



brightly

Rewriting tomorrow.

Brightly introduction

Ilpo Rouhola
+358 50 487 2920
ilpo.rouhola@brightlyworks.com

Future intelligent solutions,
delivered today.

Brightly in a Nutshell

Trusted strategic digital solutions partner that thinks brighter

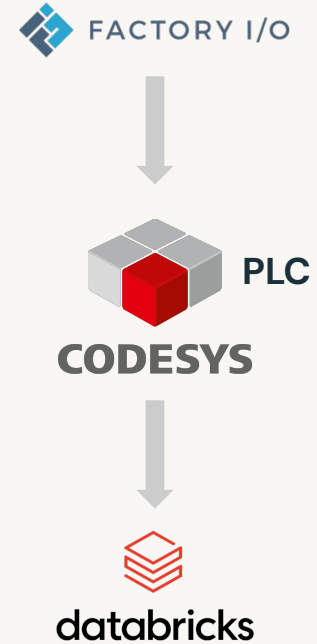




SI/Integration Partner IoT at Scale – Demo



1. Cloud Agnostic Approach

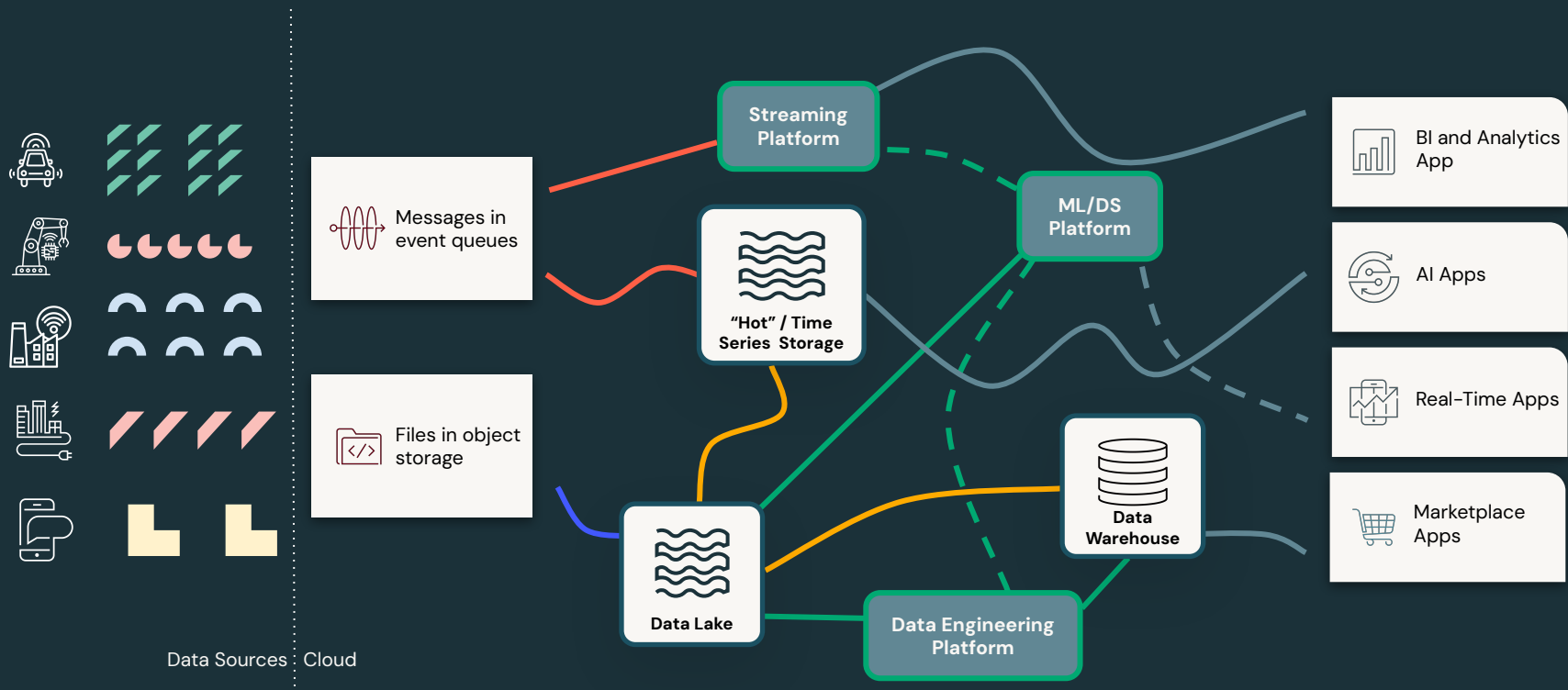


Industrial stream data management in Databricks - Jussa Klapuri

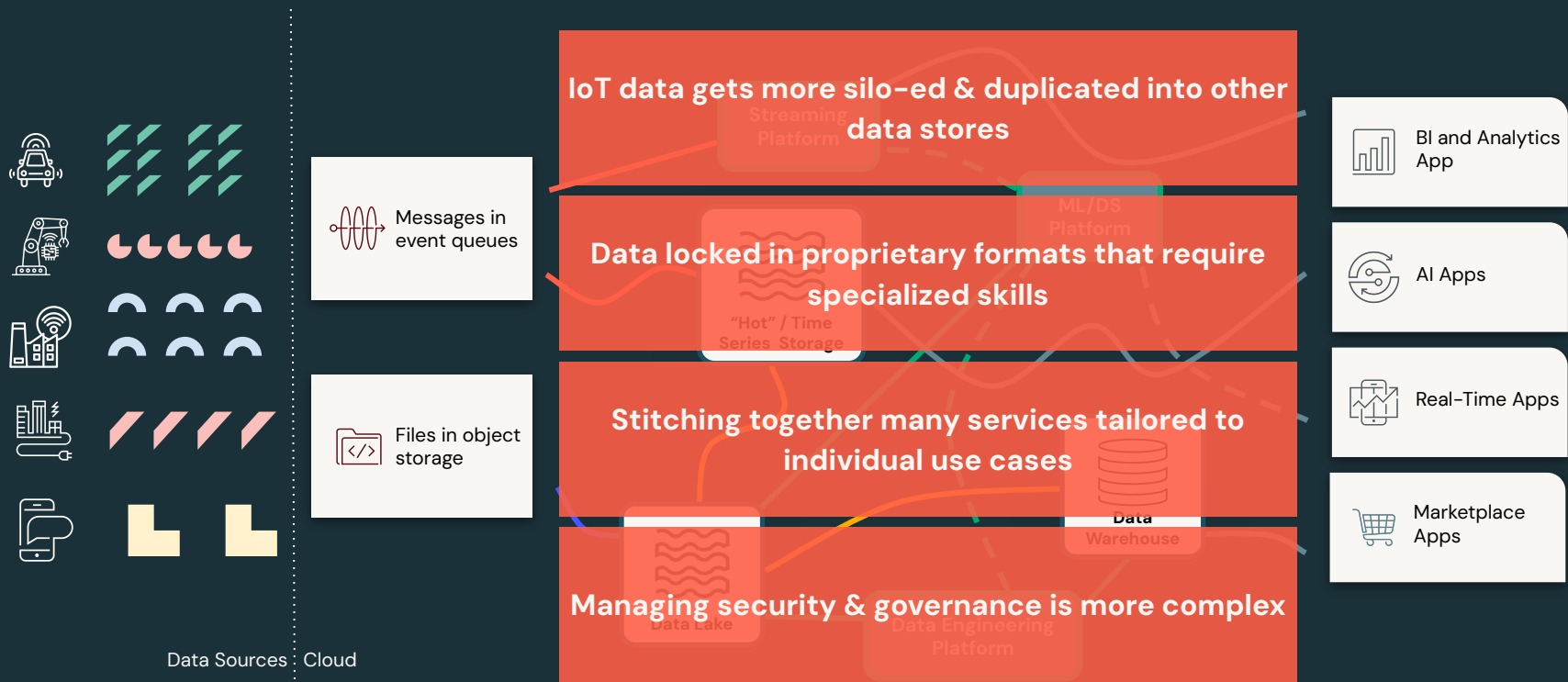
- Challenges with IOT Data
- Databricks and Lakehouse architecture
- Building Blocks
 - Spark, Delta Lake, Delta Tables
- Medallion architecture
 - Bronze, Silver, Gold
- Delta Live Tables & Data pipelines
 - Stream
 - Batch
- (Data analytics, ML/AI applications)



