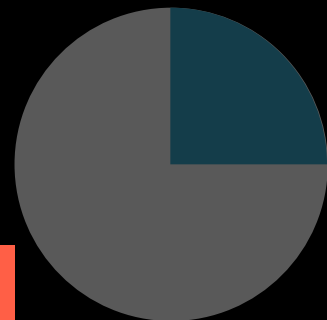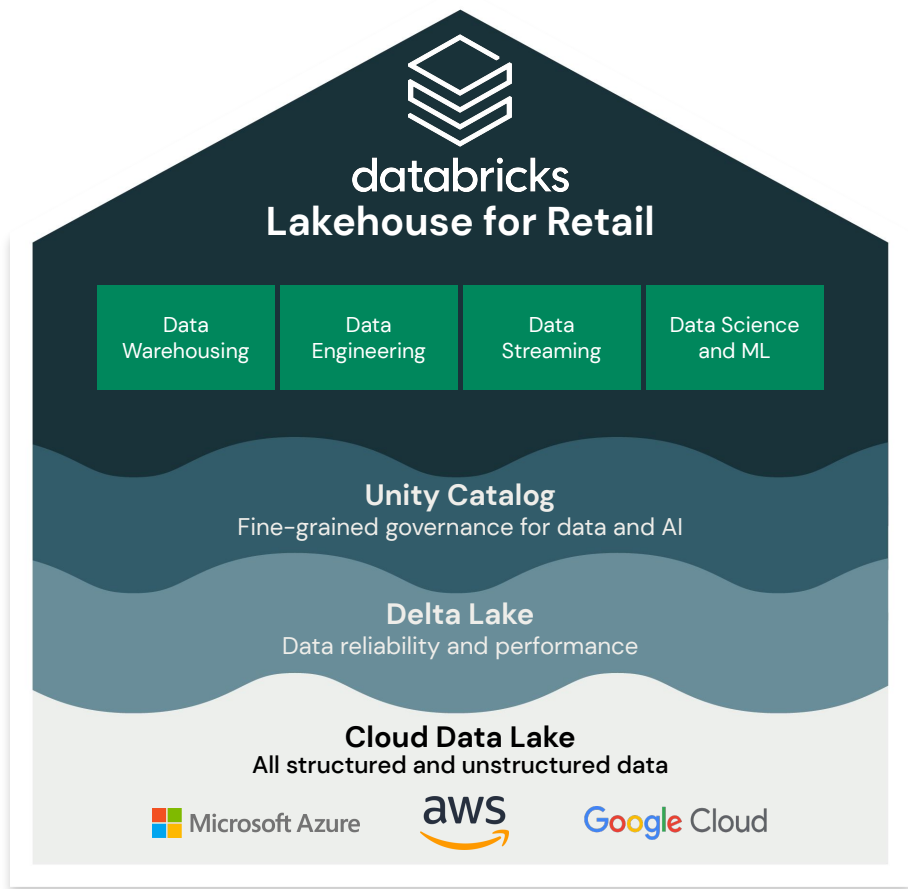# Delta Live Tables

