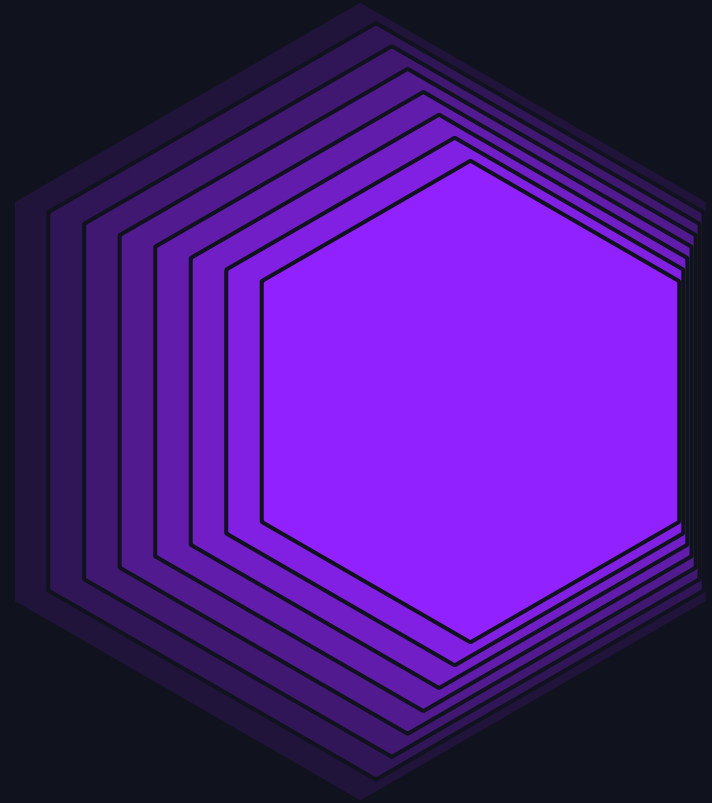# RISE OF THE MEDALLION MESH

Franco Patano, Databricks, Strategic Data & AI Advisor

1

# Top concerns in data and AI modernization*

**For leaders and practitioners across industries**

- Data security, governance, and quality

- Performance, scalability, and total cost of ownership

- Skills and expertise gap

- Data integration and migration

- Provide diverse tooling to support variety of users on a single version of the truth for all use cases

# whoami

**Franco Patano**
Strategic Data and AI Advisor

- Background in ETL, Data Warehousing, BI, and Analytics
  - MSBI Stack (SQL Server, Reporting Services, Integrations Services, Analysis Services)
  - Informatica
  - Tableau
  - Cognos
  - Excel and VBA
- Grew up in Enterprises in Chicago Area
  - Career Education Corporation (Perdoceo Education Corporation)
  - Wintrust
  - JLL (Jones Lang Lasalle)
- Dreamed of joining a startup and changing the world
  - Joined Databricks in 2019 to change the world with Lakehouse!

- Likes
  - Time Travel fiction (and non-fiction)
    - Back to the future
    - Doctor Who
    - Bill and Ted
    - Predestination
  - Making Pizza from Scratch
  - EDM
    - Daft Punk, Alesso, Cazzette, DJ Caffeine
  - Rocky, The Princess Bride, and The Fifth Element are the most EPIC films of all time
- Dislikes
  - Data Warehouses
    - Ever since I first became an analyst, all I wanted to do was get rid of that pesky data warehouse
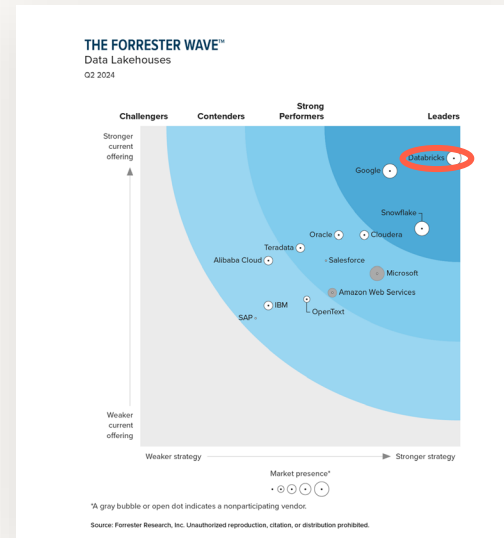
# Recognized as a leader in the industry

**Gartner MQ**
**Leader**
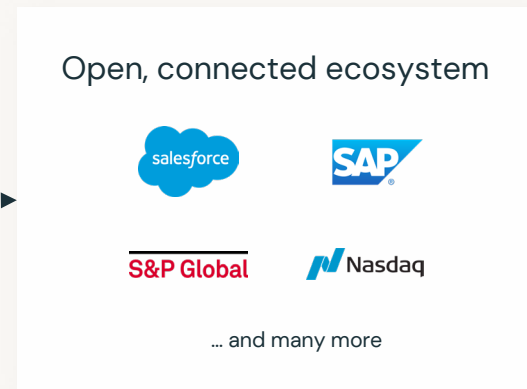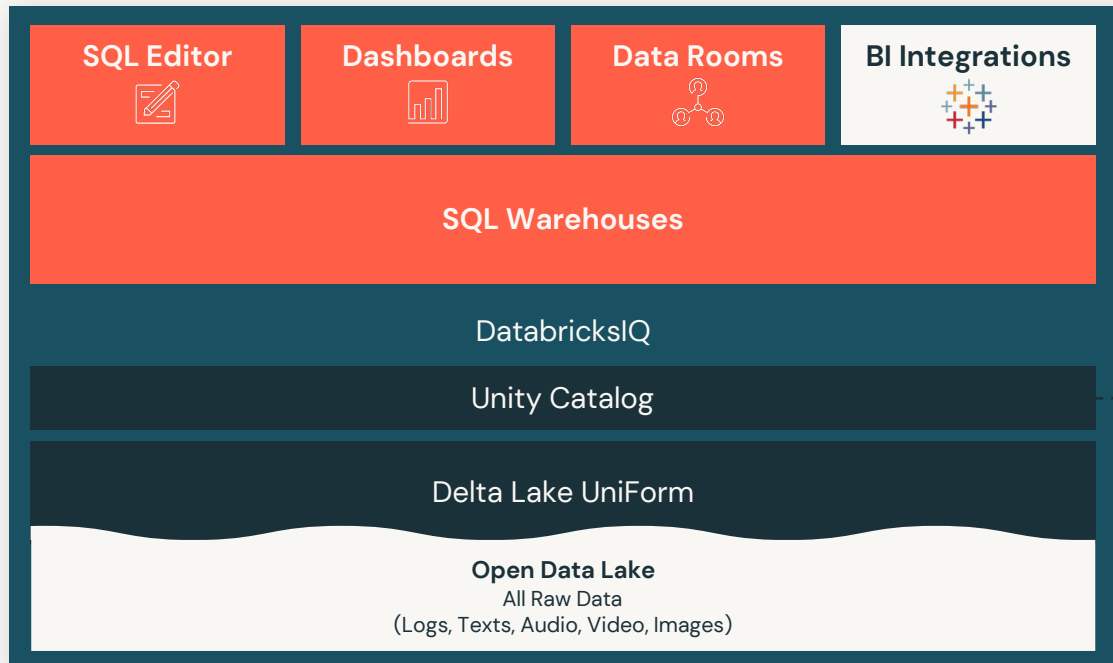
**Database Management Systems**