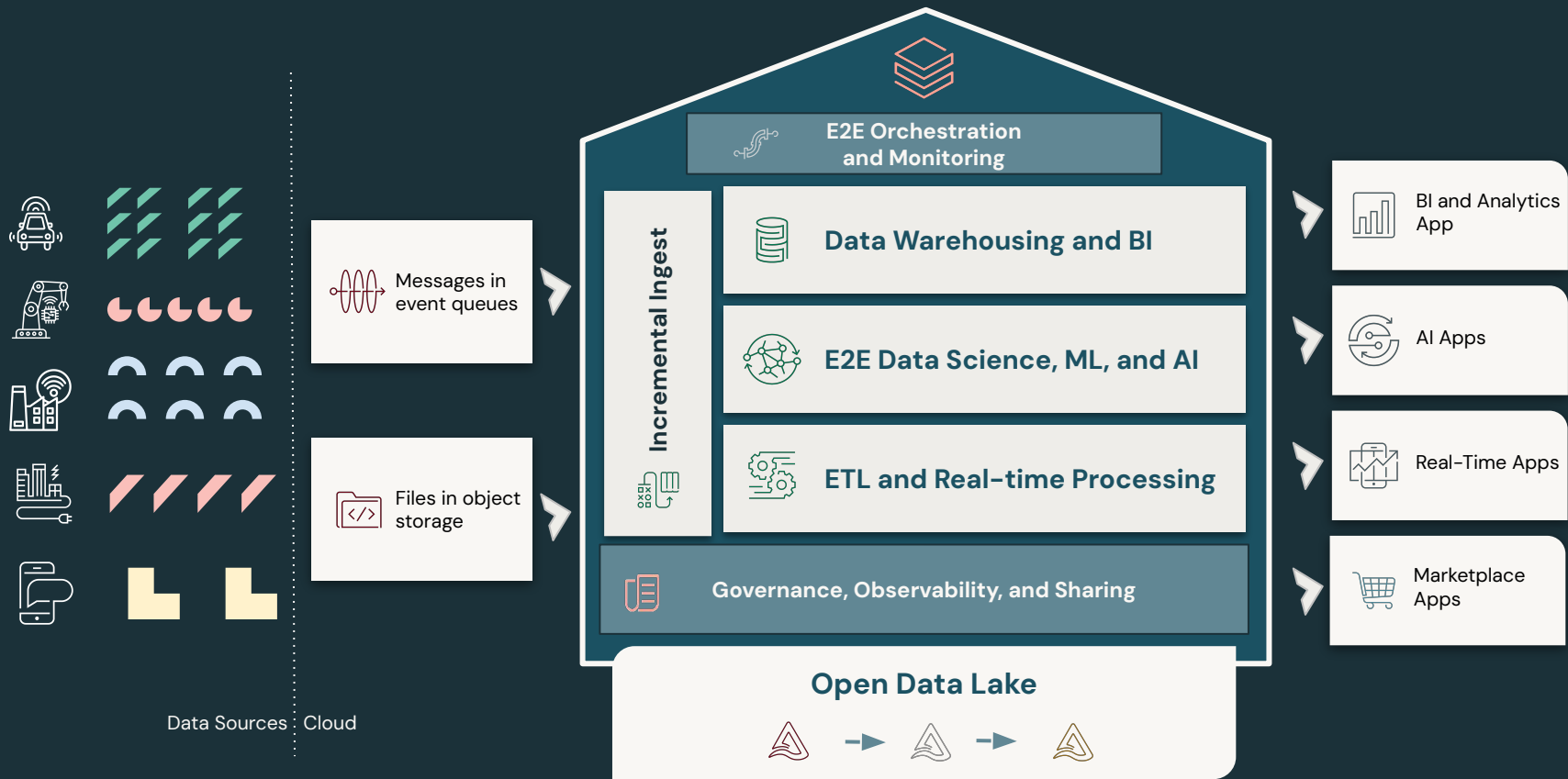
Today, Analytics and AI for IoT data is complicated..



...making it more difficult to deliver insights and bring AI at the right place at the right time