# Databricks Lakehouse for Retail

**Simple**
Unify your data warehousing and AI
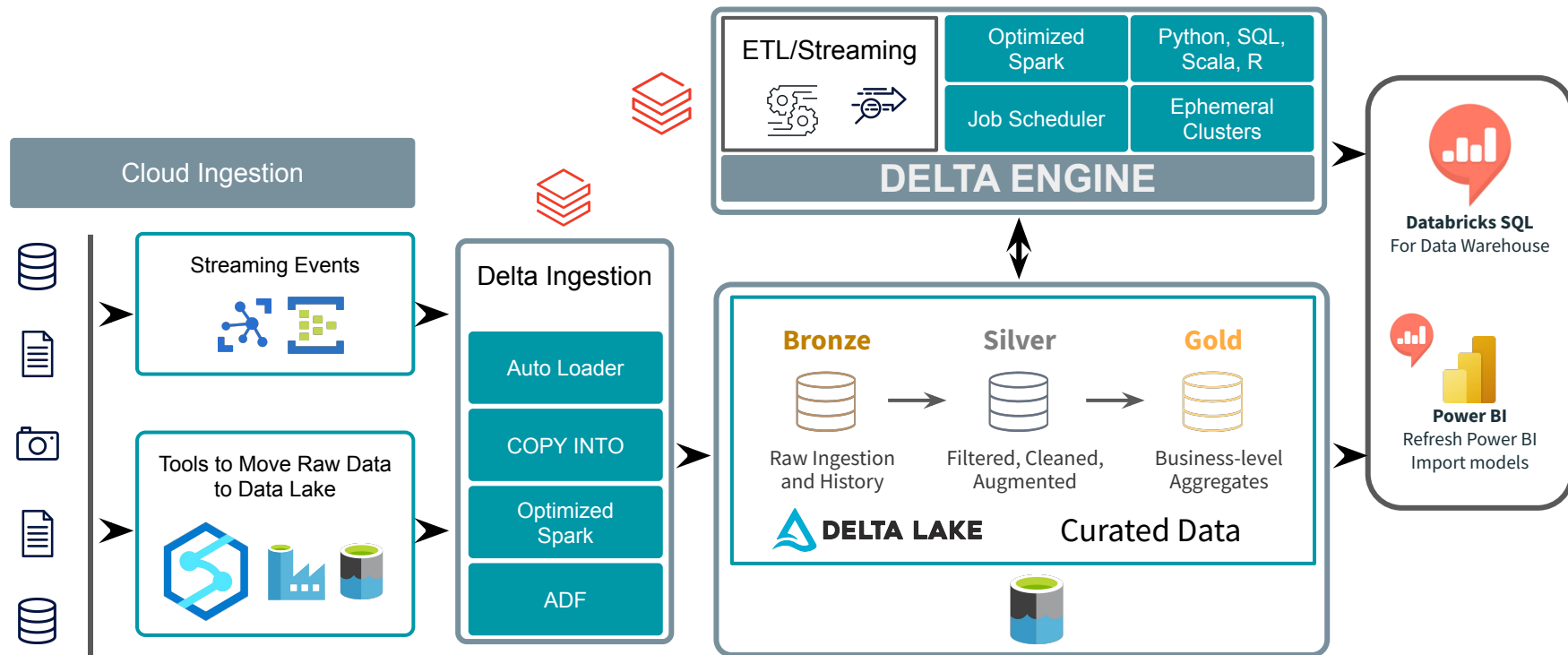use cases on a single platform

**Open**
Built on open source and open standards

**Multi-cloud**
One consistent data platform across clouds
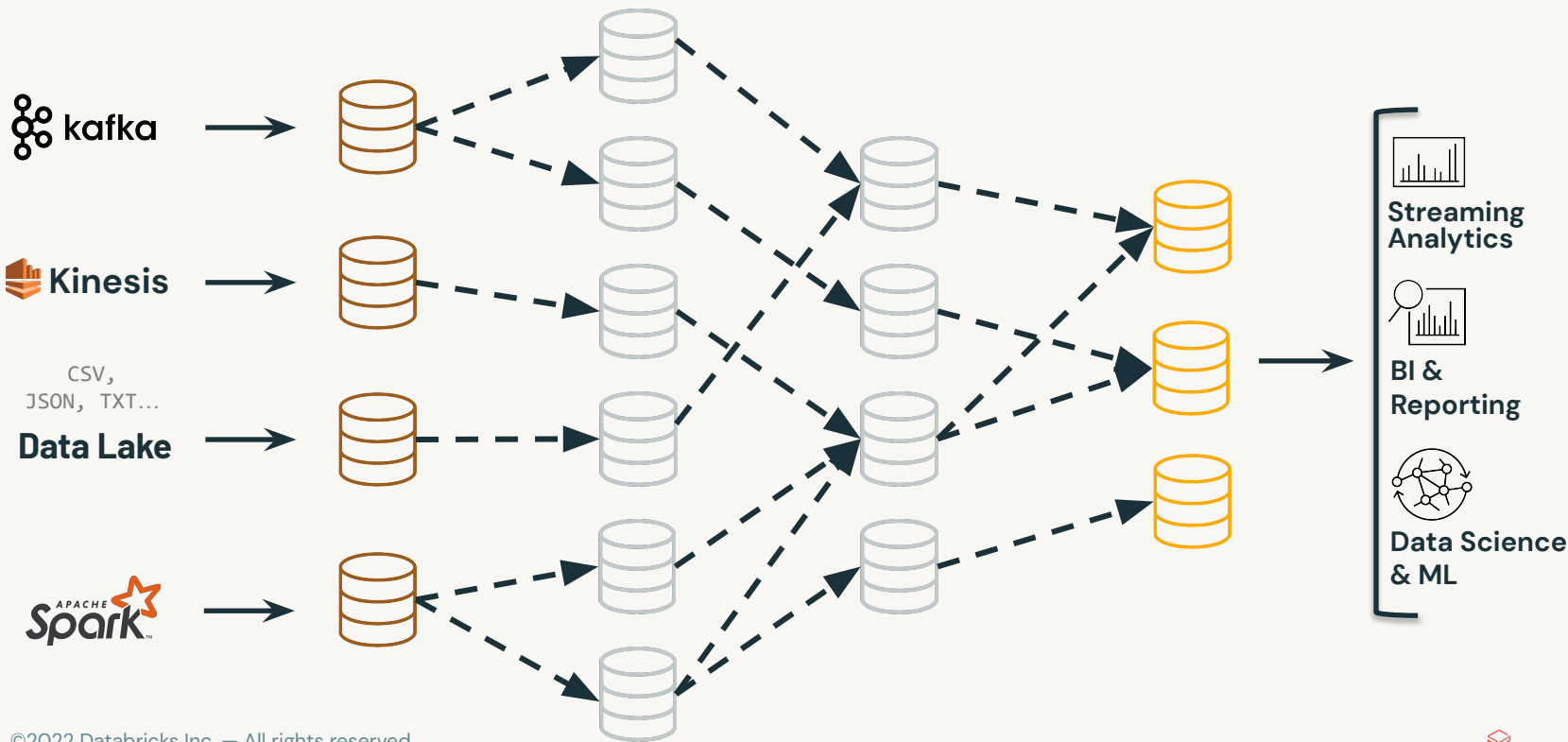
# Data Engineering and Real Time Data Applications

# Building the foundation of the Lakehouse

Greatly improve data quality for end users – on your existing data lake



kafka

Kinesis

CSV, JSON, TXT...

Data Lake

APACHE Spark

**BRONZE**

Raw Ingestion and History

**SILVER**

Filtered, Cleaned, Augmented

**GOLD**

Business–level Aggregates

Streaming Analytics

BI & Reporting

Data Science & ML

**Quality**

# But the reality is not so simple

Maintaining data quality and reliability at scale is complex and brittle

# Introducing "Workflows"

**From**

**"Jobs"**

"Cron for
Databricks Spark"

**To**

...more than Spark!

**"Workflows"** containing
- Jobs
- Delta Live Tables

**General purpose orchestration
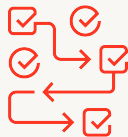for the entire Lakehouse**

# Delta Live Tables

# What is Delta Live Tables?

**Modern software engineering for ETL processing**

Delta Live Tables (DLT) is the first ETL framework that uses a simple, declarative approach to building reliable data pipelines. DLT automatically manages your infrastructure at scale so data analysts and engineers can spend less time on tooling and focus on getting value from data.

