# Product Safe Harbor Statement

This information is provided to outline Databricks' general product direction and is for informational purposes only. Customers who purchase Databricks services should make their purchase decisions relying solely upon services, features, and functions that are currently available. Unreleased features or functionality described in forward-looking statements are subject to change at Databricks discretion and may not be delivered as planned or at all.
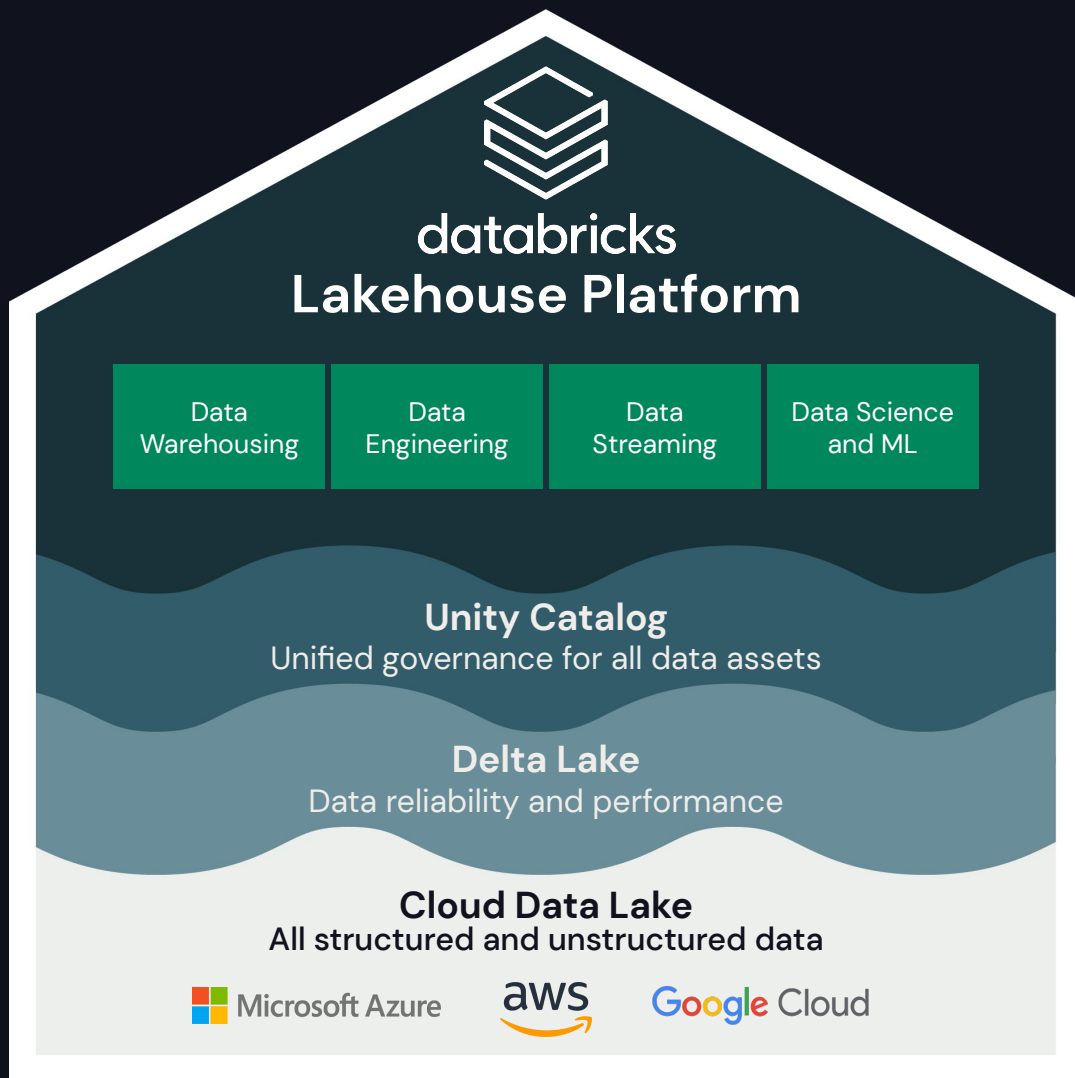
# About me

**Tao Feng**

- Staff Engineer at Databricks
- Working on Data Discovery and Lineage
- Co-Creator of Amundsen and Apache Airflow PMC
- Previously worked at Lyft, and various other tech companies

# Agenda

- Background
- Lineage Demo
- Lineage Deep Dive
- Lineage Roadmap

# Background

# What Is Data Lakehouse

Databricks
Lakehouse Platform

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

Multicloud
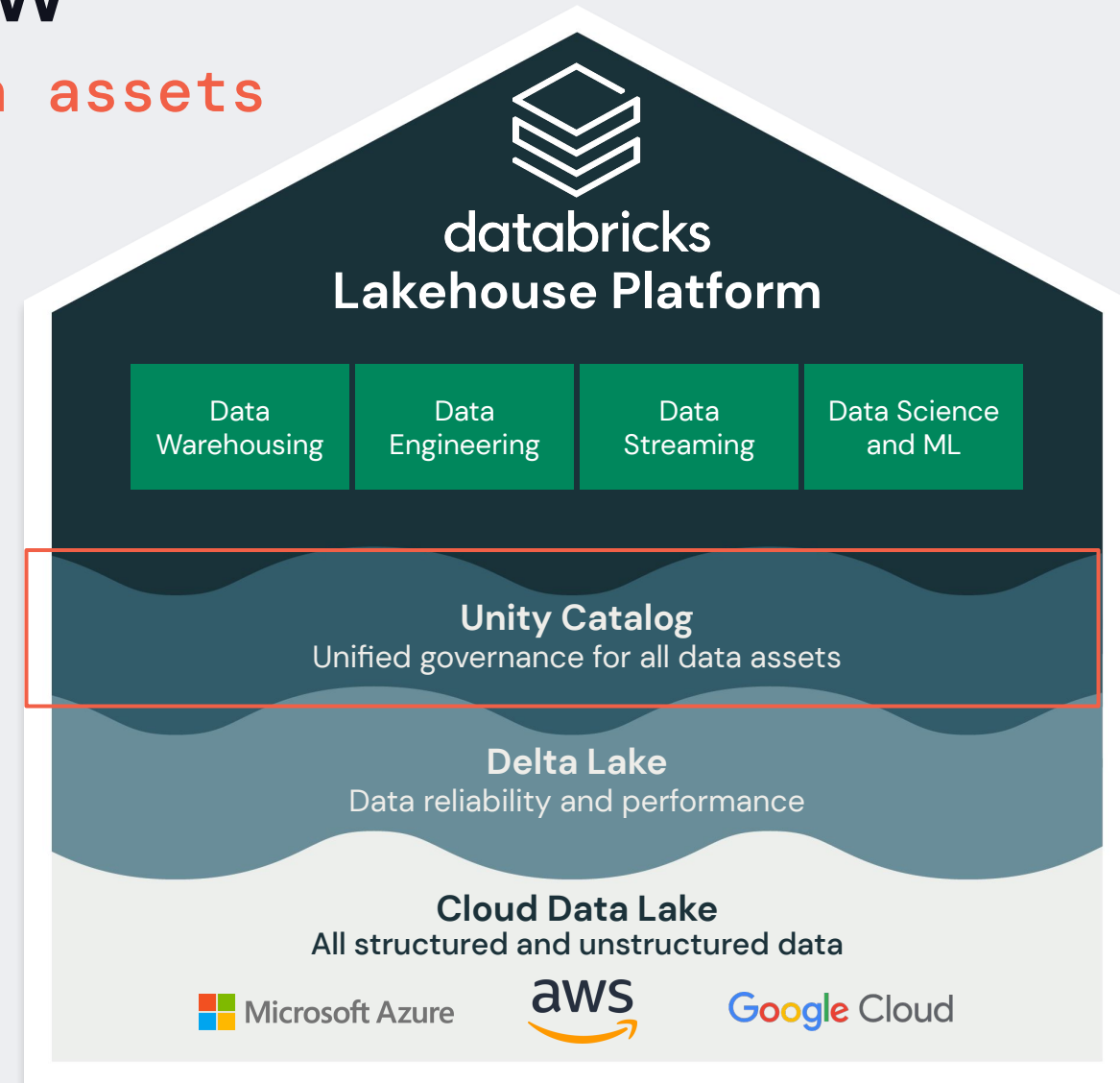One consistent data platform
across clouds

