

How CDC Integrated Complex Data to Drive Vaccination Forecasting with Databricks

Sheila Stewart, Senior Solutions Architect, Databricks
sheila.stewart@databricks.com

John Repko, Technical Program Manager, Peraton
John.Repko@Peraton.com

Jim Feters, Cloud Solution Specialist, Microsoft Azure
Jamesfe@Microsoft.com

Agenda

- Introductions
- COVID Big Data Use Case
- Benefits of Cloud Services
- Benefits of a Data Lakehouse

 Peraton



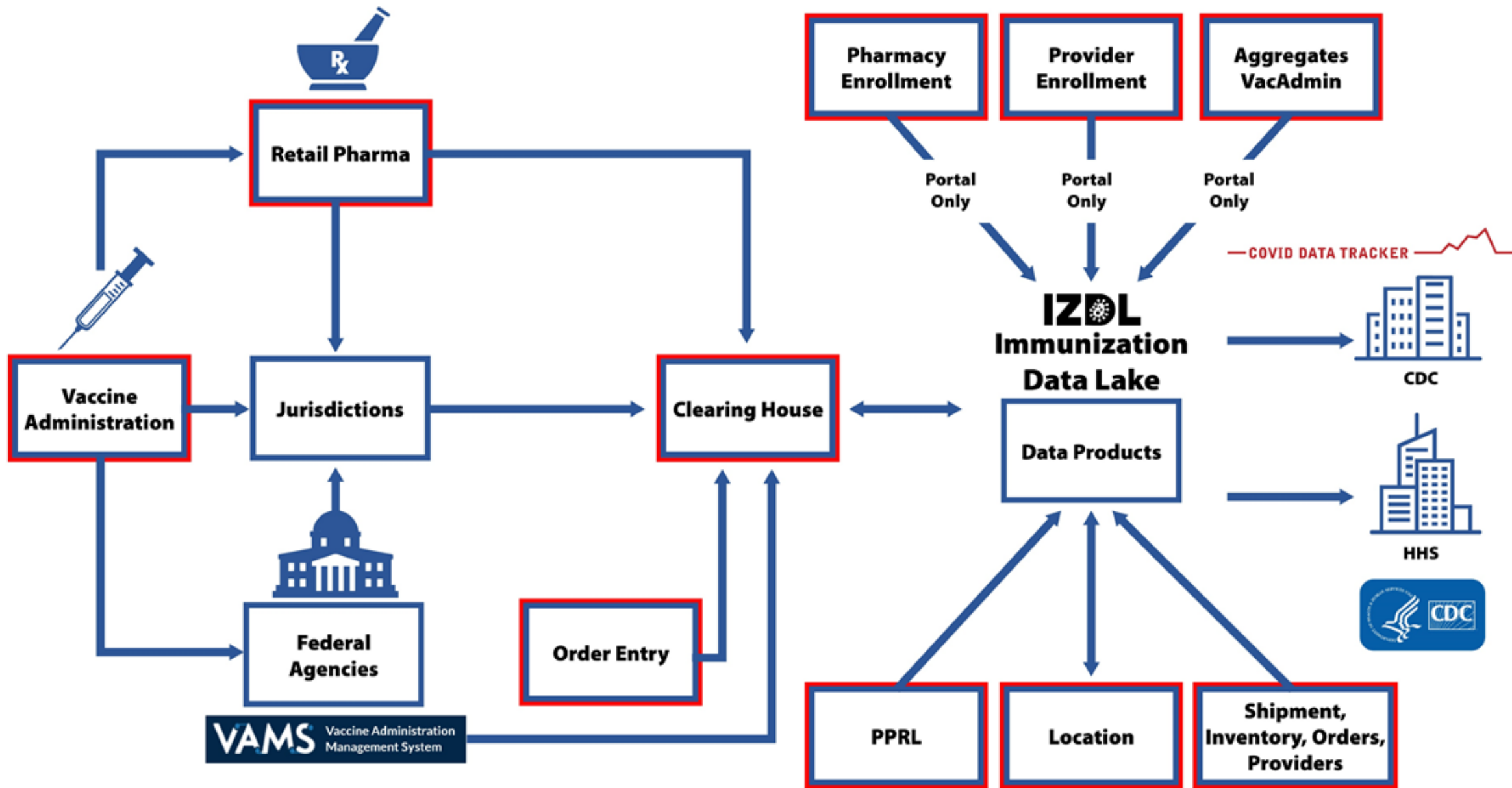
databricks



Microsoft

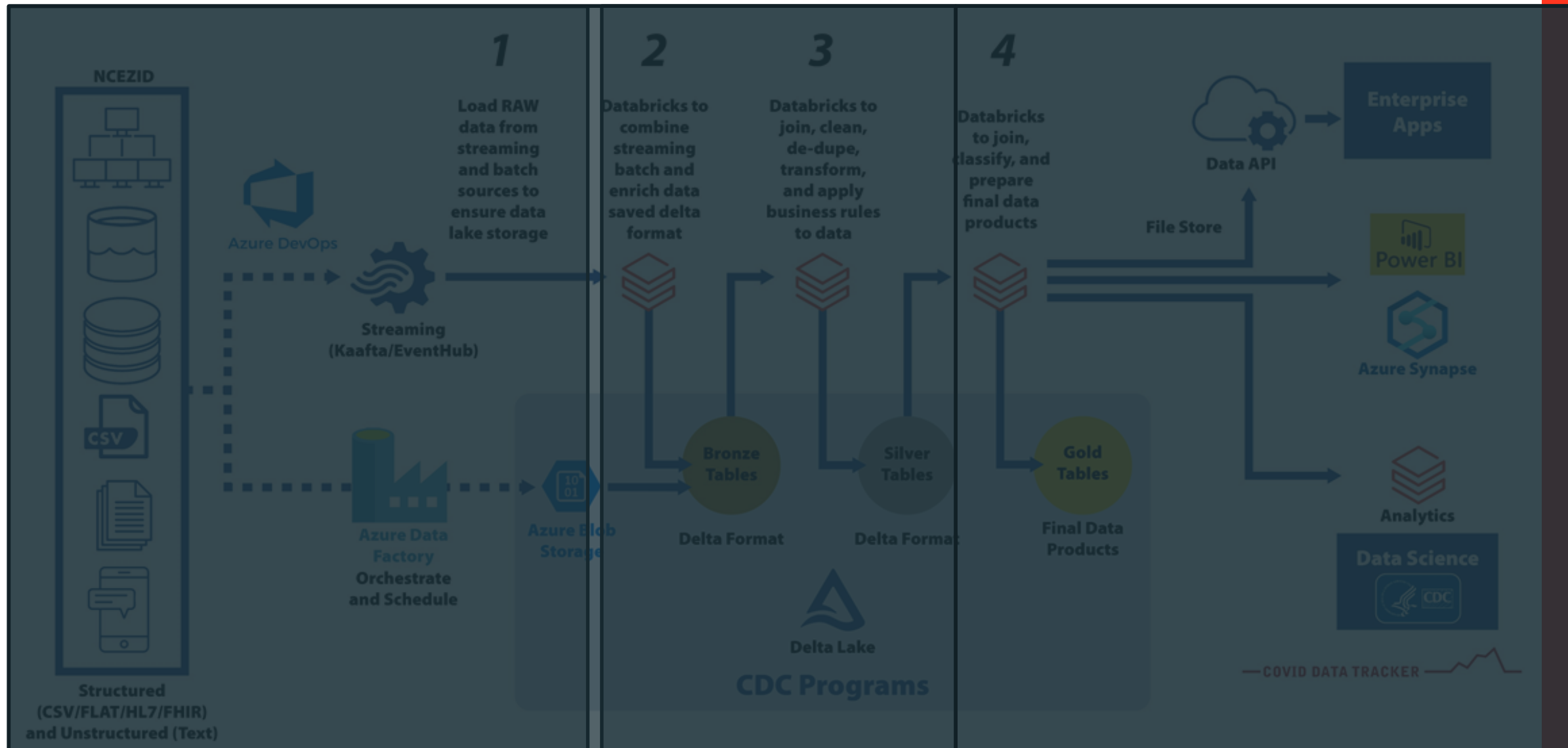
COVID Big Data Use Case

Use Case





Unified Delta Lake



Success

- **IZDL- stood up in 6 weeks from build out to Operational at CDC**
 - 8 Inbound Data Streams, 10 Outbound Data Streams and 45 Data Products**
- **Limitless Spark Processing Power**
 - At any scale for Streaming and Batch**
 - Volume – 5M+ new records per day, Billions Analyzed per day**
- **Realtime Data Streaming Pipelines**
 - Velocity – near Realtime – Hourly- Daily**
 - For very Large-Scale high throughput Parallel Processing (100K per sec)**
 - Full provenance (raw to aggregate)**
 - “What If” scenarios (Time Travel)**
 - Update/Deletes as Public Health Events**