**Accelerate ETL Development**
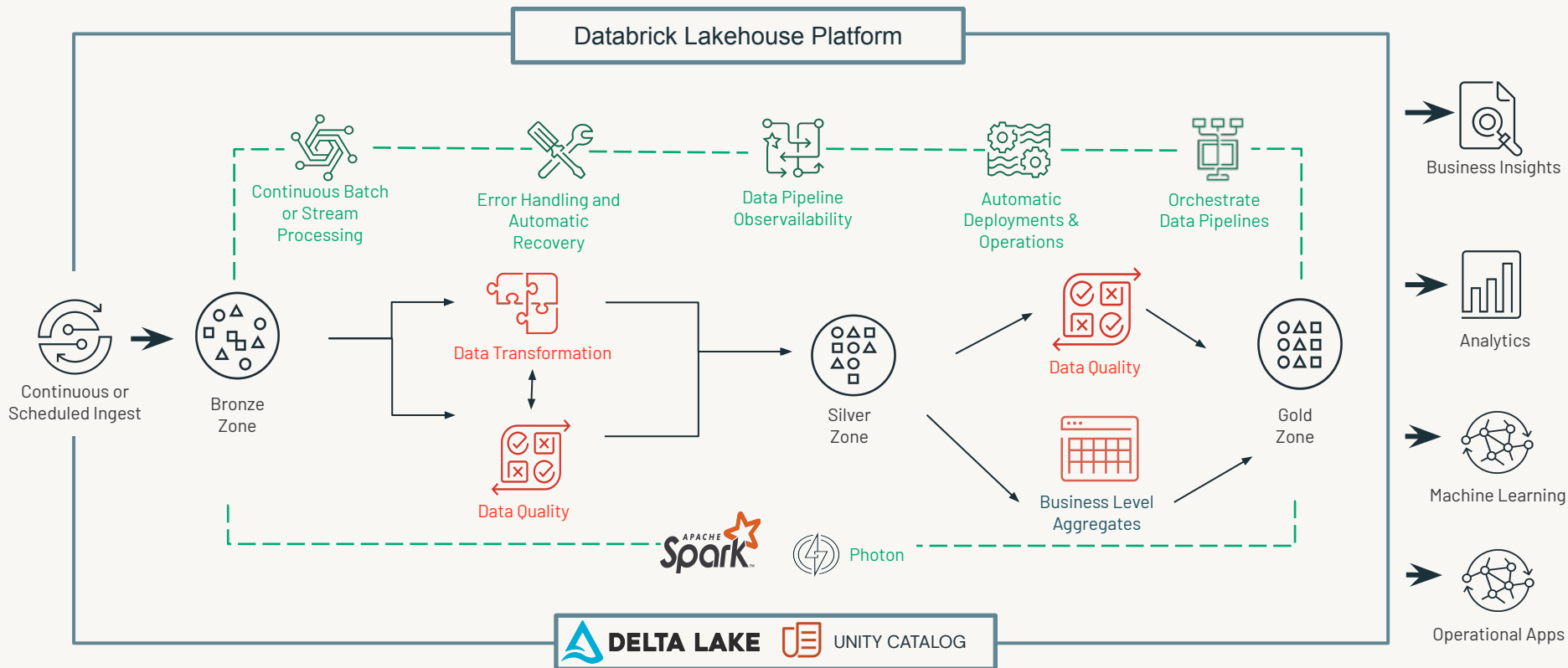
**Automatically manage your infrastructure**

**Have confidence in your data**

**Simplify batch and streaming**

# Build Production ETL Pipelines with DLT
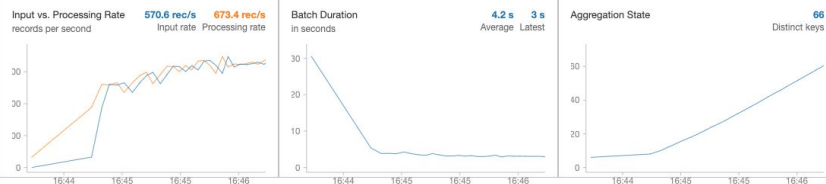
# Continuous or scheduled data ingestion