IoT Analytics and AI on the Databricks Lakehouse

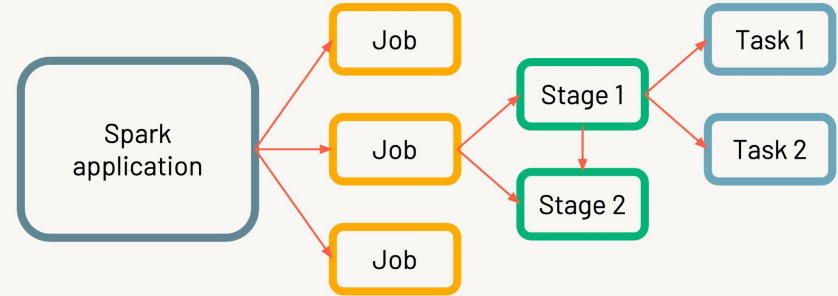


Distributed computing based on Spark

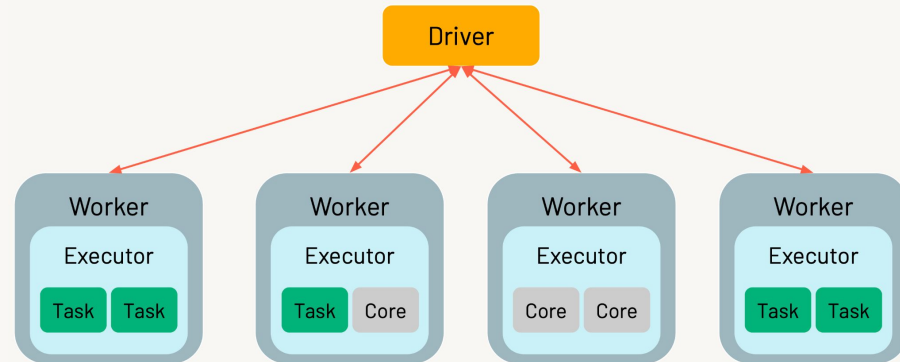
- De-facto standard unified analytics engine for big data processing
- Largest open-source project in data processing
- Technology created by the founders of Databricks

Databricks as a platform makes Spark much easier and user-friendly to use

Spark Execution



Spark Cluster



Delta Lake as Storage



ACID Transactions

Protect your data with serializability, the strongest level of isolation



Scalable Metadata

Handle petabyte-scale tables with billions of partitions and files with ease



Time Travel

Access/revert to earlier versions of data for audits, rollbacks, or reproduce



Open Source

Community driven, open standards, open protocol, open discussions



Unified Batch/Streaming

Exactly once semantics ingestion to backfill to interactive queries



Schema Evolution / Enforcement

Prevent bad data from causing data corruption



Audit History

Delta Lake log all change details providing a full audit trail



DML Operations

SQL, Scala/Java and Python APIs to merge, update and delete datasets



Delta tables = Parquet files + metadata (JSON)

Parquet format can store data from a collection of raw OPC UA produced JSON files in 10-1000x less space

Columns contain similar data, which usually can be compressed effectively.
Supports nested data.

Parquet

A columnar storage format

Name	Score	ID

Row-Oriented data on disk

Kit	4.2	1	Alex	4.5	2	Terry	4.1	3
-----	-----	---	------	-----	---	-------	-----	---