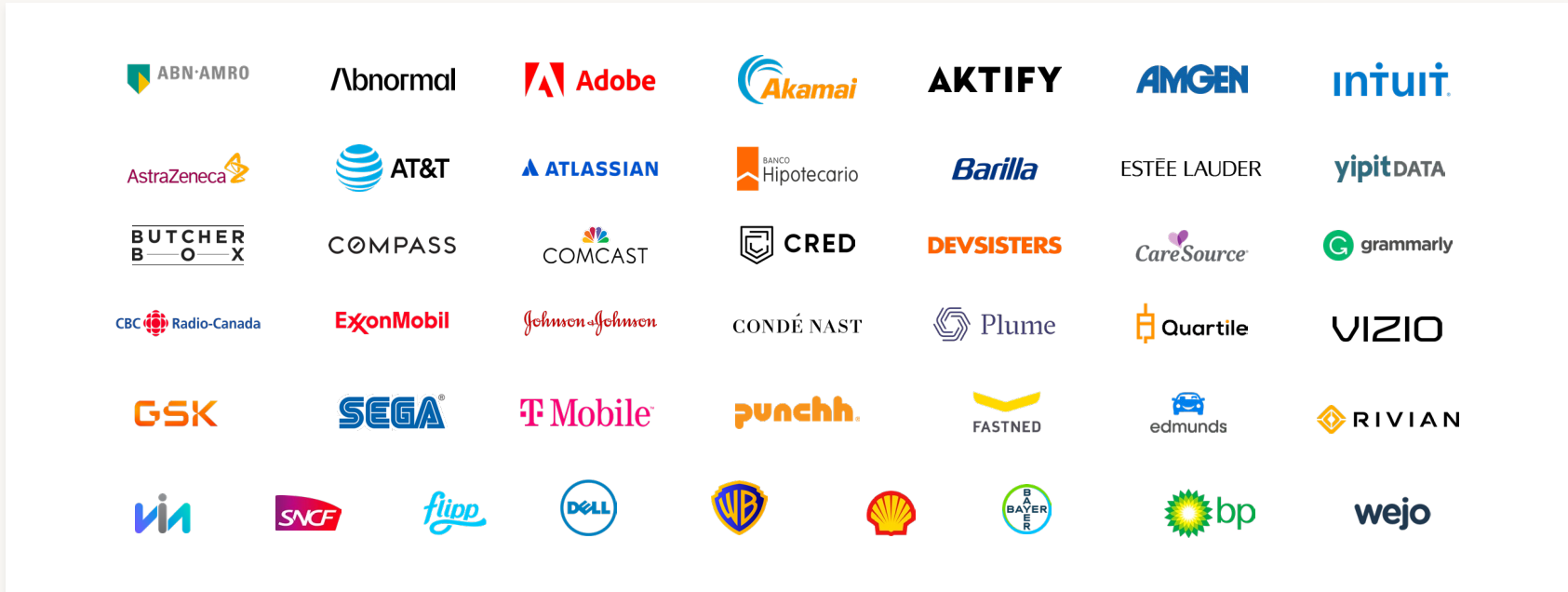**Forrester Wave Leader**

**Data Lakehouses**

# Databricks SQL

## Intelligent data warehousing on the Data Intelligence Platform

# Trusted by organizations of all sizes

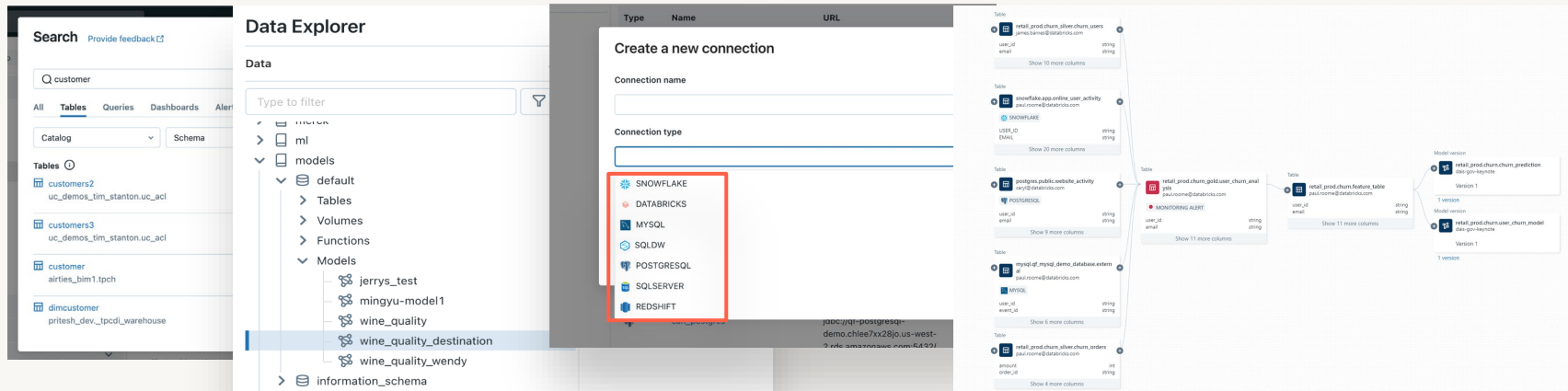4,600+ Databricks SQL customers across industries