Simple SQL Syntax for Streaming Data Ingestion

```
Cmd 4
1  CREATE INCREMENTAL LIVE TABLE sales_orders_raw
2  COMMENT "The raw sales orders, ingested from /databricks-datasets."
3  TBLPROPERTIES ("quality" = "bronze")
4  AS
5  SELECT * FROM cloud_files
6  ("/databricks-datasets/retail-org/sales_orders/",
7  "json", map("cloudFiles.inferColumnTypes", "true"));
```

- Incrementally and efficiently process new data files as they arrive in cloud storage using Auto Loader

- Automatically infer schema of incoming files or superimpose what you know with Schema Hints

- Automatic schema evolution

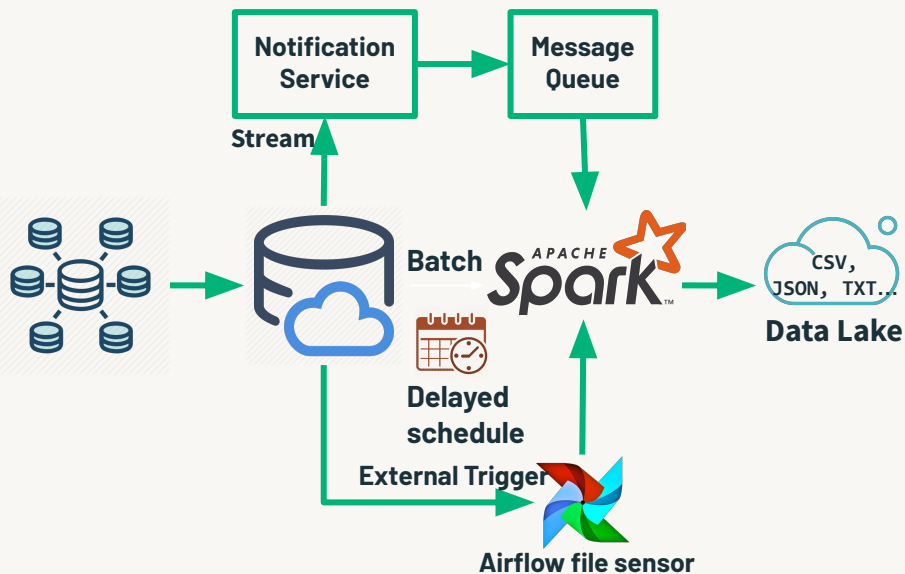- Rescue data column – never lose data again

| Schema Evolution | ✅ JSON | ✅ CSV | ✅ AVRO | Coming Soon<br>PARQUET |
|---|---|---|---|---|

# Databricks Ingest: Auto Loader

*Load new data easily and efficiently as it arrives in cloud storage*

## Before

**Notification Service** → **Message Queue**

Stream

**Batch** APACHE **Spark**

**Delayed schedule**

**External Trigger**

**Airflow file sensor**

CSV, JSON, TXT...
**Data Lake**

**Gets too complicated for multiple jobs**

## After

Auto Loader

- Pipe data from cloud storage into your data lake with Delta Lake as it arrives
- "Set and forget" model eliminates complex setup

Launch Blog Post

# Declarative SQL & Python APIs

Source

```
/* Create a temp view on the accounts table */
CREATE STREAMING LIVE VIEW account_raw AS
SELECT * FROM cloud_files("/data", "csv");
```