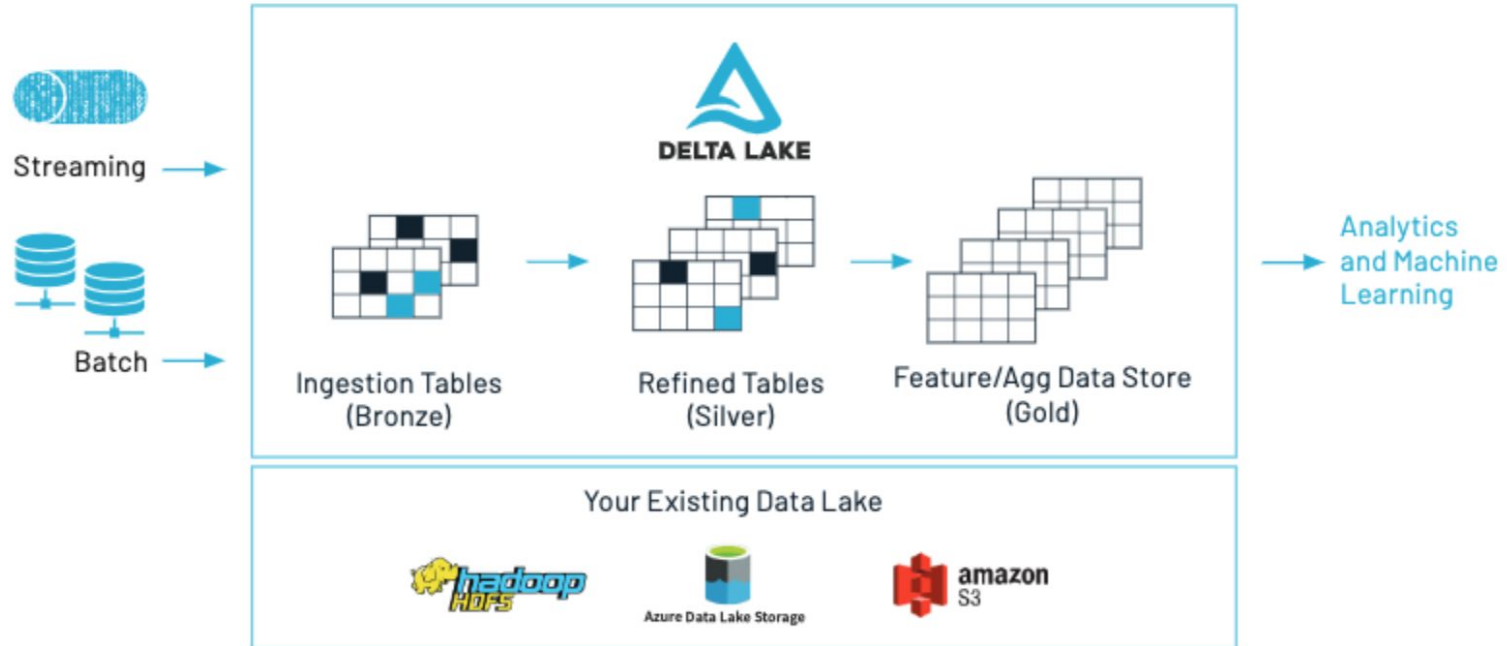
Column-Oriented data on disk

Kit	Alex	Terry	4.2	4.5	4.1	1	2	3
-----	------	-------	-----	-----	-----	---	---	---

©2022 Databricks Inc. — All rights reserved



Medallion Architecture



Delta Live Tables

The best way to do ETL on the lakehouse

```
CREATE STREAMING TABLE raw_data
AS SELECT *
FROM cloud_files ("/raw_data",
"json")
```

```
CREATE MATERIALIZED VIEW clean_data
AS SELECT ...
FROM LIVE.raw_data
```



Accelerate ETL development

Declare **SQL or Python** and DLT automatically orchestrates the DAG, handles retries, changing data



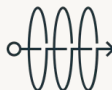
Automatically manage your infrastructure

Automates complex tedious activities like **recovery, auto-scaling, and performance optimization**



Ensure high data quality

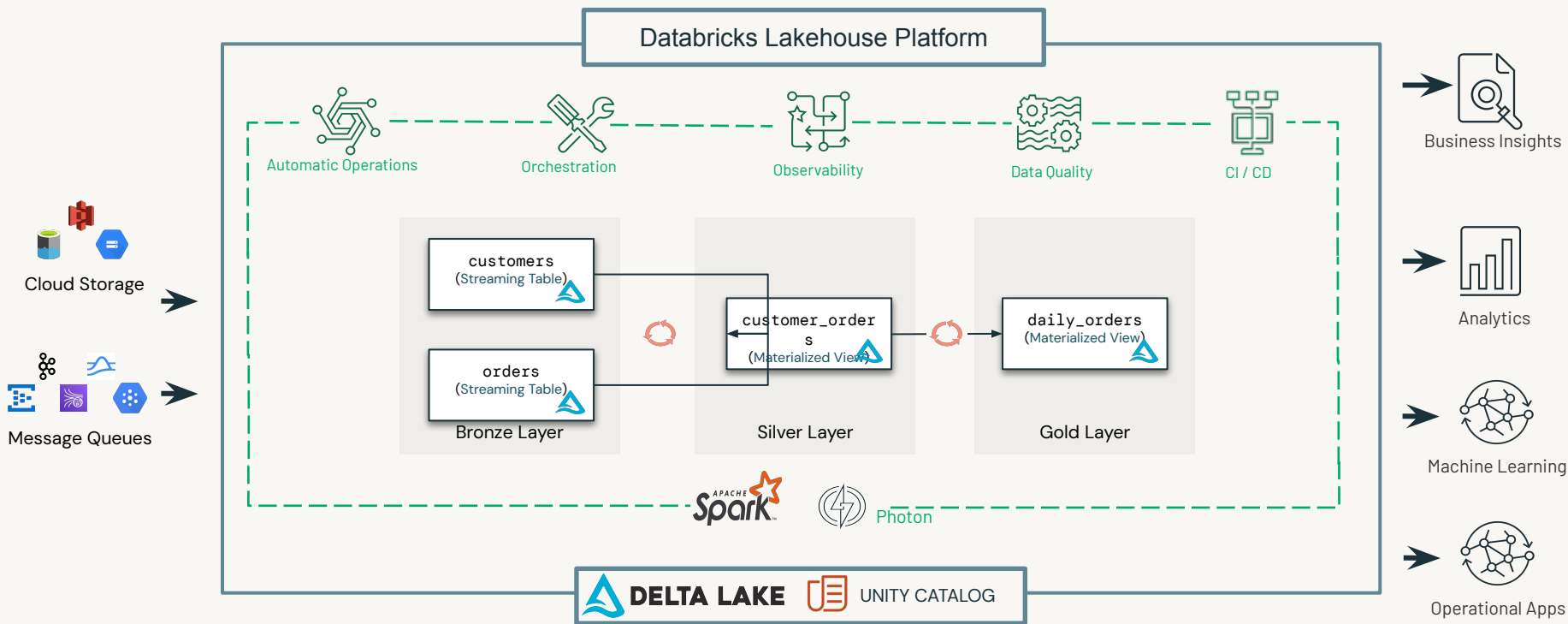
Deliver reliable data with built-in **quality controls, testing, monitoring, and enforcement**