# DATA SECURITY, GOVERNANCE, AND QUALITY

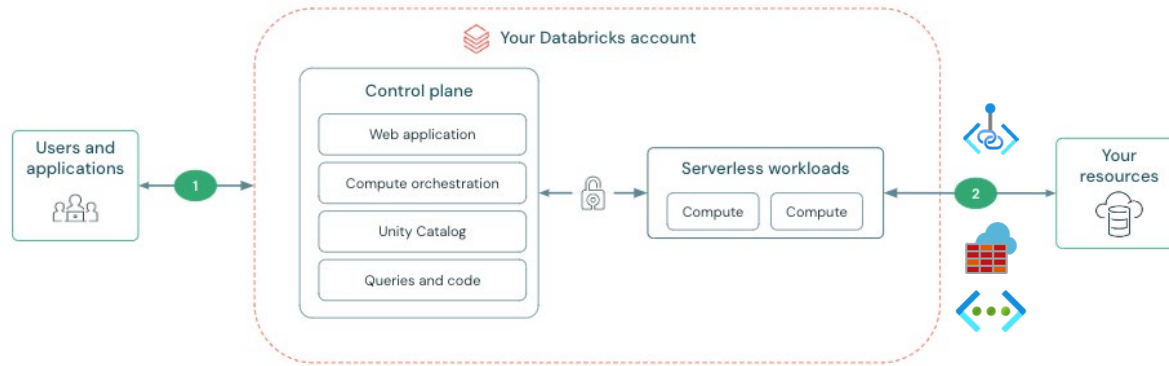# Governed and secured by Unity Catalog

## Governance for all your data and AI assets



Simplified **data discovery**, **governance**, **federation**, **lineage**, and **compliance** with enhanced **security** and **auditing** with Unity Catalog and Databricks SQL
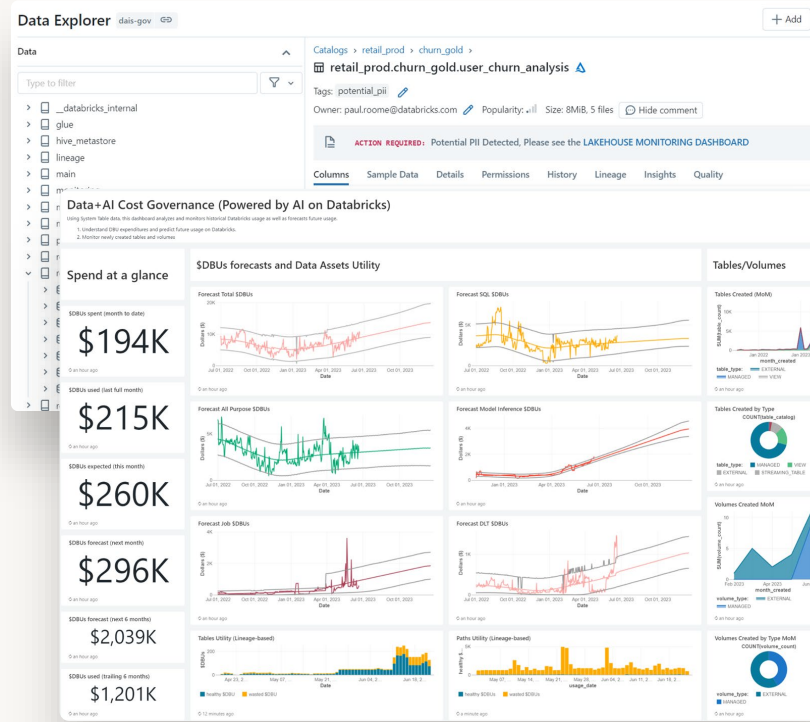
# Simple security with private connectivity

## Private Link or firewall support for stable VNET IPs



- Reduce data exfiltration risk
- Improved performance by minimizing network latency

- Meet compliance requirements for sensitive workloads
- PrivateLink included with Serverless SQL Warehouse cost on Azure

# Simple AI-powered monitoring and observability

- Receive **proactive alerts** for quality issues with data and ML models

- Access **real-time data lineage** down to the column level for efficient root cause analysis and impact assessment

- Utilize **auto-generated** dashboards to easily share data and ML quality reports with stakeholders

- Achieve complete data and AI **observability** with operational intelligence for billing, auditing, lineage, and more

# Simple data monitoring

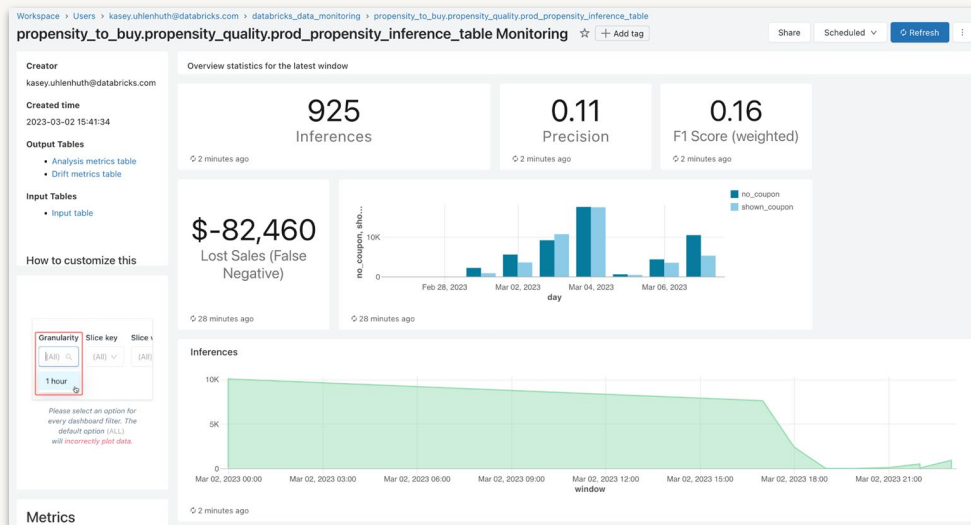Lakehouse monitoring for reliable, insightful, and simple data-to-AI-BI pipelines

**Simple:** Log inference tables automatically, and generate metric tables and SQL dashboards.

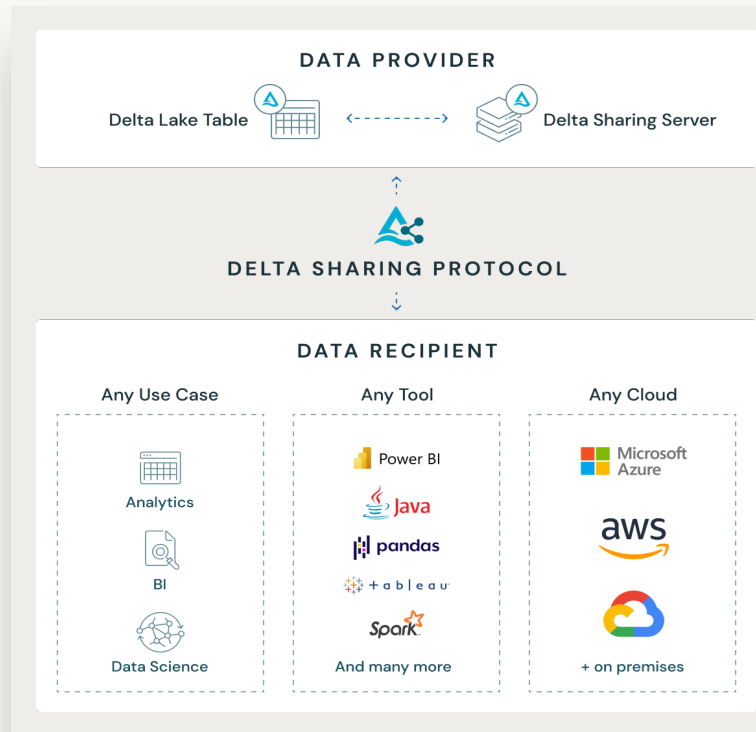**Proactive:** Automate alerts on table quality and custom metrics, and diagnose data or model issues.

**Integrated:** Track end-to-end lineage in the Unity Catalog for training data, feature tables, models, and inference logs, for simpler governance.