Open
Built on open source and
open standards

# Unity Catalog – Overview
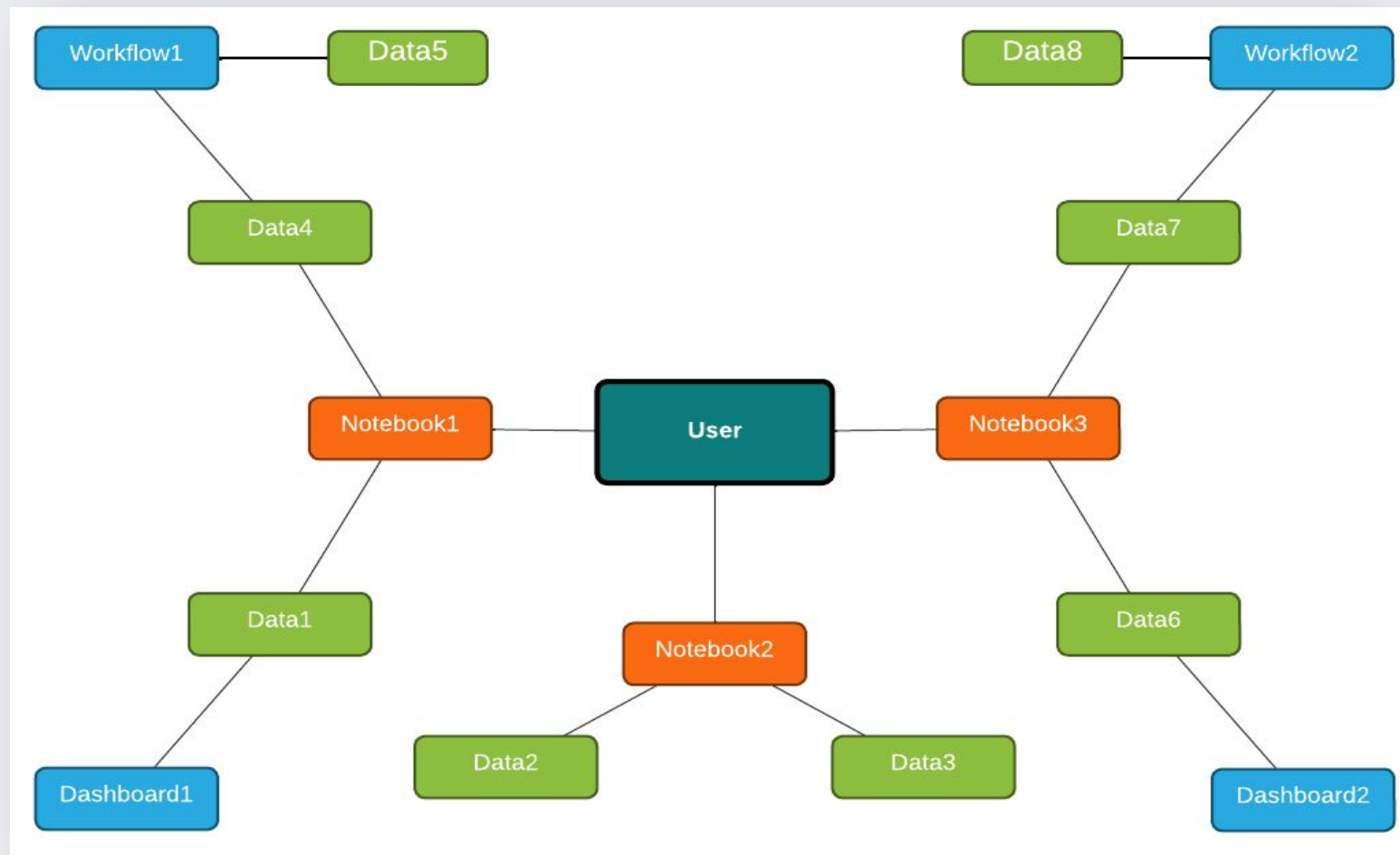## Unified governance for all data assets

- Centralized governance for Data and AI

- Built-in search and discovery

- Performance and scalability

- Automated lineage for all workloads

- Integrated with your existing tools

- 



databricks
**Lakehouse Platform**

| Data Warehousing | Data Engineering | Data Streaming | Data Science and ML |

**Unity Catalog**
Unified governance for all data assets

**Delta Lake**
Data reliability and performance

**Cloud Data Lake**
All structured and unstructured data

Microsoft Azure    aws    Google Cloud

# What Is Data Lineage

# What Is Data Lineage

- Lineage is a **graph** that connects different **data entities** in the Lakehouse and tracks their **dependencies**

# Lakehouse Use Cases

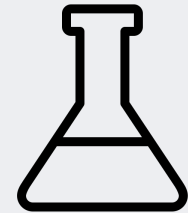# Lakehouse Personas

Analysts
Data Scientists
Product Managers
Compliance Legals
Engineers
Machine Learning Practitioners
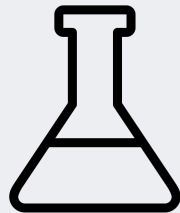
# Use Case 1: Data Discovery and Analysis