Unify batch and streaming

Get the simplicity of SQL with freshness of streaming with one **unified API**

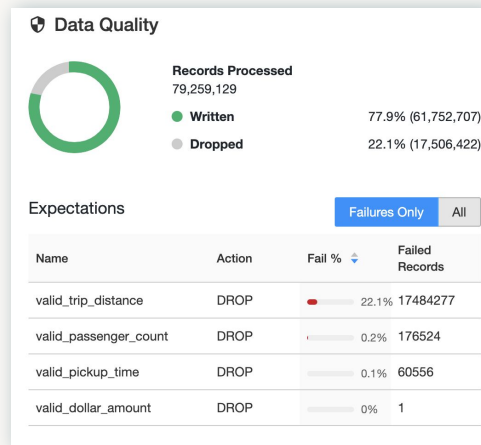
Build Production ETL Pipelines with DLT



Data quality validation and monitoring

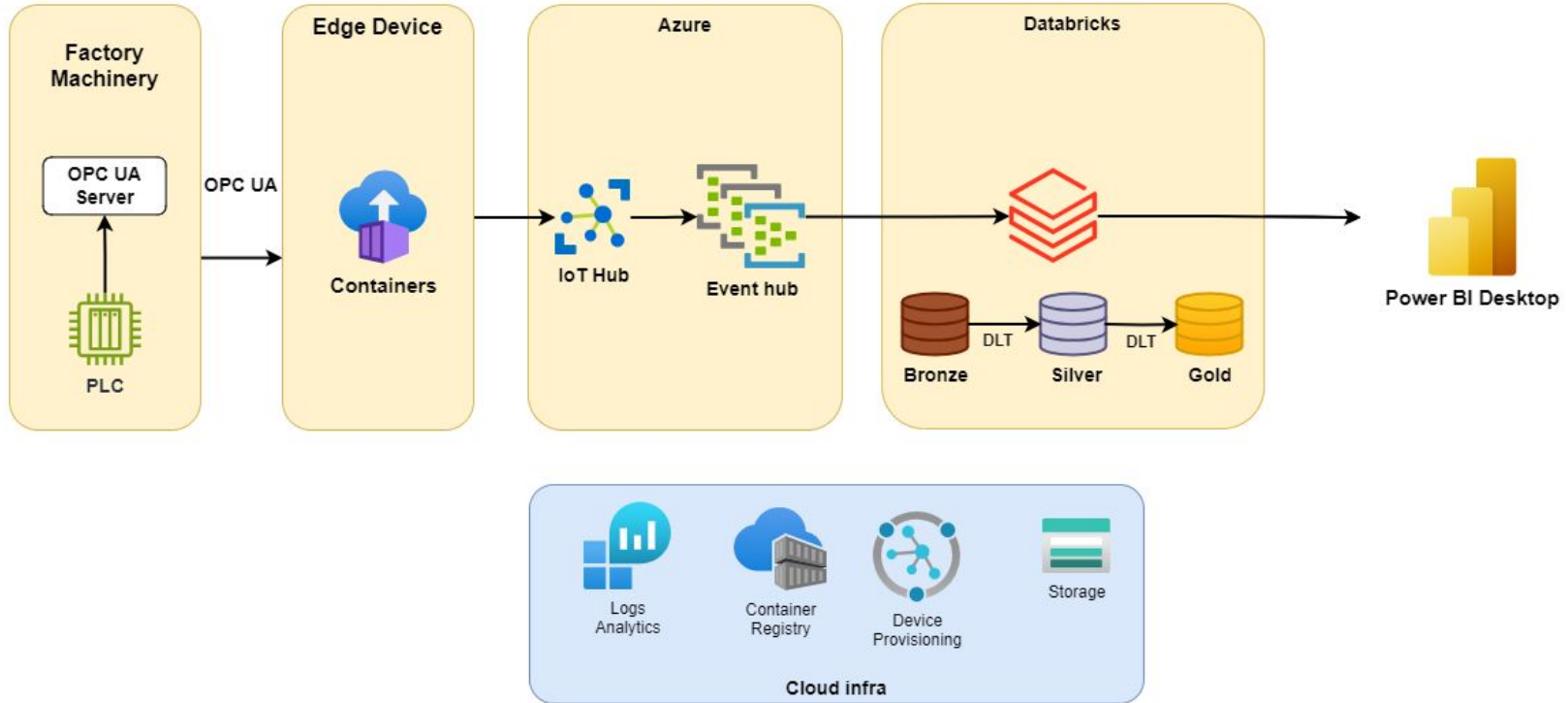
- Define data quality and integrity controls within the pipeline with data expectations
- Address data quality errors with flexible policies: fail, drop, alert, quarantine(future)
- All data pipeline runs and quality metrics are captured, tracked and reported

```
/* Stage 1: Bronze Table drop invalid rows */  
CREATE STREAMING LIVE TABLE fire_account_bronze AS  
( CONSTRAINT valid_account_open_dt EXPECT (account_dt is not null  
and (account_close_dt > account_open_dt)) ON VIOLATION DROP ROW  
COMMENT "Bronze table with valid account ids"  
SELECT * FROM fire_account_raw ...
```