# Simple and open data sharing

- **Avoid vendor lock-in** with open-source Delta Sharing for seamless data sharing across clouds, regions, and platforms, without replication

- Share **more than just data:** notebooks, ML models, dashboards, applications

- Explore and monetize data products through an **open marketplace**

- Collaborate securely on sensitive data with **scalable data clean rooms**
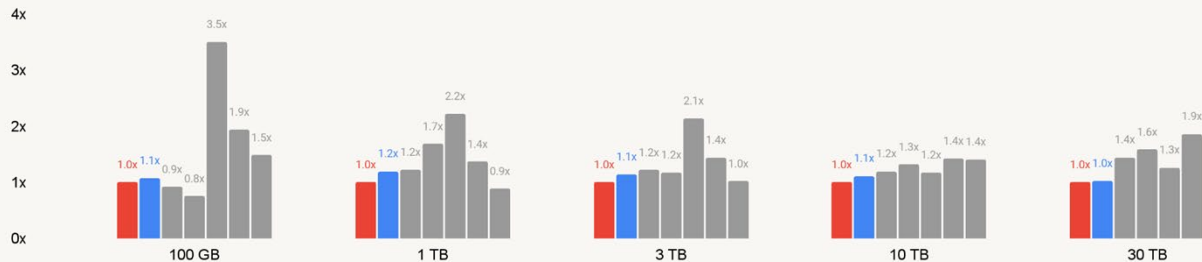


**DATA PROVIDER**

Delta Lake Table ←--------→ Delta Sharing Server

**DELTA SHARING PROTOCOL**

**DATA RECIPIENT**

| Any Use Case | Any Tool | Any Cloud |
|---|---|---|
| Analytics | Power BI | Microsoft Azure |
| BI | Java | aws |
| Data Science | pandas | Google Cloud |
| | tableau | + on premises |
| | Spark | |
| | And many more | |

# PERFORMANCE, SCALABILITY, AND TOTAL COST OF OWNERSHIP

# World-class performance & TCO

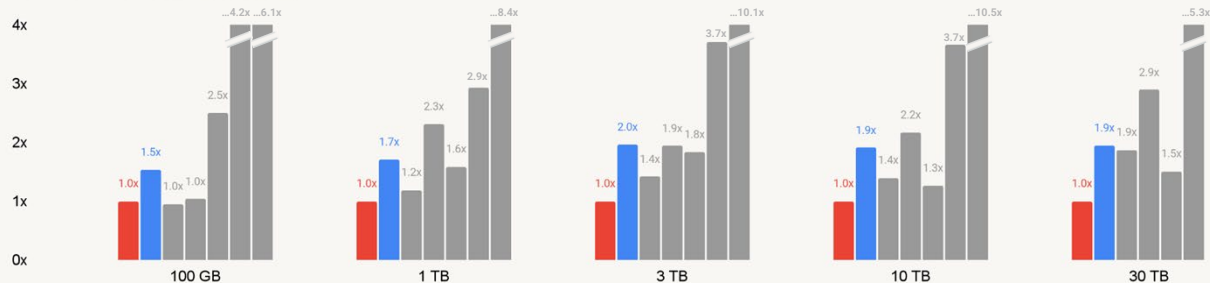## Meets or beats the price/performance of major CDWs across scales!



TPC-DS Elapsed Time (Lower is Better)

TPC-DS Total Cost (Lower is Better)

PERF

COST

Legend:
- Databricks
- CDW 1
- CDW 2
- CDW 3
- CDW 4
- CDW 5
- CDW 6

# World-class performance for BI

High concurrency with low latency



Lower is better

Latency msec / Concurrent Users

- Databricks
- Leading Multi-cloud CDW

Adevinta

"Our analysts rely on Databricks SQL to derive **business intelligence . . . we have 30% better performance** and have reduced costs by 20% on average."

—Allard de Boer

Global Director of Analytics,

Adevinta

# Simple performance with Predictive Optimization

**AI-optimized Delta table layouts for best price–performance**

Runs OPTIMIZE, VACUUM, ANALYZE, Liquid clustering

AI model prioritizes tables to maximize ROI

Out-of-box observability with system tables

**Customer storage costs**
(actual data from Preview customer)

"Databricks' Predictive Optimizations intelligently optimized our Unity Catalog storage, which **saved us 50% in annual storage costs** while **speeding up our queries by >2x**. It learned to prioritize our largest and most-accessed tables. And, it did all of this **automatically**, saving our team valuable time."

**—Anker**

# SKILLS AND EXPERTISE GAP

# Intelligent experiences with natural language
Powered by DatabricksIQ



**SQL Editor**

**Data Science / Engineers**



**Lakeview Dashboards**

**Analysts**



**Project Genie**

**Coming Soon!**

**Business Users**

# Assistant in SQL Editor

## Developers can create, explain and fix SQL code



**Assistant in Query Editor**

# Assistant in Lakeview

## Analysts can generate visualizations using natural language

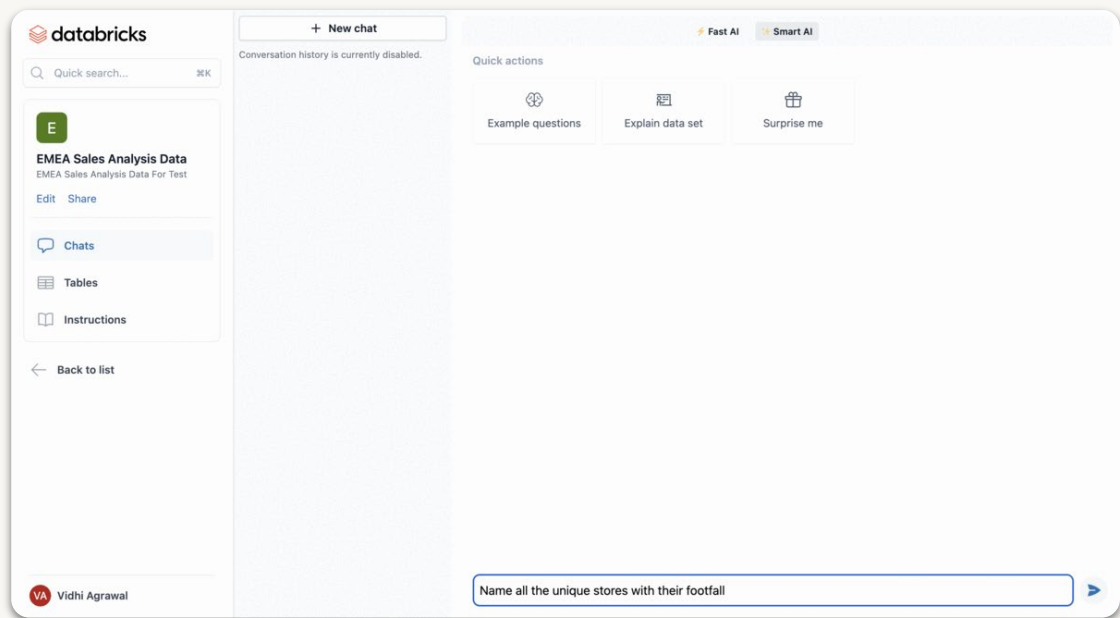

**Assistant
in Lakeview**

# Project Genie (Coming Soon)

## Business users can query and visualize data using natural language



**Project Genie**

# DATA INTEGRATION AND MIGRATION

# Simple migrations at your pace



## Federate
**with Lakehouse Federation**

Lower the barrier to entry and get started fast
Unity can secure and govern the sources, and track the consumption
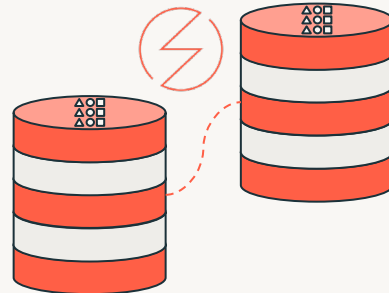
## Materialize
**with Materialized Views**