**Data Consumer**
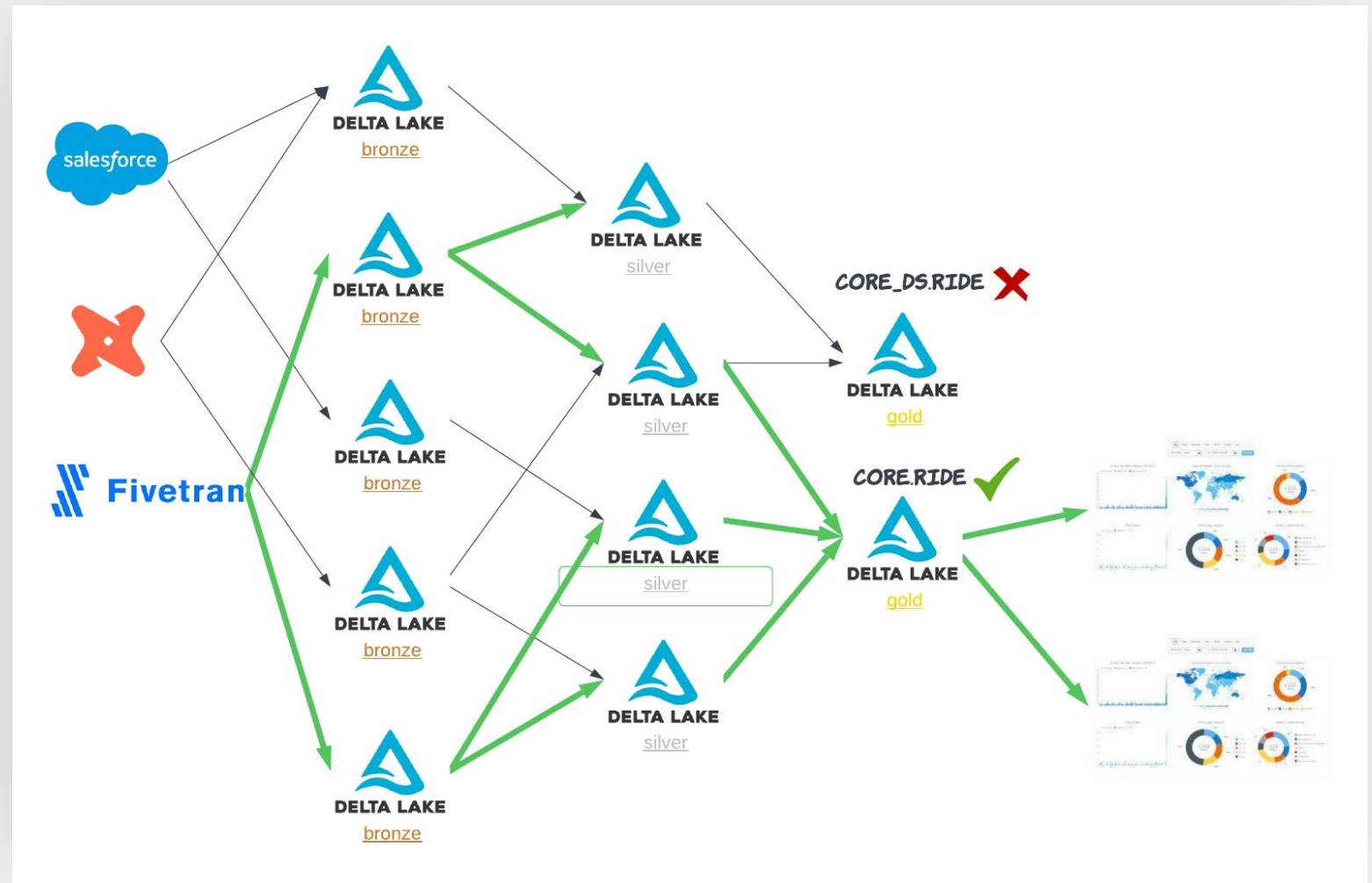
Analysts

Data Scientists

Product Managers

Machine Learning Practitioners

- Explore and understand data with its context and document and its origins
- Find out who are using the data
- Figure out what are the data producers and the consumers
- Make data driven decisions

# Use Case 1: Data Discovery and Analysis

- How do I discovery the **trustworthy** data to use for my analysis?

# Use Case 2: Change Management
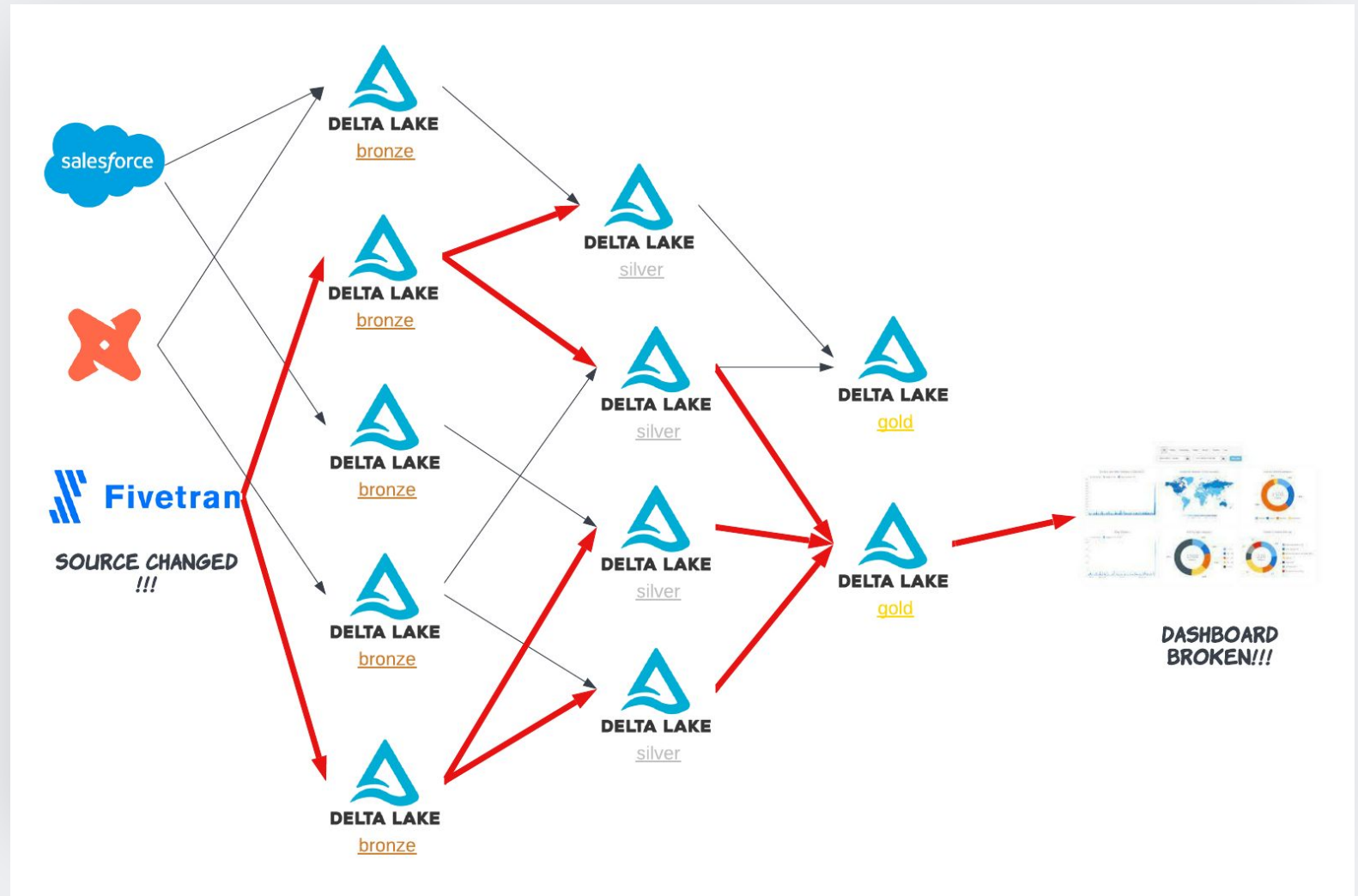
**Data Producer**

Engineers

- Identify data quality issue of a data artifact (e.g. dashboard, dataset)
- Deprecate column with downstream usage / Impact analysis
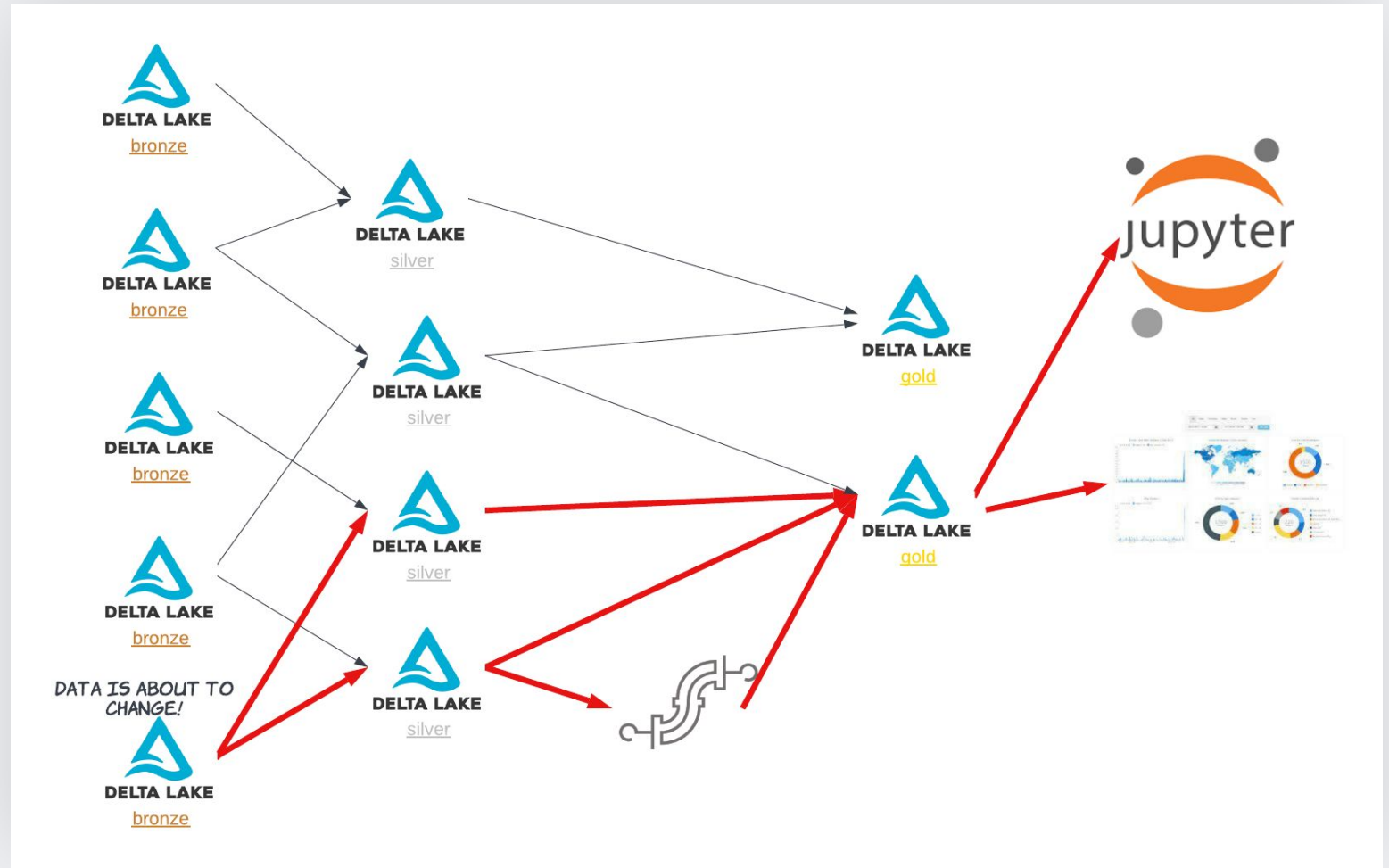- Triage unused dashboards / artifacts and reduces costs