Bronze

```
/* Stage 1: Bronze Table drop invalid rows */
CREATE STREAMING LIVE TABLE account_bronze AS
COMMENT "Bronze table with valid account ids"
SELECT * FROM fire_account_raw ...
```
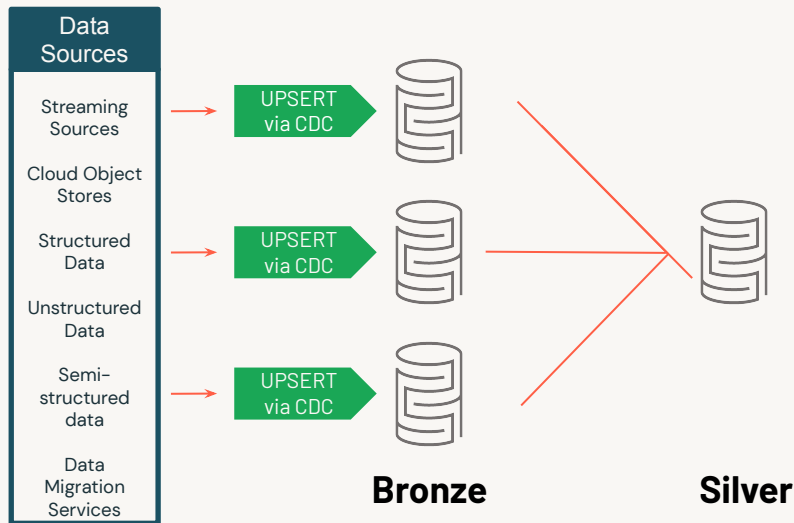
Silver

```
/* Stage 2:Send rows to Silver, run validation rules */
CREATE STREAMING LIVE TABLE account_silver AS
COMMENT "Silver Accounts table with validation checks"
SELECT * FROM fire_account_bronze ...
```

Gold

- Use intent-driven declarative development to abstract away the **"how"** and define **"what"** to solve

- Automatically generate **lineage** based on table dependencies across the data pipeline

- Automatically checks for errors, missing dependencies and syntax errors

# Change data capture (CDC)



Data Sources:
- Streaming Sources
- Cloud Object Stores
- Structured Data
- Unstructured Data
- Semi-structured data
- Data Migration Services
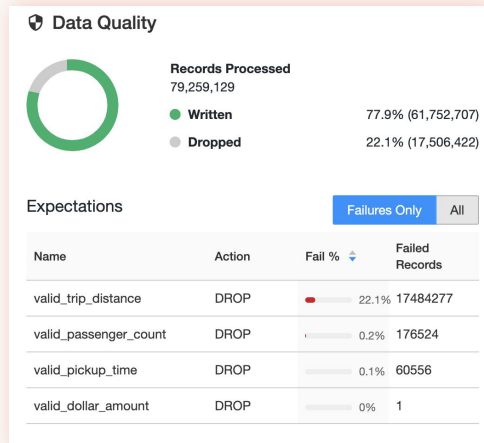
UPSERT via CDC → **Bronze** → **Silver**

- Stream change records (inserts, updates, deletes) from any data source supported by DBR, cloud storage, or DBFS