When the users and business find value, lower the stress on federated data sources with Materialization in the cloud
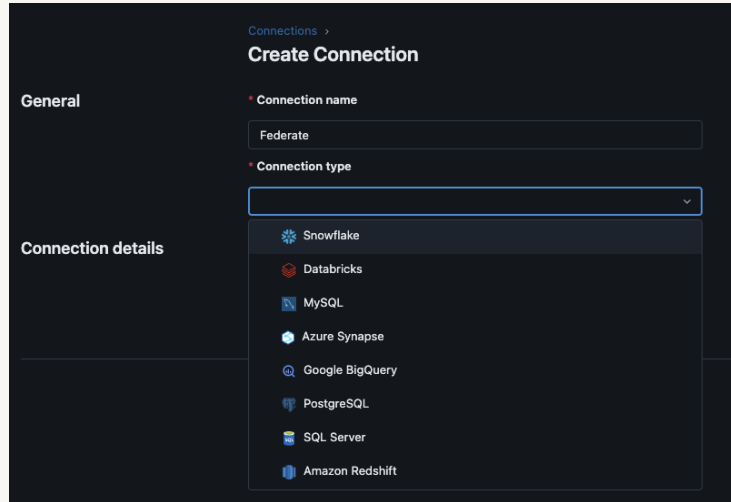
## Ingest
**with Change Data Capture**

When the business and technical folks are aligned, lower the costs and migrate with low latency and high throughput incremental CDC

**IN PRODUCTION TODAY**

**Announced at Summit**

**ALL OF THIS IS POWERED BY UNITY CATALOG**

# Simple data access with Federation



Discover, query, and govern *all* your data with **a unified view**, unified **engine** and **governance**—no matter where it lives
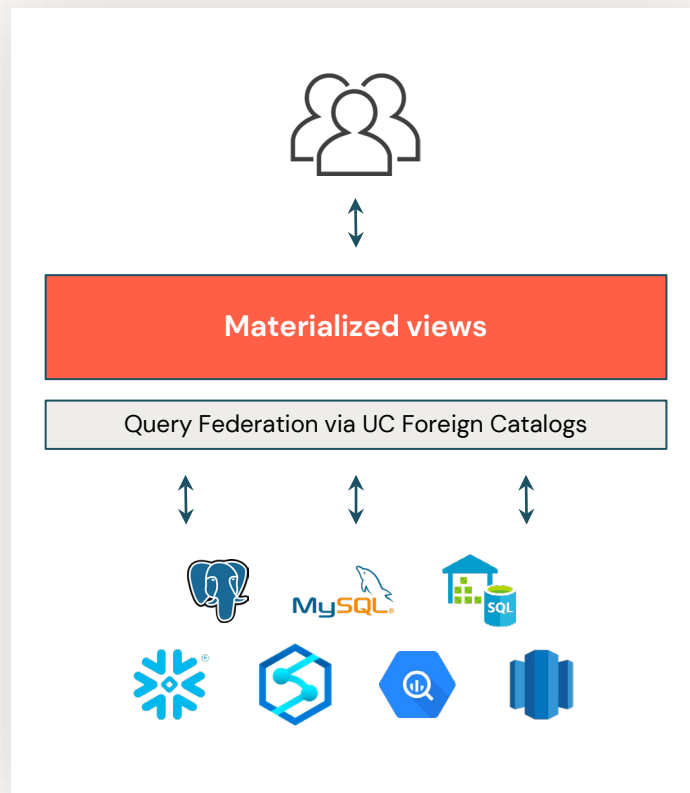
# Simple and fast performance

**Accelerating federated workloads**

**Federation 💛 Materialized views:**

- Consistent latency & concurrency for data outside of the Lakehouse

- Accelerate cross-source joins and complicated transformation logic

- Offload access to underlying databases via materialized views to avoid high/concurrent loads on operational databases



Materialized views

Query Federation via UC Foreign Catalogs

# Simple data ingestion

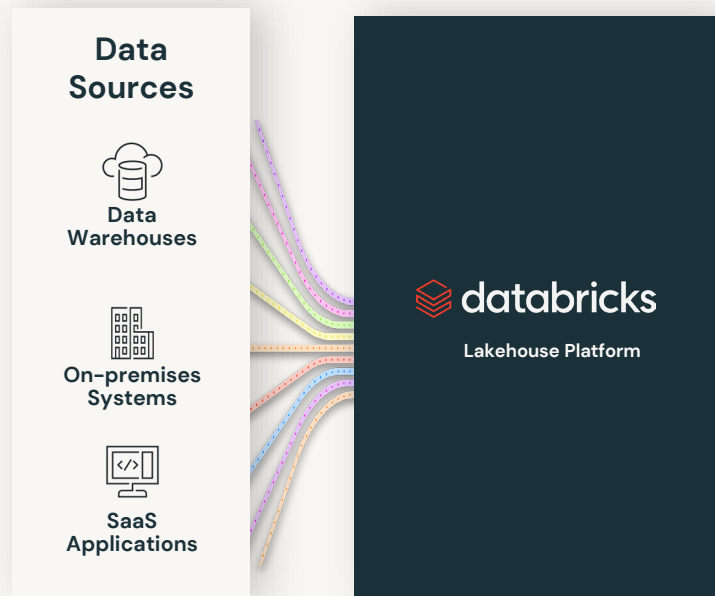Built-in connectors for common database and SaaS sources