# Use Case 2: Change Management

- Biz dashboard is broken, how to identify the **culprit**?

# Use Case 2: Change Management

- Now we need to deprecate a column of a source table, what entities we need to change? And who should we notify?

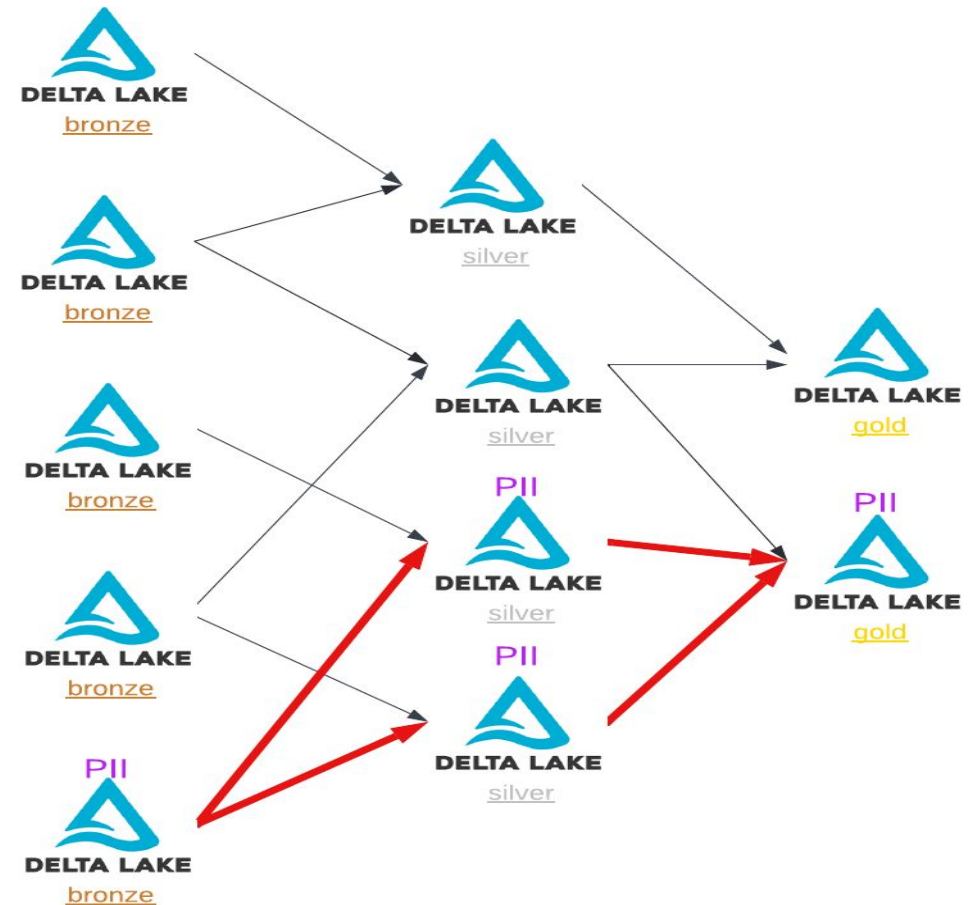# Use Case 3: Data Governance

**Compliance**



Compliance
Legals

- Identify PII or other sensitive information within the lakehouse

# Use Case 3: Data Governance

- How to identify all the tables that have PII?

# Why Is Lineage Important

# Why Is Data Lineage Important

| Compliance | Discovery | Data Observability |
|---|---|---|
| > **Regulatory** requirements to verify data lineage<br><br>> Track the **spread of sensitive data** across datasets | > Understand **context** and **trustworthiness** of data before using in analytics<br><br>> **Prevent duplicative** work and datasets | > Track down **issues / discrepancies** in reports by tracing back the data<br><br>> Analyze **impact of proposed changes** to downstream reports eg column deprecation |

# Demo

# Demo

**DATA+AI**
SUMMIT 2022

# Deep Dive

# Data Lineage With Unity Catalog



ETL / Job

Ad-hoc

Cluster or SQL Endpoint

Lineage service

Table and column lineage

Explore lineage in UI

Catalog partners

**1** Code (any language) is submitted to a cluster or SQL Warehouse endpoint

**2** Lineage service analyzes lineage events emitted from the cluster / SQL Endpoint
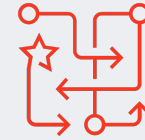
Assembles column and table level lineage