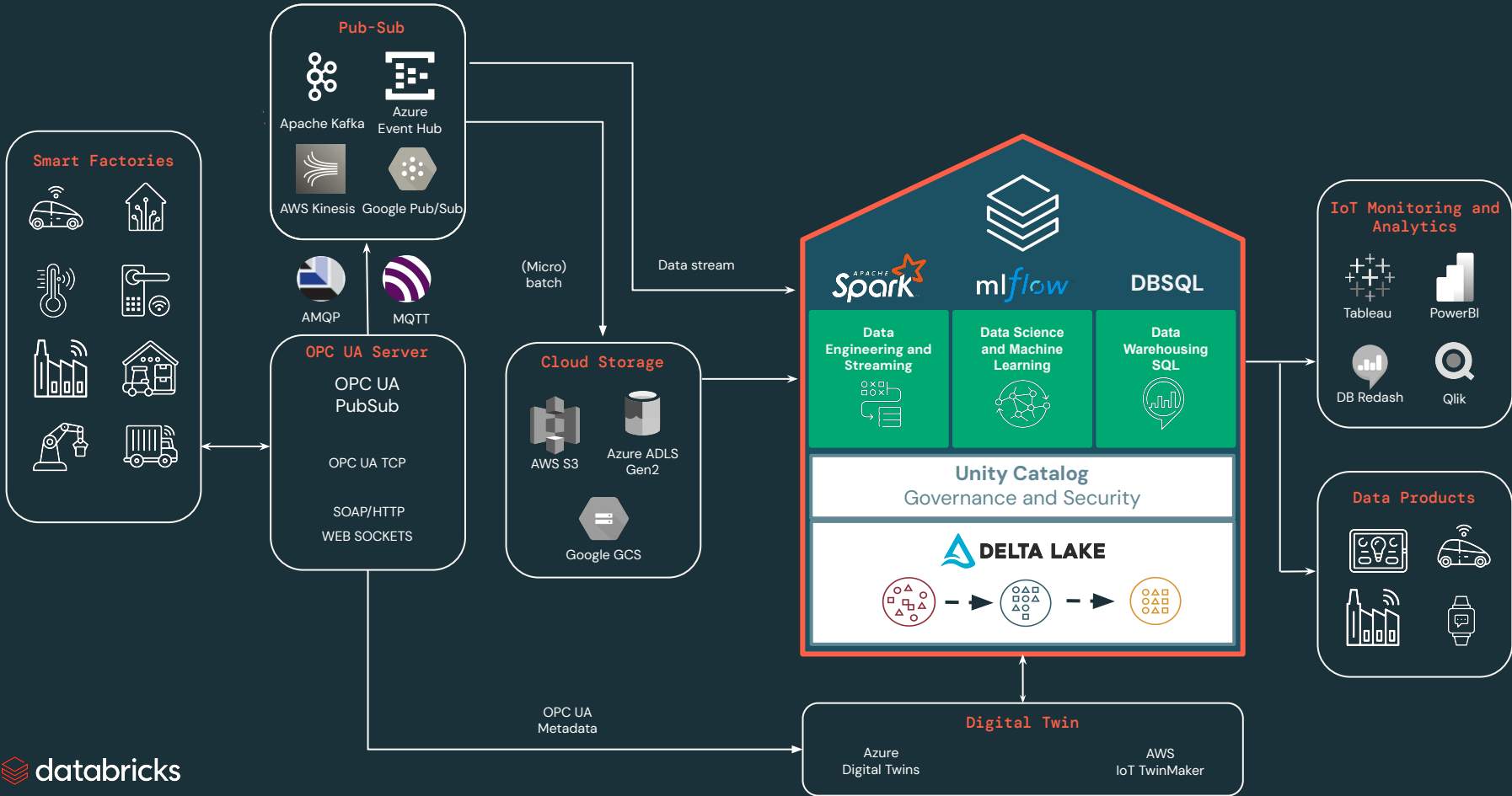
PATH TO SUCCESS

Scalable architecture for industrial data & analytics

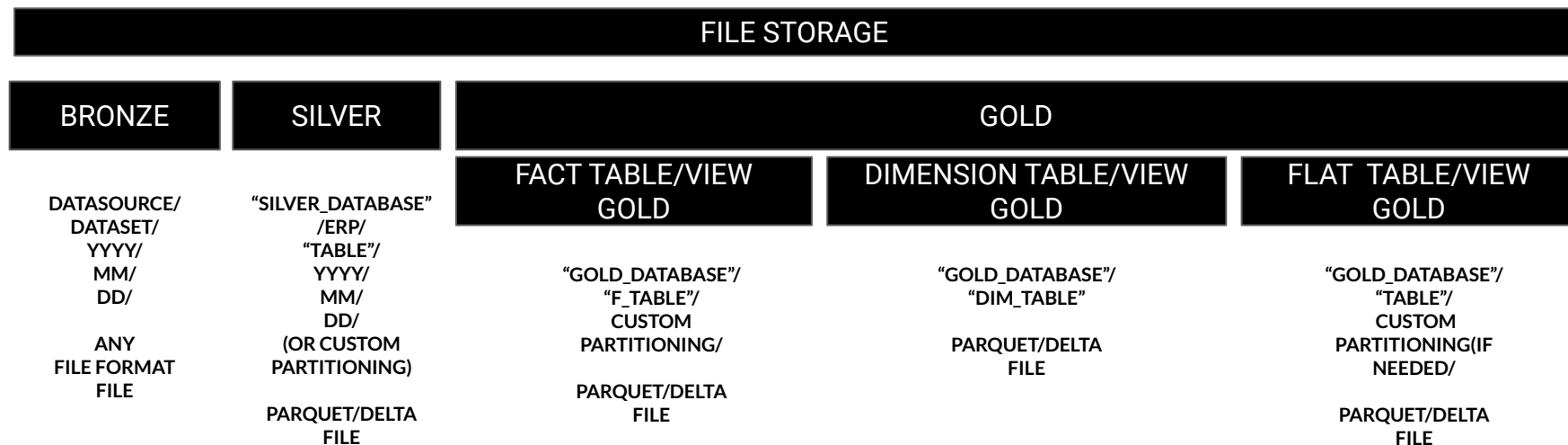


Thank you





Lakehouse architecture, transformations

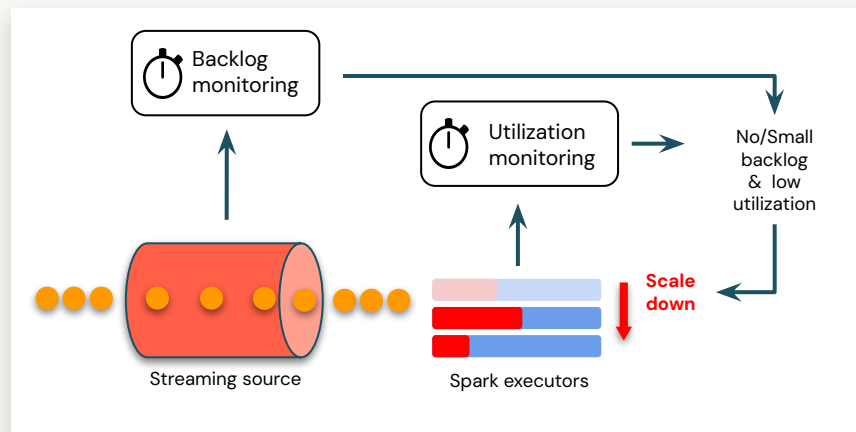


Enhanced Autoscaling

Save infrastructure costs while maintaining end-to-end latency SLAs for streaming workloads

Problem

Optimize infrastructure spend when making scaling decisions for streaming workloads



- Built to handle streaming workloads which are spiky and unpredictable
- Shuts down nodes when utilization is low while guaranteeing task execution
- Only scales up to needed # of nodes

AWS	Azure	GCP
Generally Available	Generally Available	Public Preview GA Coming Soon