**Quickly and reliably replicate data across on-prem, cloud databases, and data platforms into the Lakehouse**

Fully integrated into the Lakehouse

Low-maintenance API or no-code UI
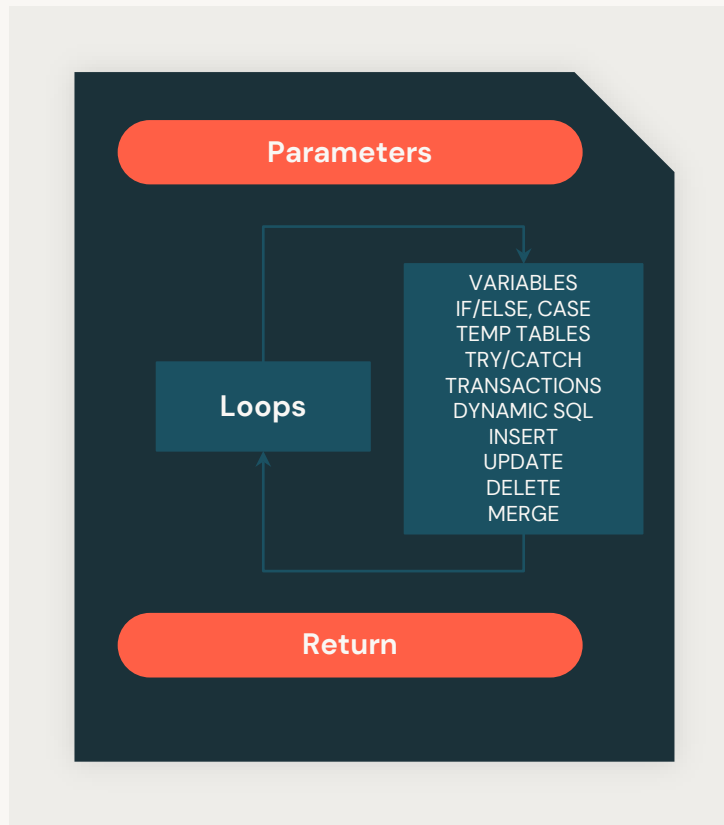
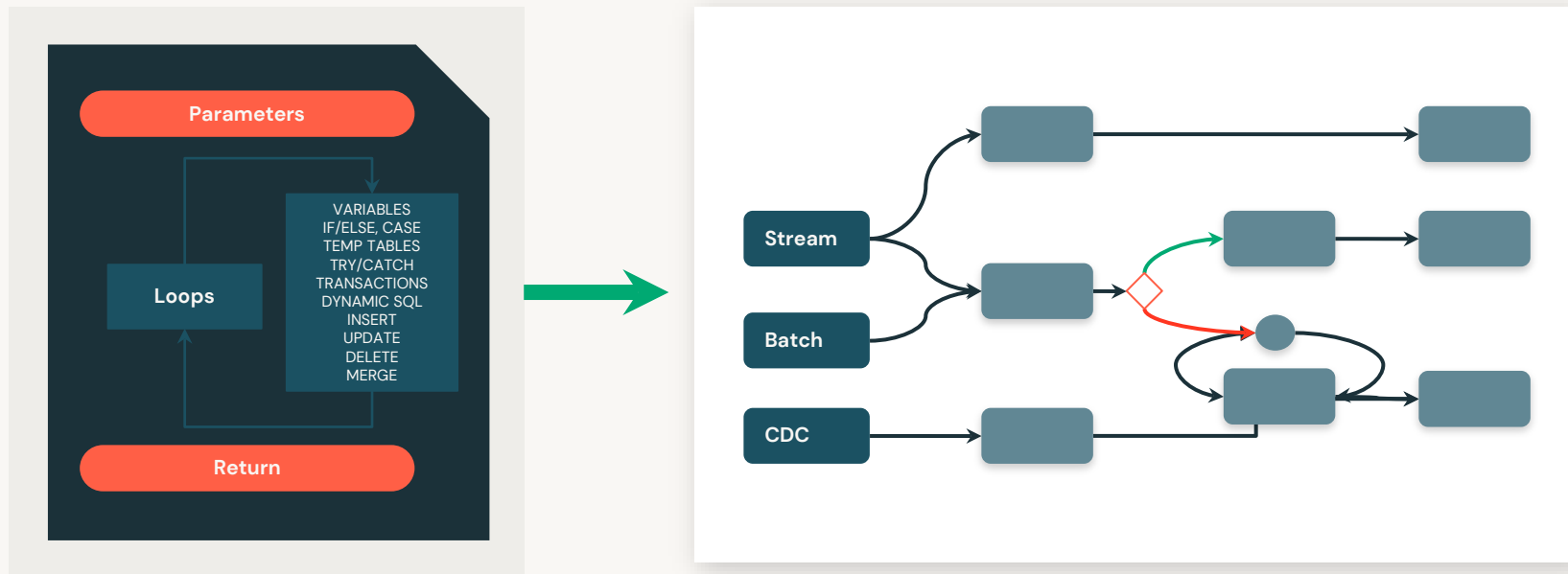Utilizing Arcion's advanced CDC technology

**Private Preview in Q1**



Data Sources

Data Warehouses

On-premises Systems

SaaS Applications

databricks
Lakehouse Platform

# Warehousing technical debt
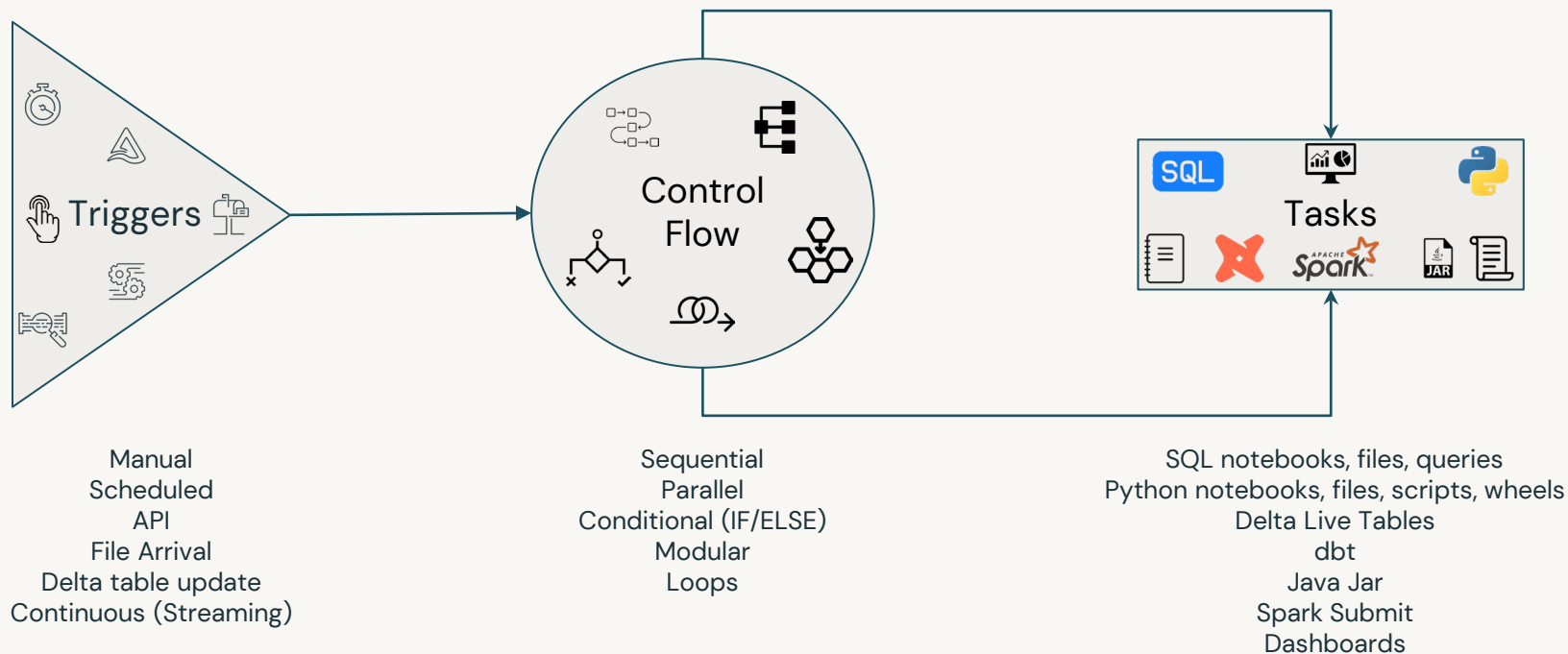
## Stored procedures

- **Why**
  - Performance & cost
    - Runs on warehouse, compared to external ETL tools