**3** Presented to the end user graphically in Databricks

Lineage can be exported via API and imported into other catalog partners

# Data Lineage With Unity Catalog

**Auto-capture** runtime data lineage across all languages

Track lineage down to the **table and column level**

Govern access by **Unity Catalog** Permission Model

Pipelines to surface lineage in **near real-time**

Table/Column/Entity Lineage **Graph** Visualization

# Lineage Example

- Define car sales revenue table, aggregate table and an analysis view
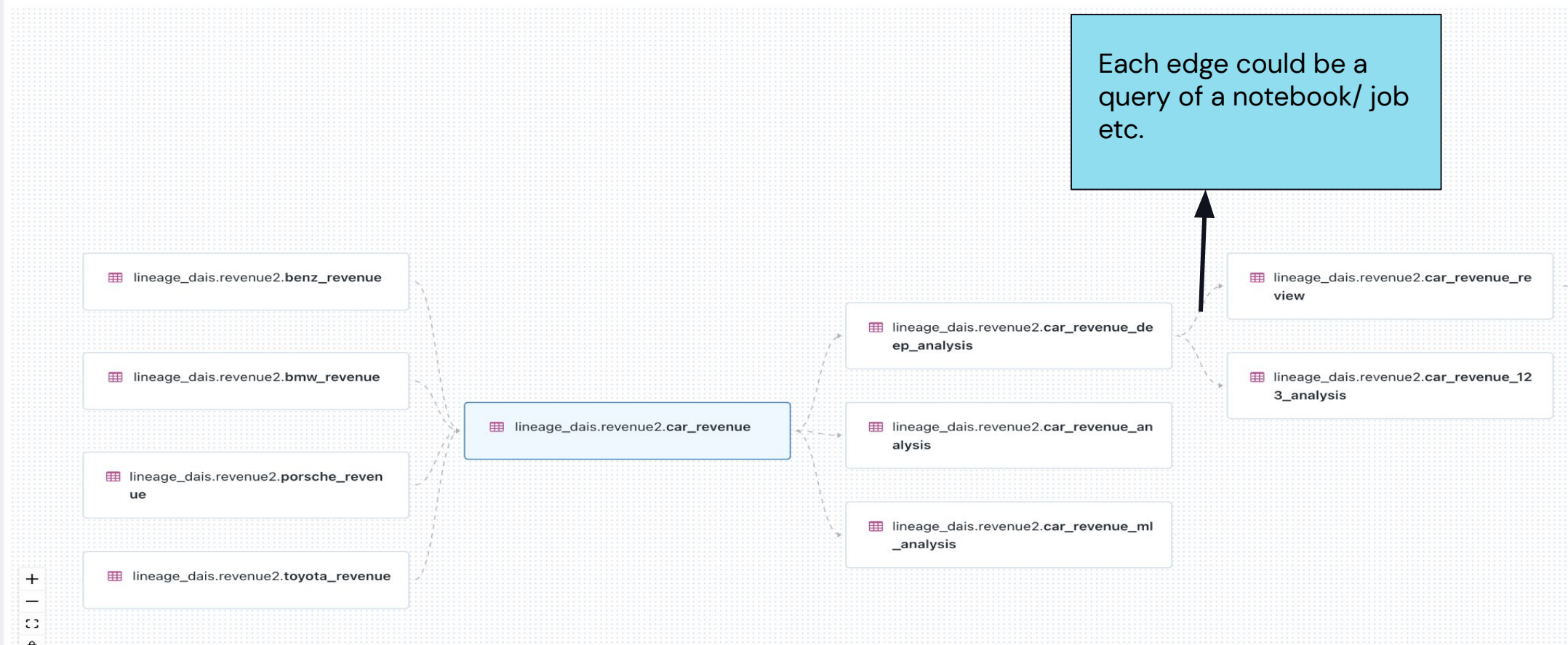
```sql
CREATE schema if not exists revenue2;

CREATE TABLE if not exists revenue2.asian_revenue  (id INT, name STRING, revenue INT);
CREATE TABLE if not exists revenue2.porsche_revenue  AS (select * from revenue2.asian_revenue);
CREATE TABLE if not exists revenue2.benz_revenue  (id INT, name STRING, revenue INT);
CREATE TABLE if not exists revenue2.bmw_revenue  (id INT, name STRING, revenue INT);
CREATE TABLE if not exists revenue2.toyota_revenue  (id INT, name STRING, revenue INT);
CREATE TABLE if not exists revenue2.car_revenue AS (select * from revenue2.porsche_revenue UNION ALL select * from revenue2.benz_revenue UNION ALL select * from
revenue2.bmw_revenue UNION ALL select * from revenue2.toyota_revenue);

CREATE TABLE if not exists revenue2.car_revenue_analysis AS (select * from revenue2.car_revenue);
CREATE TABLE if not exists revenue2.car_revenue_deep_analysis AS (select * from revenue2.car_revenue);
CREATE TABLE if not exists revenue2.car_revenue_ml_analysis AS (select * from revenue2.car_revenue);
CREATE TABLE if not exists revenue2.car_revenue_latest_analysis AS (select * from revenue2.car_revenue_analysis);
CREATE TABLE if not exists revenue2.car_revenue_analysis_snapshot AS (select * from revenue2.car_revenue_analysis);
create view if not exists revenue2.car_revenue_review as select id, name, revenue from revenue2.car_revenue_deep_analysis;
insert into revenue2.car_revenue_agg_analysis select id, name, revenue from revenue2.car_revenue_review;
```
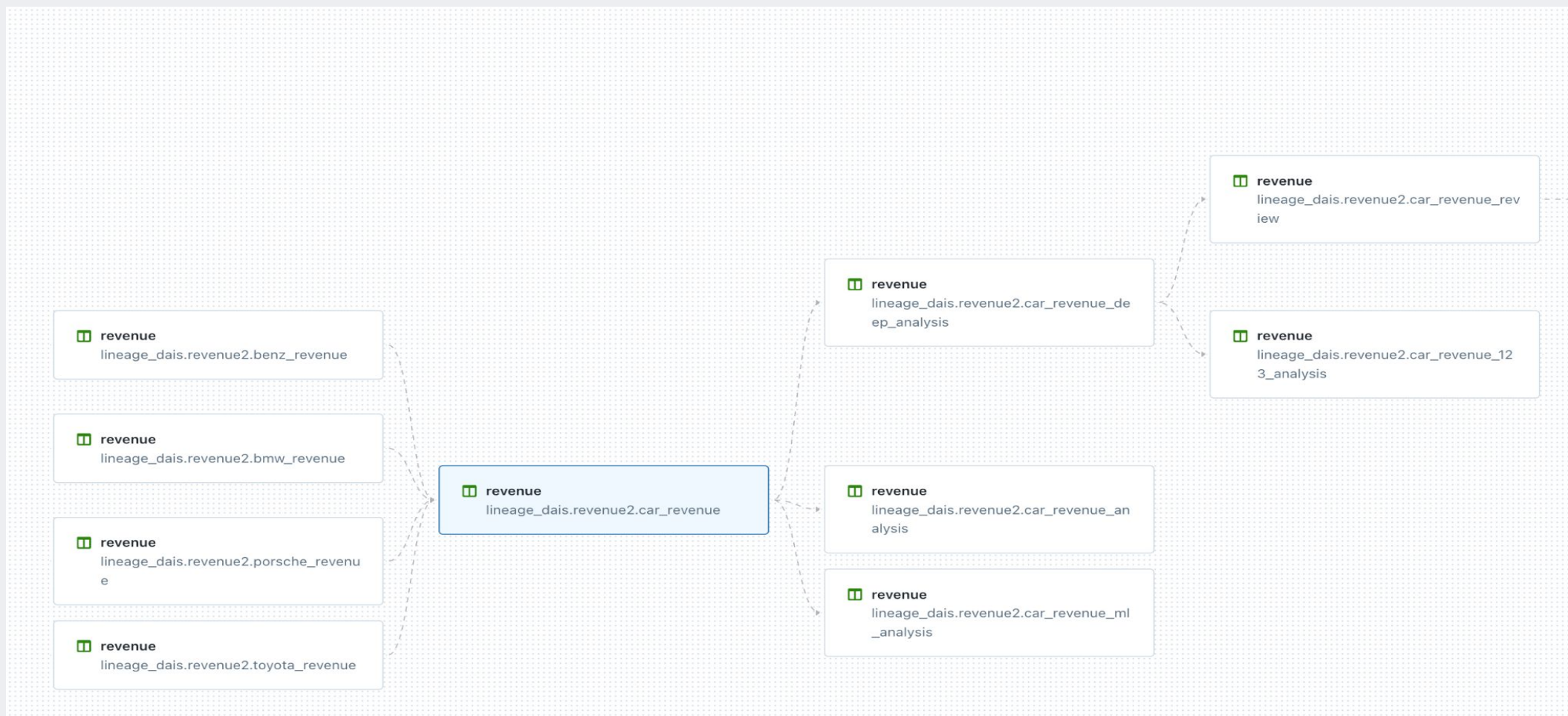
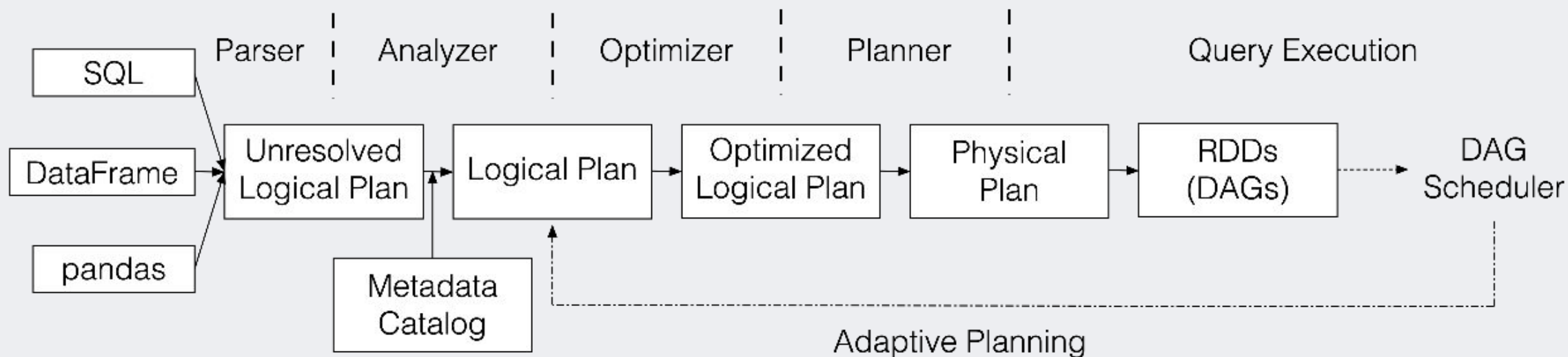# Lineage Example
## Table level lineage