Success (cont)

- **Horizontal scaling for storage**

Blob, Synapse -> Petabytes, Parquet (blazingly fast consumption)

- **Delivering Data to Public Health and multiple CDC Centers at Scale**

COVID Data Tracker & Data.CDC.Gov

Bulk Exports -> Terabyte Slices of Data

Data APIs (Synch and Asynch)

IZDL, CELR, DAART, EZDX, PHLIP, – In production or nearly so

Benefits of Cloud Services

Azure capabilities for public health



Data Interoperability

Standardizing, structuring and persisting data according to leading health standards (HL7 FHIR, DICOM)



Data governance

Discovering, protecting, and controlling data



Data modernization

Establishing future-proof platform that connects all data



Cloud native applications

Supporting high performance apps of any scale



Analytics and insights

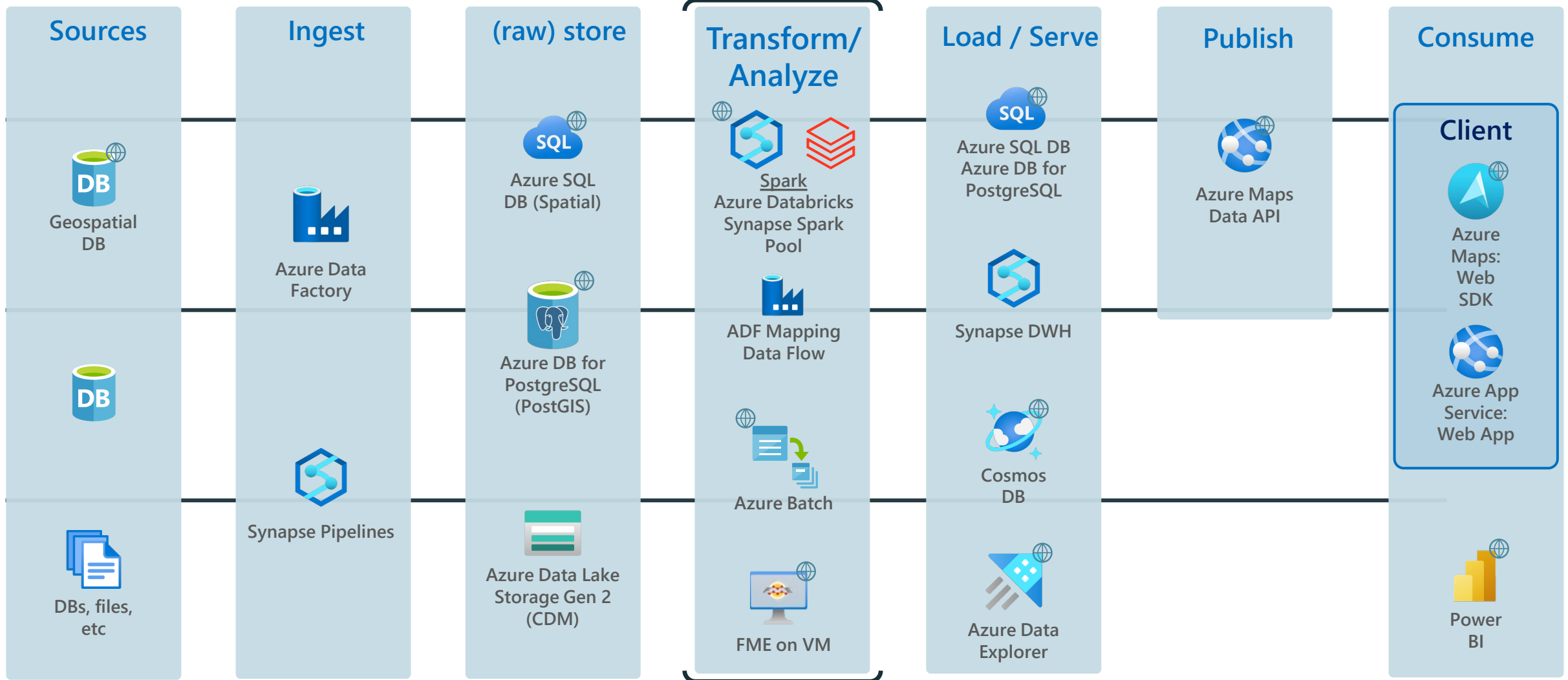
Understanding customer behaviours to generate insights