- Simple, declarative "APPLY CHANGES INTO" API for SQL or Python

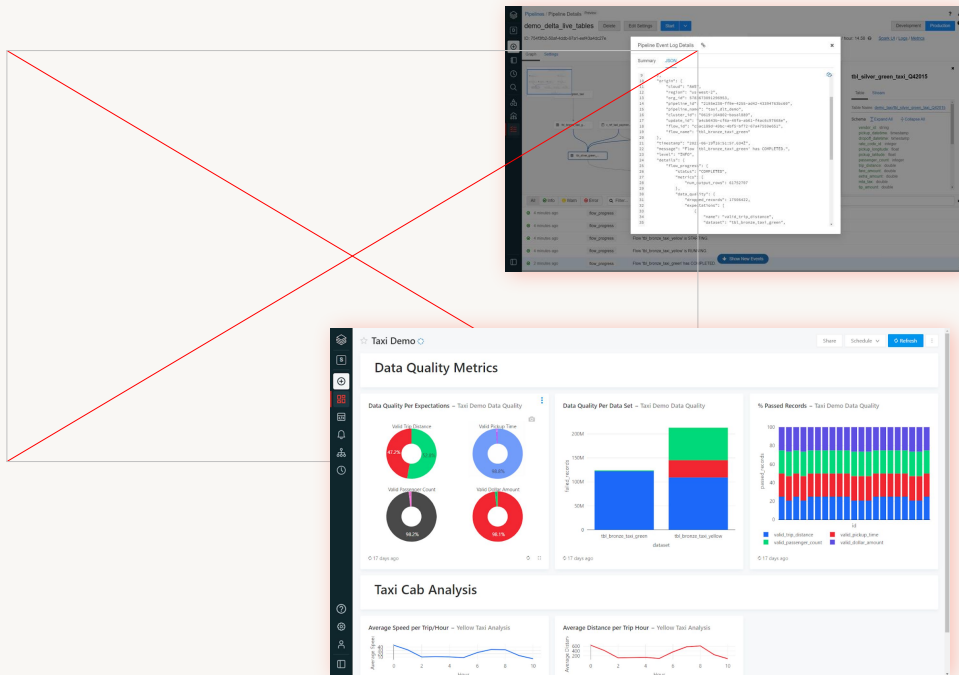- Handles out-of-order events

- Schema evolution

- SCD2 support

# 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 (acconut_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 ...
```



🛡 Data Quality

Records Processed
79,259,129

● Written          77.9% (61,752,707)
○ Dropped          22.1% (17,506,422)

Expectations                    [Failures Only] [All]

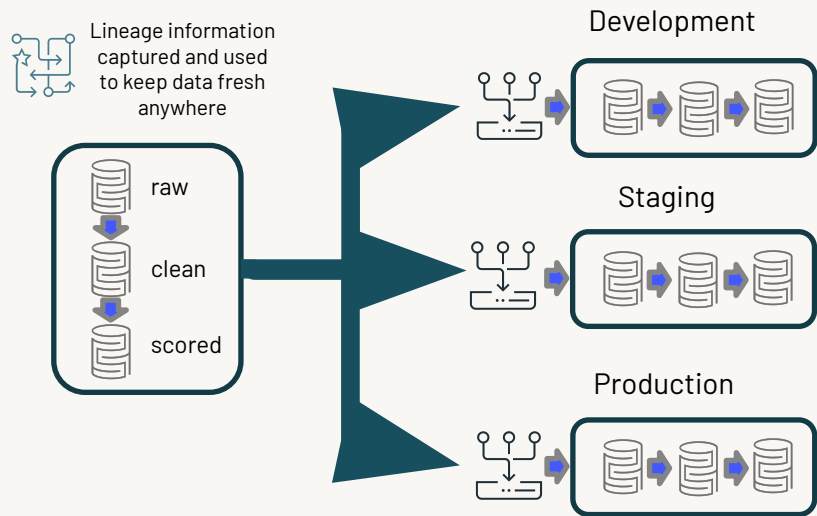| Name | Action | Fail % ⇕ | Failed Records |
|------|--------|----------|----------------|
| valid_trip_distance | DROP | 22.1% | 17484277 |
| valid_passenger_count | DROP | 0.2% | 176524 |
| valid_pickup_time | DROP | 0.1% | 60556 |
| valid_dollar_amount | DROP | 0% | 1 |

# Data pipeline observability



- High-quality, high-fidelity lineage diagram that provides visibility into how data flows for impact analysis

- Granular logging for operational, governance, quality  and status of the data pipeline at a row level

- Continuously monitor data pipeline jobs to ensure continued operation

- Notifications using Databricks SQL

# Automated ETL development lifecycle

- Develop in environment(s) separate from production with the ability to easily test it before deploying – entirely in SQL

- Deploy and manage environments using parameterization

- Unit testing and documentation

- Enables metadata–driven ability to programatically scale to 100s of tables/pipelines dynamically

Lineage information captured and used to keep data fresh anywhere

raw

clean

scored

Development

Staging

Production

# Automated ETL operations



- Reduce down time with automatic error handling and easy replay

- Eliminate maintenance with automatic optimizations of all Delta Live Tables

- Auto-scaling adds more resources automatically when needed.

# DLT Demo