# Lineage Example

## Column level lineage

# How to capture the lineage of a query?

# Life of a Query in Spark



- SQL Language

- Dataset/DataFrame/Pandas APIs: richer, language-integrated and user-friendly interfaces

# Tree: Abstractions of Users' Programs

```
INSERT INTO
lineage_dais.revenue2.car_revenue_agg_analysis

SELECT id,

    NAME,

    revenue

FROM
lineage_dais.revenue2.car_revenue_review;


lineage_dais.revenue2.car_revenue_dee
p_analysis
```
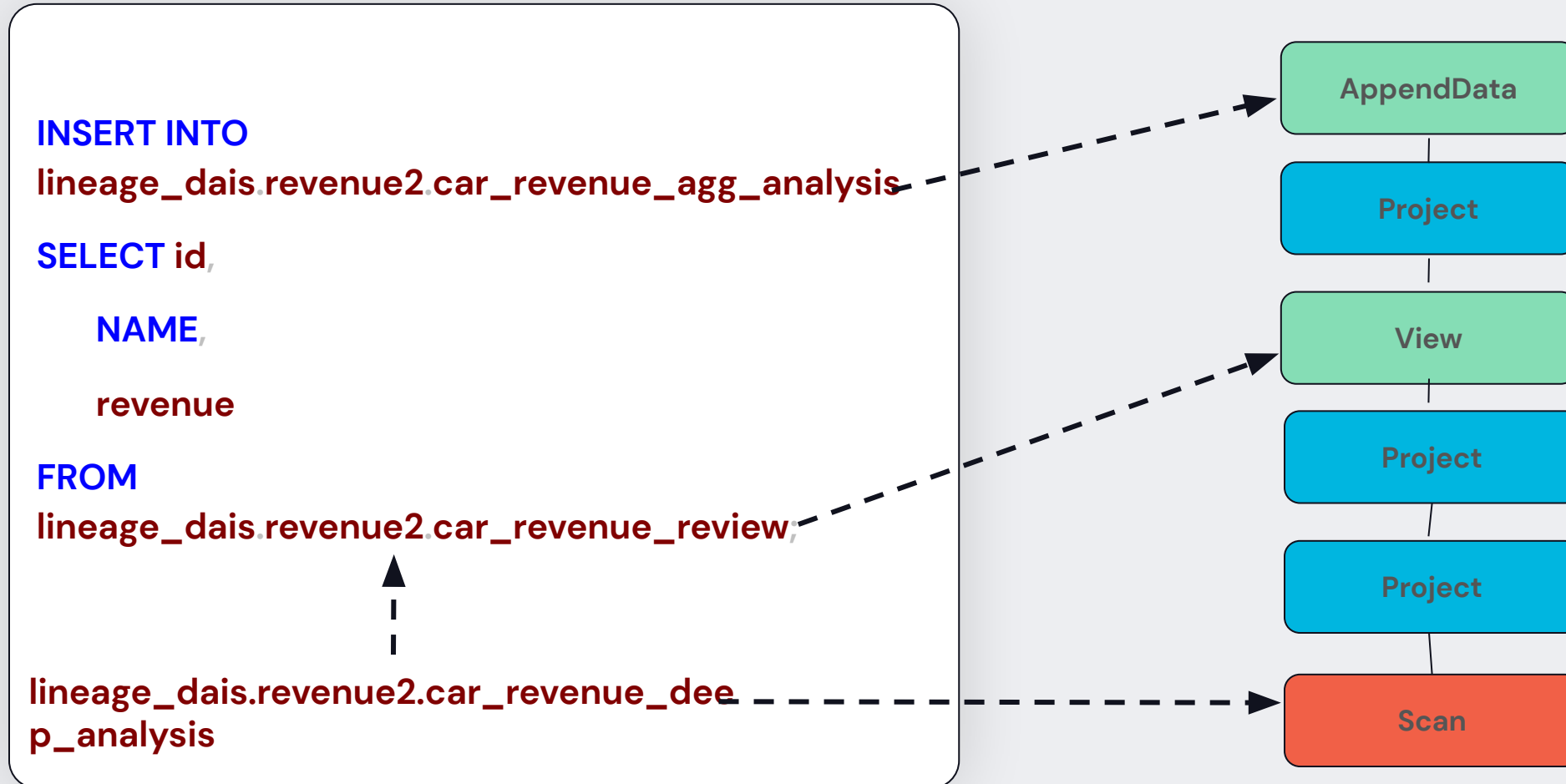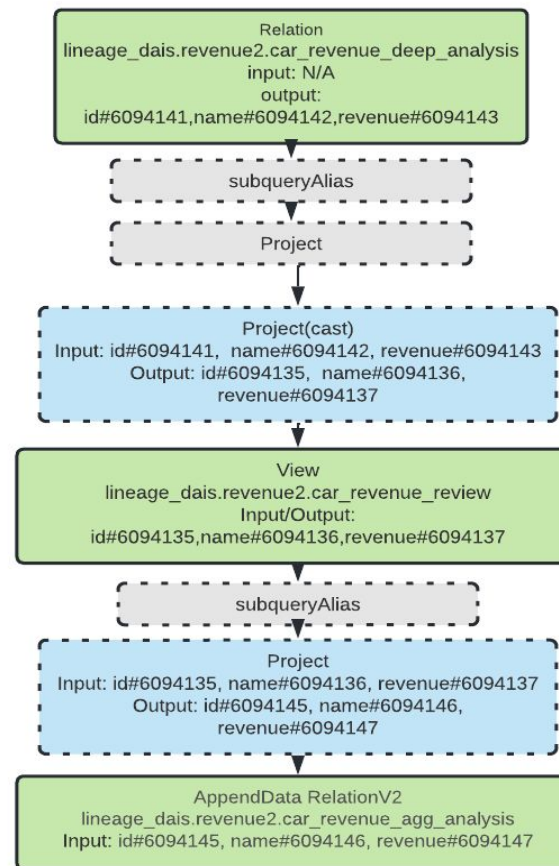
- **Expression:** An expression represents a new value, computed based on input values
- **Attribute**: A column of a dataset (e.g id) or a column generated by a specific data operation

# Tree: Abstractions of Users' Programs

**INSERT INTO**
lineage_dais.revenue2.car_revenue_agg_analysis

**SELECT id**,

    **NAME**,

    revenue

**FROM**
lineage_dais.revenue2.car_revenue_review,

lineage_dais.revenue2.car_revenue_deep_analysis

AppendData

Project

View

Project

Project

Scan

# Analyzed Logical Plan

**INSERT INTO**
lineage_dais.revenue2.car_revenue_agg_analysis

**SELECT id,**

    **NAME,**

    revenue

**FROM**
lineage_dais.revenue2.car_revenue_review;

lineage_dais.revenue2.car_revenue_deep_analysis

# Capture Lineage

- **Criteria**
  - Capture table name with 3L namespace (catalog.schema.table)
  - Capture persistent view with source table dependence
  - Capture target column and its source columns for column level lineage
  - Don't capture lineage unless the command executed successfully
  - Shouldn't throw exception during command execution

# Capture Lineage

- **Challenges**
  - No existing data for table / column level heuristic validation
  - No generic solution to support all Spark use cases
  - Support across different Spark versions