Data science

Powering AI experiences

Azure Cloud Services



Azure Data Factory
Synapse Pipelines

Transform/ Analyze

Spark
Azure Databricks
Synapse Spark
Pool



ADF Mapping
Data Flow



Azure Batch



FME on VM

Load / Serve



Azure SQL DB
Azure DB for
PostgreSQL



Synapse DWH



Cosmos
DB



Azure Data
Explorer

Client

- Azure Maps: Web SDK
- Azure App Service: Web App

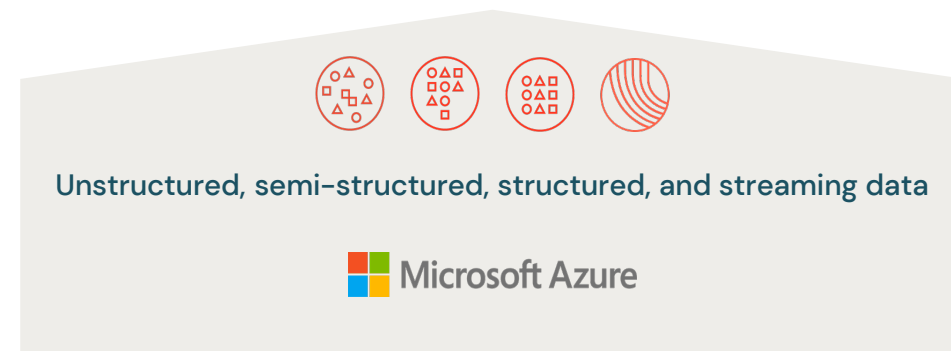
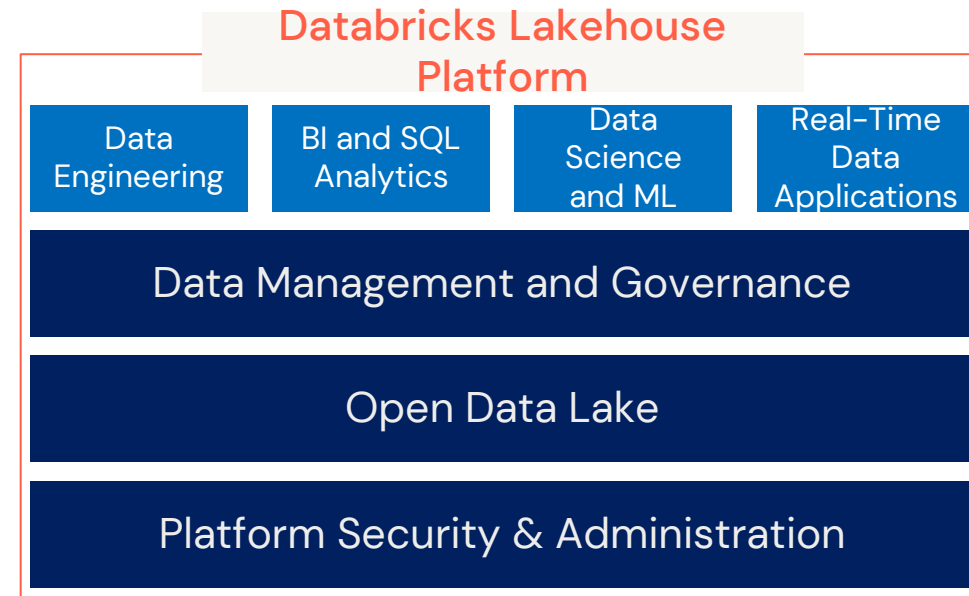
Benefits of a Data Lakehouse

The Databricks Lakehouse Platform

- ✓ Simple
- ✓ Open
- ✓ Collaborative

What is a Lakehouse?

- One simple platform to unify all of your data, analytics, and AI workloads by combining the best features of data warehouses and data lakes.



Unstructured, semi-structured, structured, and streaming data

 Microsoft Azure



Takeaways

Close teamwork between Databricks, Peraton, and Azure produced a very successful, quick-turnaround, reproducible cloud data management strategy

- Successful Forecasting is based on Rich, Complete, and Curated Data
- Whether historical data or real time collection is used for forecasting, data management needs planning and implementation
- It always takes a team of different technologies, processes, and dedicated people to be successful
- **THANK YOU to for the support from the CDC, Peraton, Microsoft Azure, and Databricks**

Questions?

Please come visit us at our booths!