- **What**
  - Powerful SQL script that encapsulates complex and reusable database operations

- **Pain**
  - Difficult to debug
  - Costly in the cloud
  - Migration pain with vendor lock-in
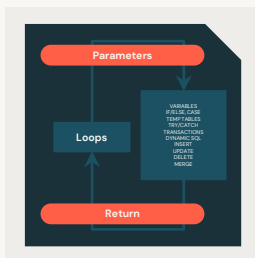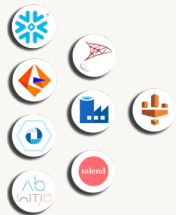  - Anti-Patterns in the cloud

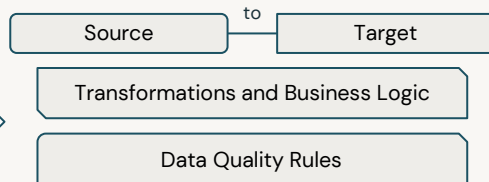# Refactor monolithic black box procedures to modularized Workflows

# Simple authoring of Workflows



| Triggers | Control Flow | Tasks |
|---|---|---|
| Manual | Sequential | SQL notebooks, files, queries |
| Scheduled | Parallel | Python notebooks, files, scripts, wheels |
| API | Conditional (IF/ELSE) | Delta Live Tables |
| File Arrival | Modular | dbt |
| Delta table update | Loops | Java Jar |
| Continuous (Streaming) | | Spark Submit |
| | | Dashboards |

# Simple Migrations with AI



Use AI to extract Source to Target mappings and transformations and business logic from existing ETL + Warehouse

| Source | to | Target |
|--------|----|--------|

Transformations and Business Logic

Data Quality Rules
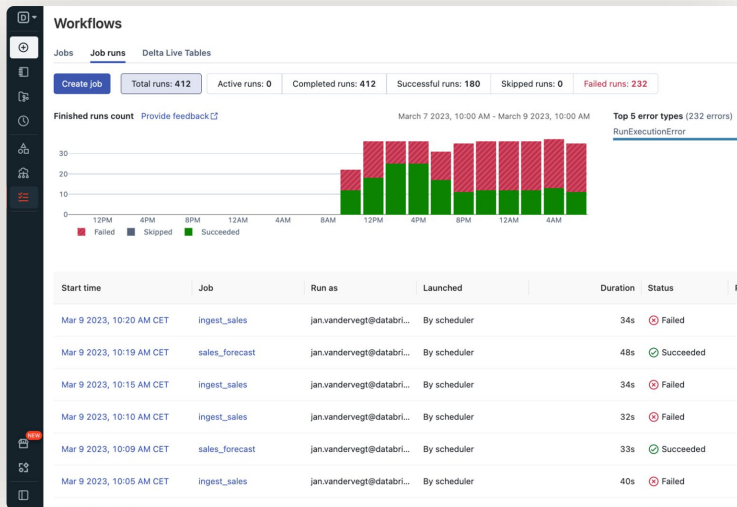
Map to a metadata driven framework

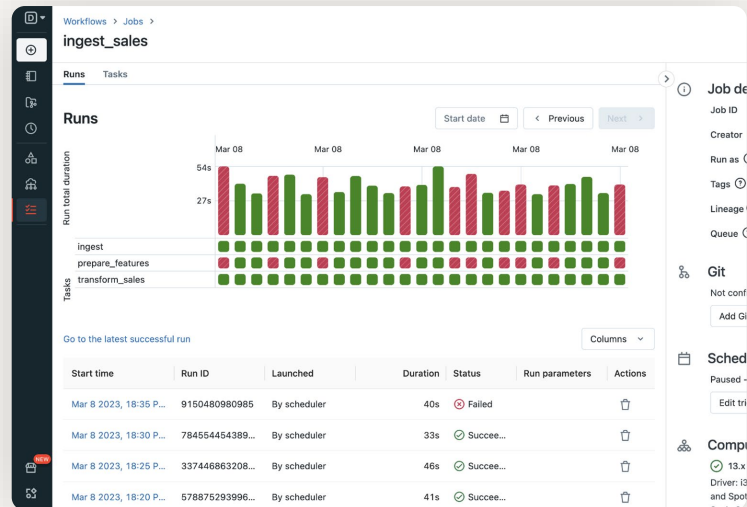Use AI to generate workflow DAGs using Python support for Databricks Asset Bundles

- Reduce risk of migration friction with LLM as judges to ensure transformation logic validity
- Accelerate time to value
- Backwards compatibility with Hive Federation, Delta Uniform or Delta Sharing to dependent systems
- Gain visibility with unified Governance and Orchestration
- Leverage existing Solutions Integrator frameworks

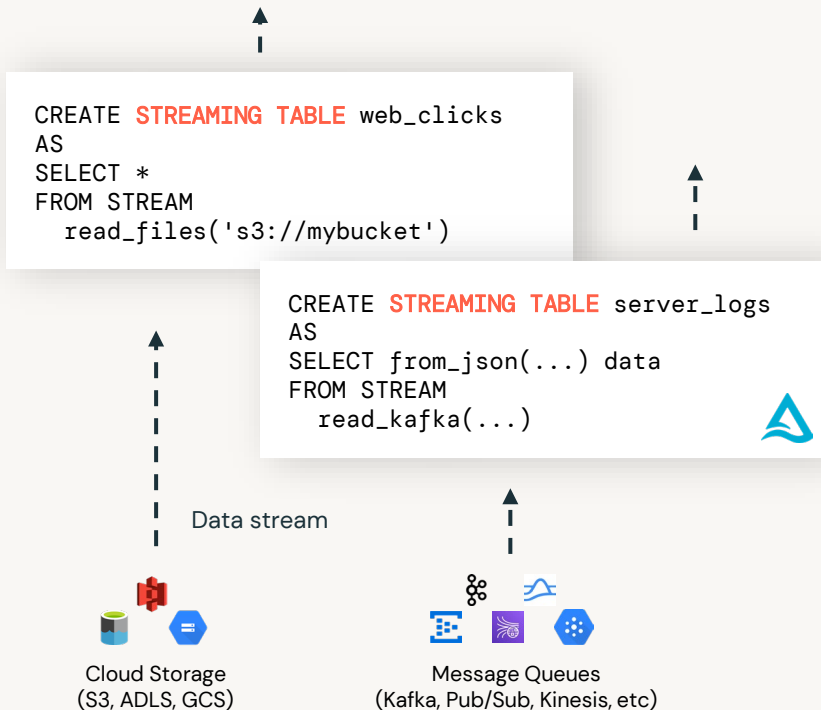# Simple real–time monitoring



A simple and intuitive monitoring UI provides real-time metrics and detailed analytics for every workflow run