# How to capture the entity lineage?

# Entity Lineage

- Capture all entity artifacts lineage within the lakehouse
  - Initial entities: notebook, Job and DBSQL Dashboard
- Two main use cases:
  - **Transformer** : Notebook and Job consumes bunches of input tables and produces output tables
  - **Consumer**: Notebook and DBSQL dashboard only consumes tables and produces artifacts(e.g BI Dashboard)

# How to capture the entity lineage

- Capture data lineage along with entity type + entity ID
- Capture the lineage even for the read query
- Respect entity ACL
- Link back to the original entity

# Details

# #1 How to represent entities in the lineage

- Lineage is like a graph
- Entity type + Unique ID / name + Scope to uniquely identify the entity in the lineage graph
  - Workspace level entity: Workspace ID + Entity Type + Entity ID
  - Account level entity (shared across workspace): Metastore ID + Table + Column

# #2 What is the approach to collect lineage

| Pattern | Description | Example | Key Benefit | Key Challenge |
|---------|-------------|---------|-------------|---------------|
| Manual linked by User | Manual added and described how datasets are linked | | | Does not scale |
| Inferred from SQL / SQL parsing | Programmatic extracting lineage with SQL dialect | https://github.com/uber/query parser | Accurate, supports all sql dialect | SQL is easier, but can't support Spark which has multi lang API |
| **Auto capture during runtime** | **Language agnostic, auto captured in query execution phase** | **Lineage in Unity catalog** | ● **Correct**<br>● **Handle multi-lang** | |

# #3 What is the latency to update lineage



| ETL Batch Update | Neal realtime Update |
|---|---|

**ETL Batch Update**

Lineage Source → Crawler → Lineage graph

Scheduler

Pros
- The graph could be fully rebuilt

**Neal realtime Update**

Lineage Source → Message queue → Lineage graph

Pros
- The graph only keeps 1 level upstream / dowstream of entity
- Provide lineage timely to end users

# #4 What is the permission model of lineage

- Table Lineage leverages Unity Catalog permission model to view lineage
- Entity(notebook, workflow, dashboard) Lineage subjects to the entity permission ACL model

# Roadmap

# In Context Lineage

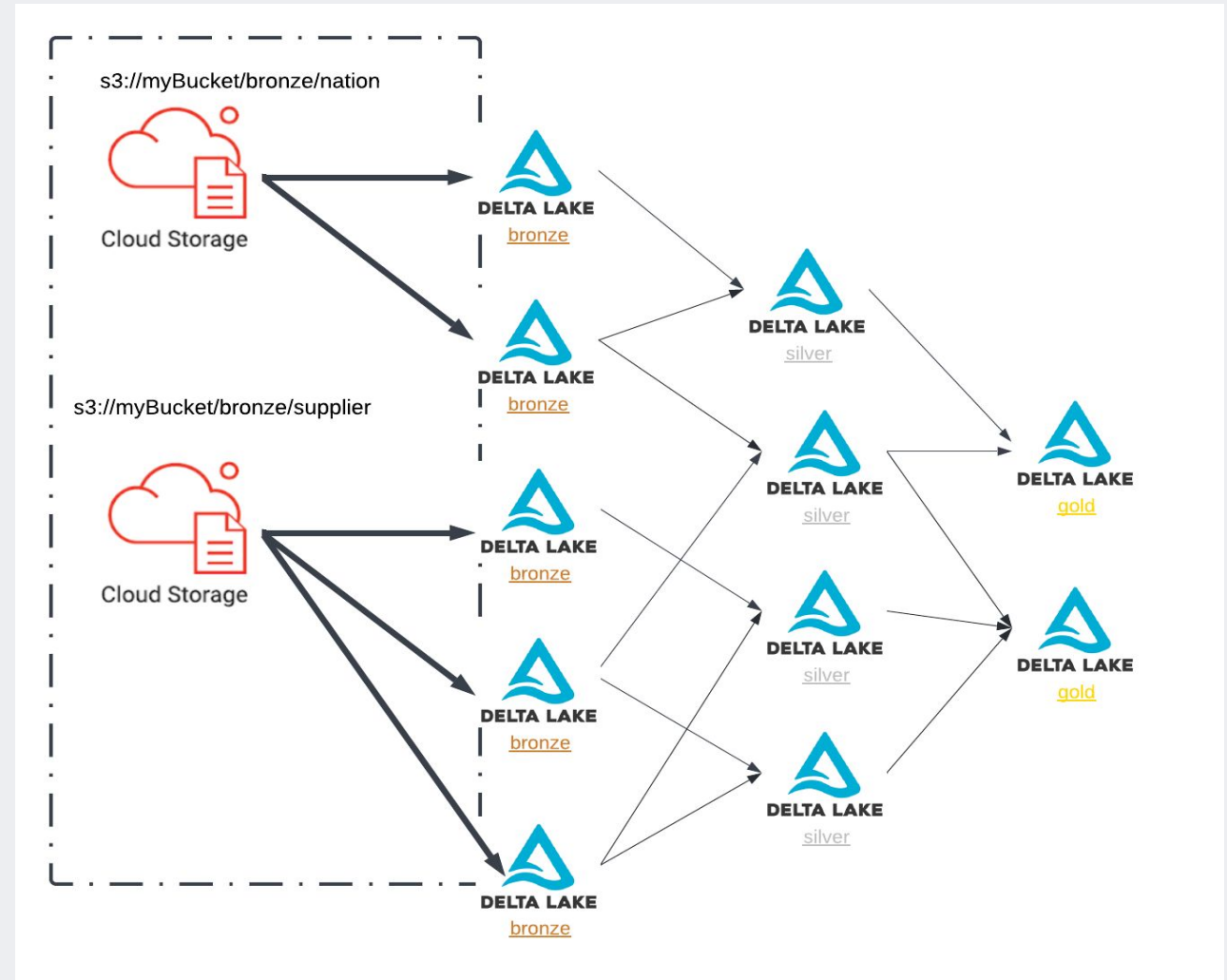- Lineage should be context aware.
- Lineage in job page
- Debugging and impact analysis

# Tag Propagation

- Propagate tag through lineage automatically

# Files Lineage

- Capture file lineage of the `first-mile ETL` which the ETL writes to cloud storage first
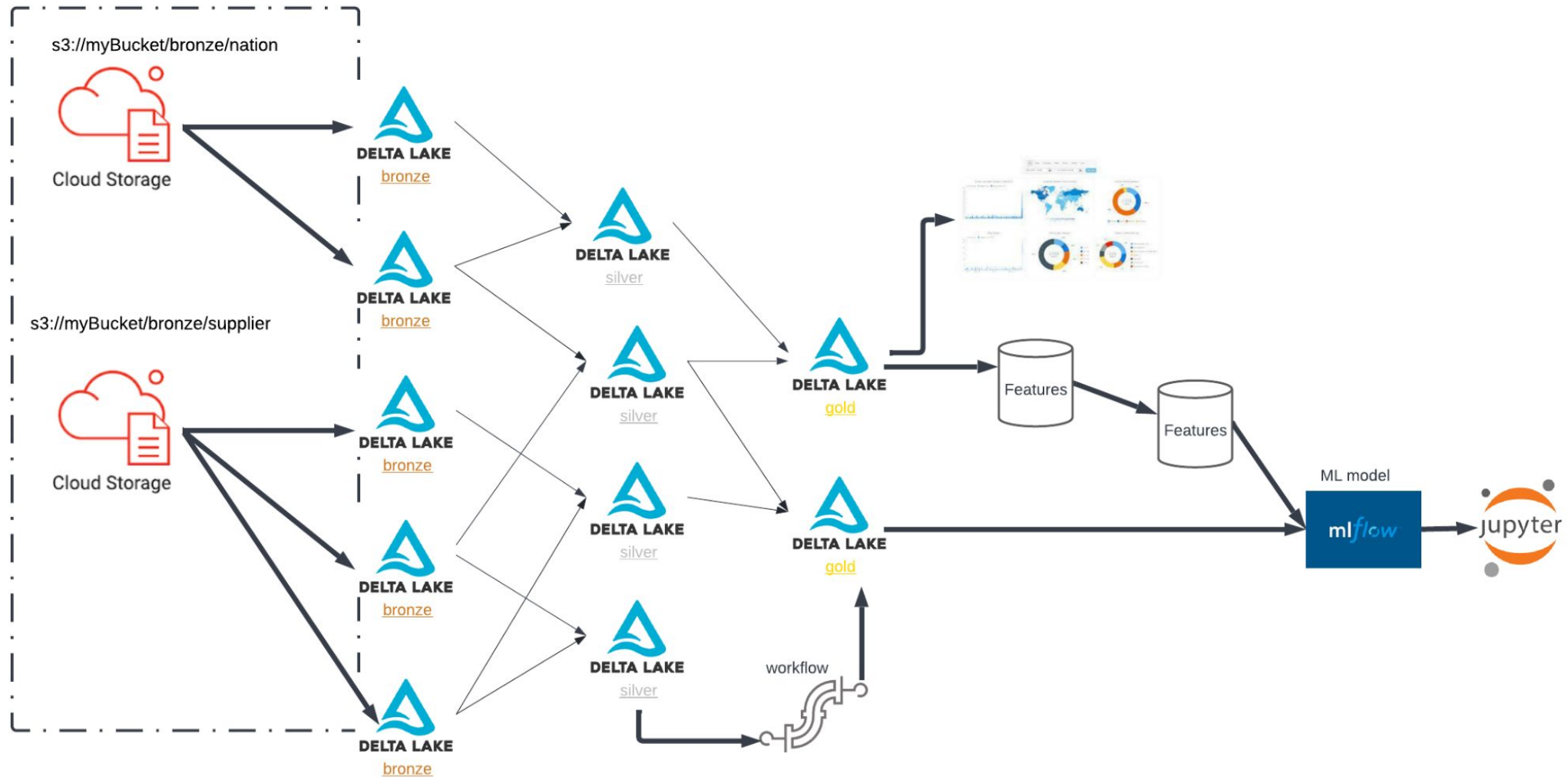
# ML Lineage

