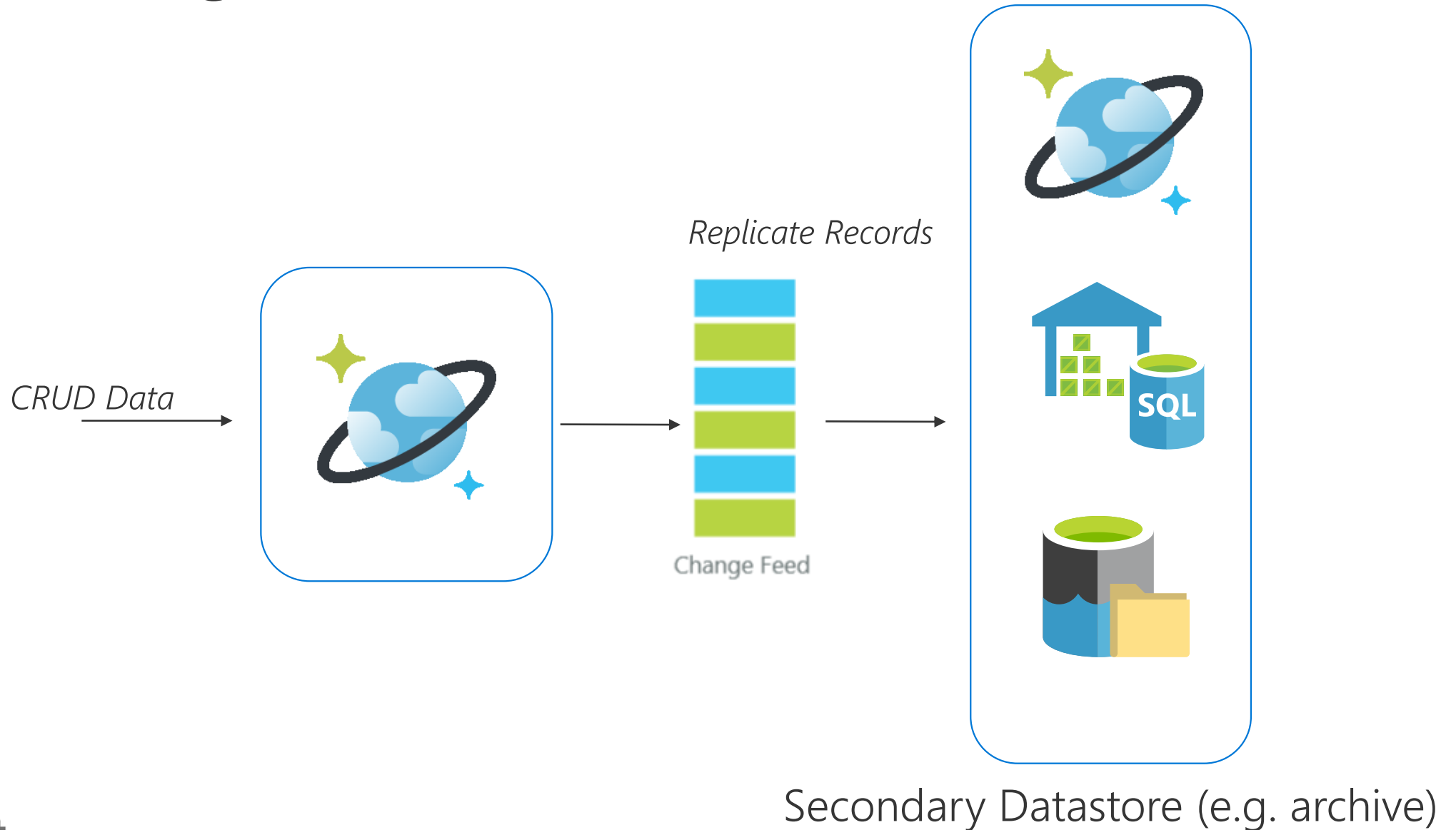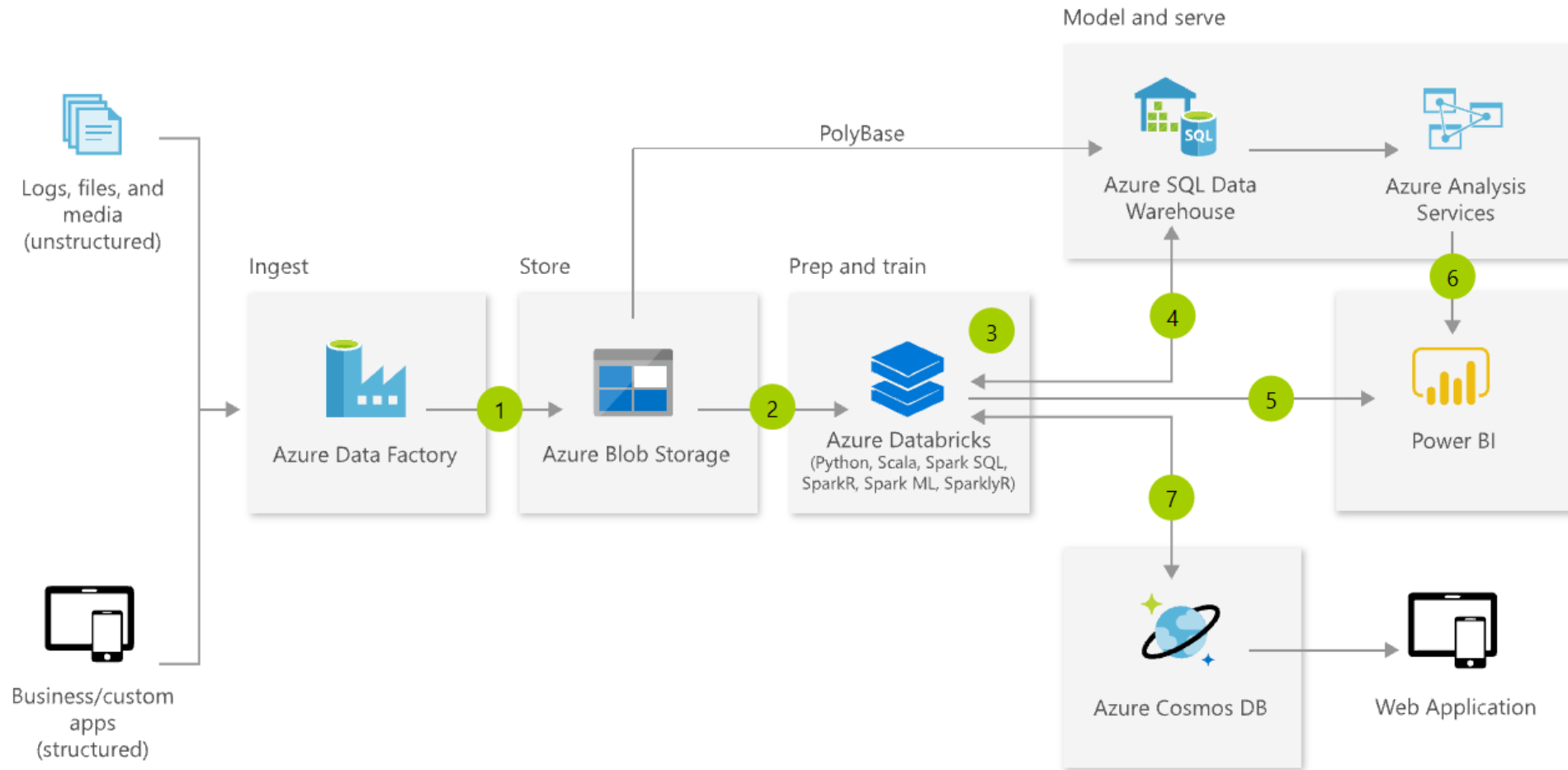| Consistency Level | Guarantees |
|---|---|
| Strong | Linearizability (once operation is complete, it will be visible to all) |
| Bounded Staleness | Consistent Prefix.<br>Reads lag behind writes by at most k prefixes or t interval<br>Similar properties to strong consistency (except within staleness window), while preserving 99.99% availability and low latency. |
| Session | Consistent Prefix.<br>Within a session: monotonic reads, monotonic writes, read-your-writes, write-follows-reads<br>Predictable consistency for a session, high read throughput + low latency |
| Consistent Prefix | Reads will never see out of order writes (no gaps). |
| Eventual | Potential for out of order reads. Lowest cost for reads of all consistency levels. |

# Event Sourcing for Microservices

**sqlbits**

New Event

Persistent
Event Store

Change Feed

Trigger Action
From Change Feed

Microservice
#1

Microservice
#2

...

Microservice
#N

Microsoft

# Materializing Views

# Replicating Data



CRUD Data

*Replicate Records*

Change Feed

Secondary Datastore (e.g. archive)

# Common architectures

# Modern DW for big data

# Modern DW for big data continued



Advanced analytics on big data

# Real-time analytics

# Real-time analytics continued

## Real time analytics

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

Sensors and IoT
(unstructured)

Apache Kafka for
HDInsight

Azure Databricks

Spark Streaming

Cosmos DB

Real-time apps

Media
(unstructured)

Logs (unstructured)

Files
(unstructured)

Business/custom apps
(structured)

Azure Data Factory

Azure Data Lake Storage

PolyBase

Azure SQL Data
Warehouse

Power BI

Azure Analysis
Services

Microsoft