Drill down to understand which tasks are failing and why. Troubleshoot issues before your customers are impacted

# Simple streaming with SQL

```
CREATE STREAMING TABLE web_clicks
AS
SELECT *
FROM STREAM
  read_files('s3://mybucket')
```

```
CREATE STREAMING TABLE server_logs
AS
SELECT from_json(...) data
FROM STREAM
  read_kafka(...)
```

Data stream

Cloud Storage
(S3, ADLS, GCS)

Message Queues
(Kafka, Pub/Sub, Kinesis, etc)

## Benefits:

1. **Unlock real-time use cases.** Ability to support real-time analytics/BI, machine learning and operational use cases with streaming data.

2. **Better scalability.** More efficiently handle high volumes of data via incremental processing vs. large batches.

3. **Enable more practitioners.** Simple SQL syntax makes data streaming accessible to all data engineers and analysts.

# Simple and fast BI with Materialized Views

```
CREATE MATERIALIZED VIEW customer_orders
AS
SELECT
  customers.name,
  sum(orders.amount),
  orders.orderdate
FROM orders
  LEFT JOIN customers ON
    orders.custkey = customers.c_custkey
GROUP BY
  name,
  orderdate;
```

Results are pre-computed and incrementally refreshed

customers
(Table)

orders
(Table)

## Benefits:

1. **Accelerate BI dashboards.** Much faster to query data that is pre-computed vs querying base tables.

2. **Reduce data processing costs.** MV results are refreshed incrementally avoiding the need to completely rebuild the view when new data arrives.

3. **Improve data access control.** More tightly govern what data can be seen by consumers by controlling access to base tables.

# Ingest and Transform with Medallion Mesh

* https://future.a16z.com/emerging-architectures-modern-data-infrastructure/

# Minimal data domain blueprint



* https://future.a16z.com/emerging-architectures-modern-data-infrastructure/

37

# Vertical view into the domains and layers

# Unified data modeling in Medallion Mesh

# UNIFIED PLATFORM FOR ALL USERS AND USE CASES

# Medallion Mesh on Data Intelligence Platform



AI
Hugging Face
OpenAI

Lakehouse Monitoring and Observability

Azure Data Factory | Databricks Workflows | Orchestration | CI/CD | MLOps/LLMOps

Unity Catalog | Data and AI Governance | Feature Store | Model Registry

DatabricksIQ | Discovery & Search | Intelligence Engine | Performance Optimization

**Databricks Data Intelligence Platform**

Databricks AI

Feature Engineering | mlflow | AutoML | Model Serving
Notebooks | Feature Serving
Vector Search

ETL and Process Engine

Lakehouse Federation | Auto Loader | Spark
Structured Streaming | Delta Live Tables | Photon

Data Warehousing
DB SQL

Cloud Storage

Lakehouse Federation For Snowflake Marketplace

Collaboration
Databricks Marketplace | Delta Sharing

bronze → silver → gold DELTA LAKE

Iceberg via Delta Uniform

Snowflake

Apps

Power BI

BI Tools

Operational Database

3rd party

Sensors and IoT (unstructured)

CDC

Streaming

RDBMS (structured)

Business Apps (structured)

Federation

Other clouds

Media (unstructured)

Ingest

Files / Logs (semi-.structured)

# Simple data understanding

## Real-time relational modelling visualization

Easily visualize **table-to-table** relationships leveraged by **BI tools** and **DatabricksIQ**

Modify and create PK/FK Constraints with UI

Visualize table to table relationships with a built-in Entity

Relationship Diagram (ERD)

Built on Unity Catalog

View PK, FK columns via the schema browser

# Simple and unified visibility into data and AI

- **Discover and classify** structured and unstructured data, files, notebooks, ML models, and dashboards at one place

- Consolidate and query data from **other databases and data warehouses** using a **single point of access,** without moving or copying the data

- Build better **understanding of your data estate** with automated lineage, tags and auto-generated data insights

- Boost productivity by searching, understanding and gaining insights from your data and AI assets, using **natural language**

# Simple help from AI in your SQL

## Write SQL to get insight from unstructured text data via LLMs

### SQL AI ANALYZE SENTIMENT

```
> SELECT ai_analyze_sentiment('I am happy');
  positive

> SELECT ai_analyze_sentiment('I am sad');
  negative
```

### AI SQL CLASSIFY

```
SELECT ai_classify("My password is leaked.", ARRAY("urgent", "not urgent"));
urgent

SELECT
  description,
  ai_classify(description, ARRAY('clothing', 'shoes', 'accessories', 'furniture')) AS category
FROM
  products
```

### SQL AI EXTRACT

```
> SELECT ai_extract(
    'John Doe lives in New York and works for Acme Corp.',
    array('person', 'location', 'organization')
  );
{"person": "John Doe", "location": "New York", "organization": "Acme Corp."}

> SELECT ai_extract(
    'Send an email to jane.doe@example.com about the meeting at 10am.',
    array('email', 'time')
  );
{"email": "jane.doe@example.com", "time": "10am"}
```

### SQL AI FIX GRAMMAR

```
SELECT ai_fix_grammar('This sentence have some mistake');
'This sentence has some mistakes'

SELECT ai_fix_grammar('She dont know what to did.');
'She doesn't know what to do.'
```

### SQL AI MASK

```
SELECT ai_mask(
  'John Doe lives in New York. His email is john.doe@example.com.',
  array('person', 'email')
);
[MASKED] lives in New York. His email is [MASKED].

SELECT ai_mask(
  'Contact me at 555-1234 or visit us at 123 Main St.',
  array('phone', 'address')
);
Contact me at [MASKED] or visit us at [MASKED]
```

### SQL AI SIMILARITY

```
SELECT ai_similarity('Apache Spark', 'Apache Spark');
1.0

SELECT
  company_name
FROM
  customers
ORDER BY ai_similarity(company_name, 'Databricks') DESC
LIMIT 1

Databricks Inc.
```

# Simplify your SQL

**Iterate faster with AI**

Assistant does all the heavy lifting

**Cell in focus**

Focus on a cell to turn it into a SQL Editor

Focus back to notebook mode

**Run Cells in parallel**

Use Run Now to execute SQL cells in parallel, no more waiting for execution



databricks

# Simplify DEVOPS with Repos

Repos API to automate CI/CD

# Simple Prompt Engineering UI with mlflow

Adjust text prompts for Gen AI models to elicit better responses

MLflow Prompt Engineering lets you compare and analyze many models and prompts, across many inputs.

# RISE UP AND UNIFY WITH THE MEDALLION MESH

fin

DATA+AI SUMMIT