- Provide end-to-end lineage from dataset to features to MLflow models
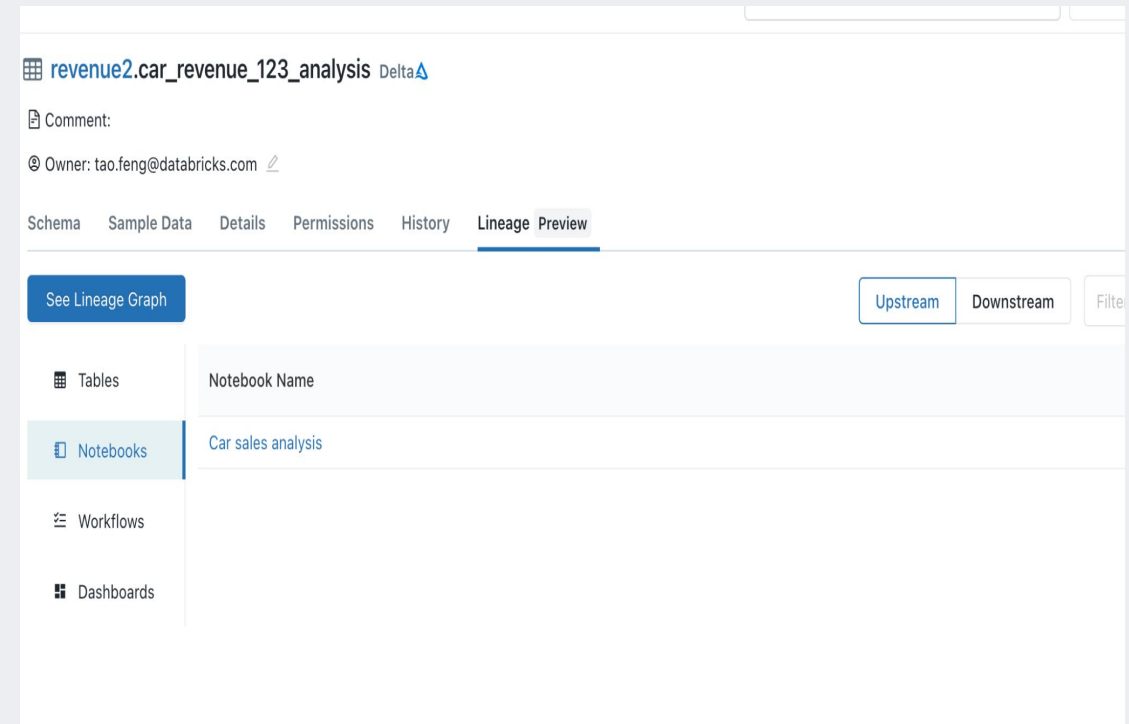
# Lakehouse E2E Lineage

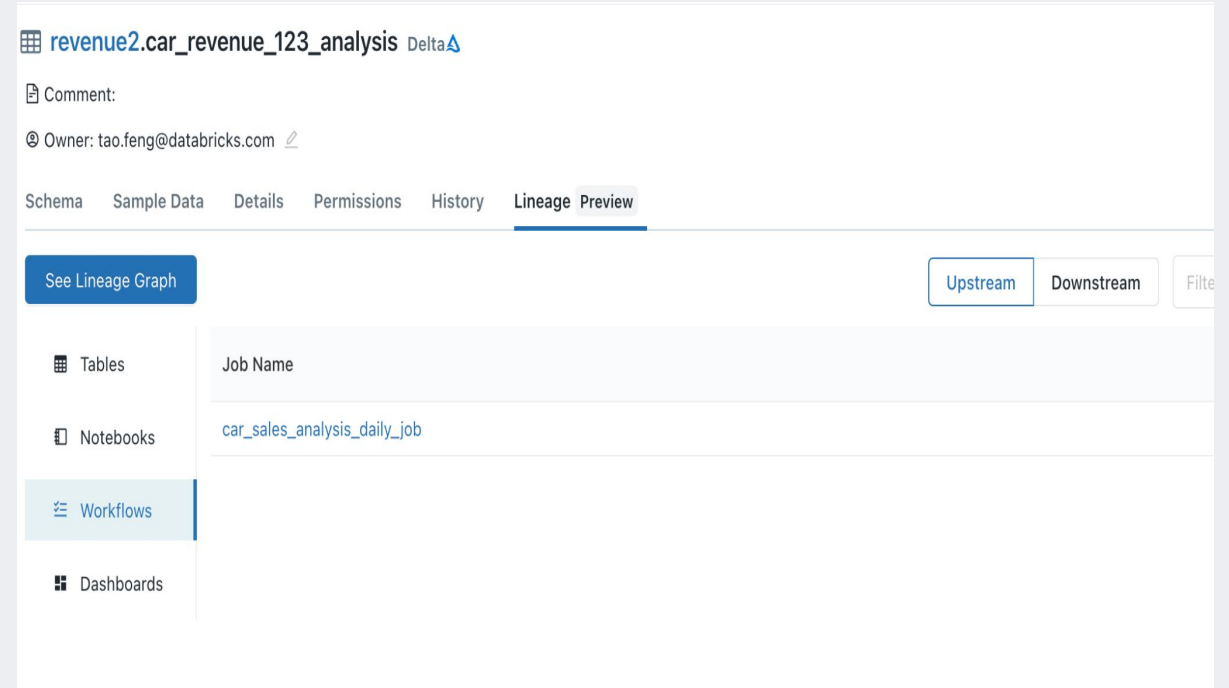# Appendix

**DATA+AI**
SUMMIT 2022

# Notebook Lineage

- Currently table centric
- Capture upstream / downstream notebooks of a given table
- Respect Notebook ACL
- Could direct back to the original notebook

**DATA+AI**
SUMMIT 2022

# Workflow Lineage

- Capture upstream / downstream workflow / job of a given table
- Could surface task level dependency if needed
- Respect workflow / job ACL
- Could direct back to workflow / job page

**DATA+AI**
SUMMIT 2022

# Dashboard Lineage

- Capture downstream dashboard of a given table

- Respect dashboard ACL

- Easy to find out impacted dashboard if the table has issues

**DATA+AI**
SUMMIT 2022