databricks
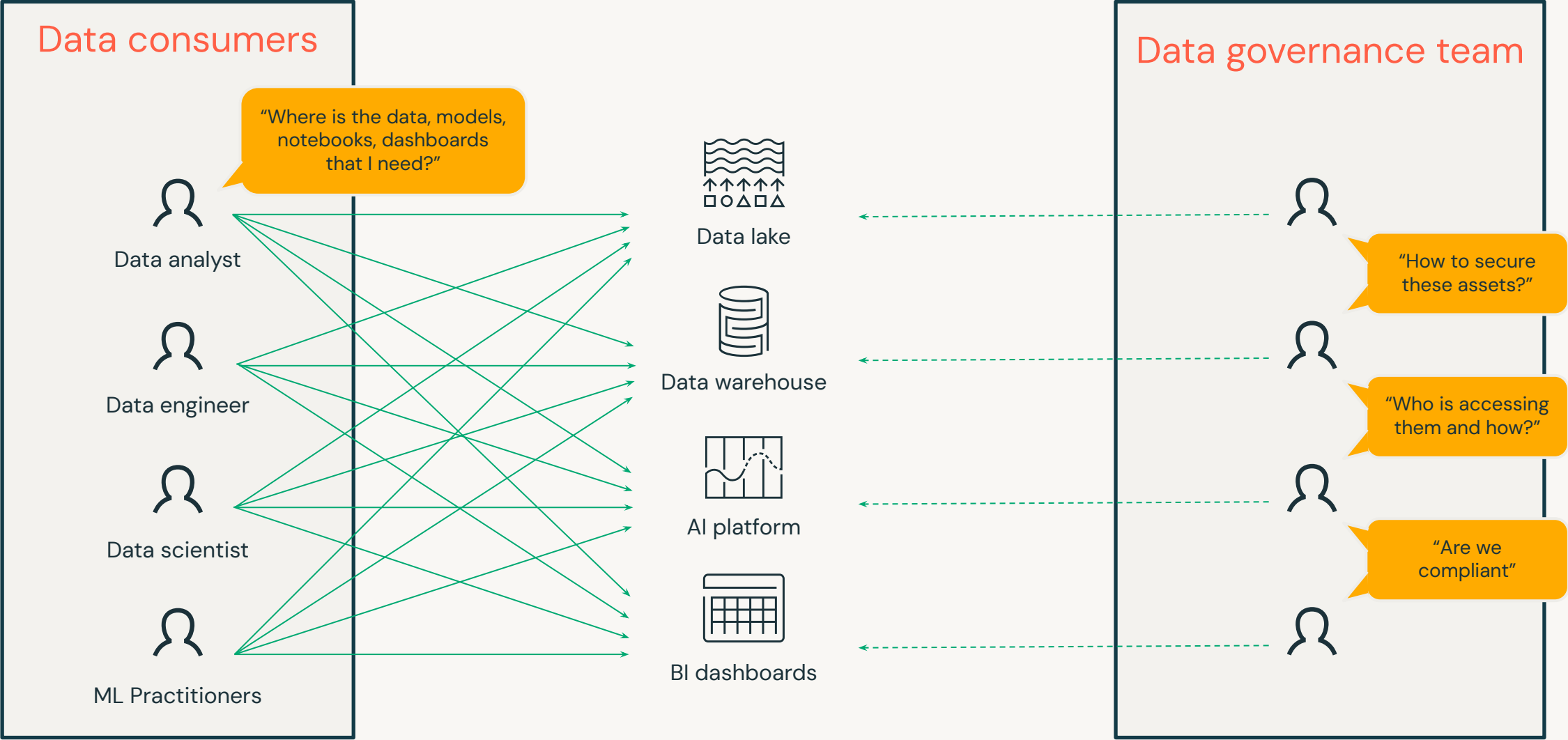
# DATA+AI WORLD TOUR

# What's new in Unity Catalog

Dalya Altaha, Solution Architect @ Databricks
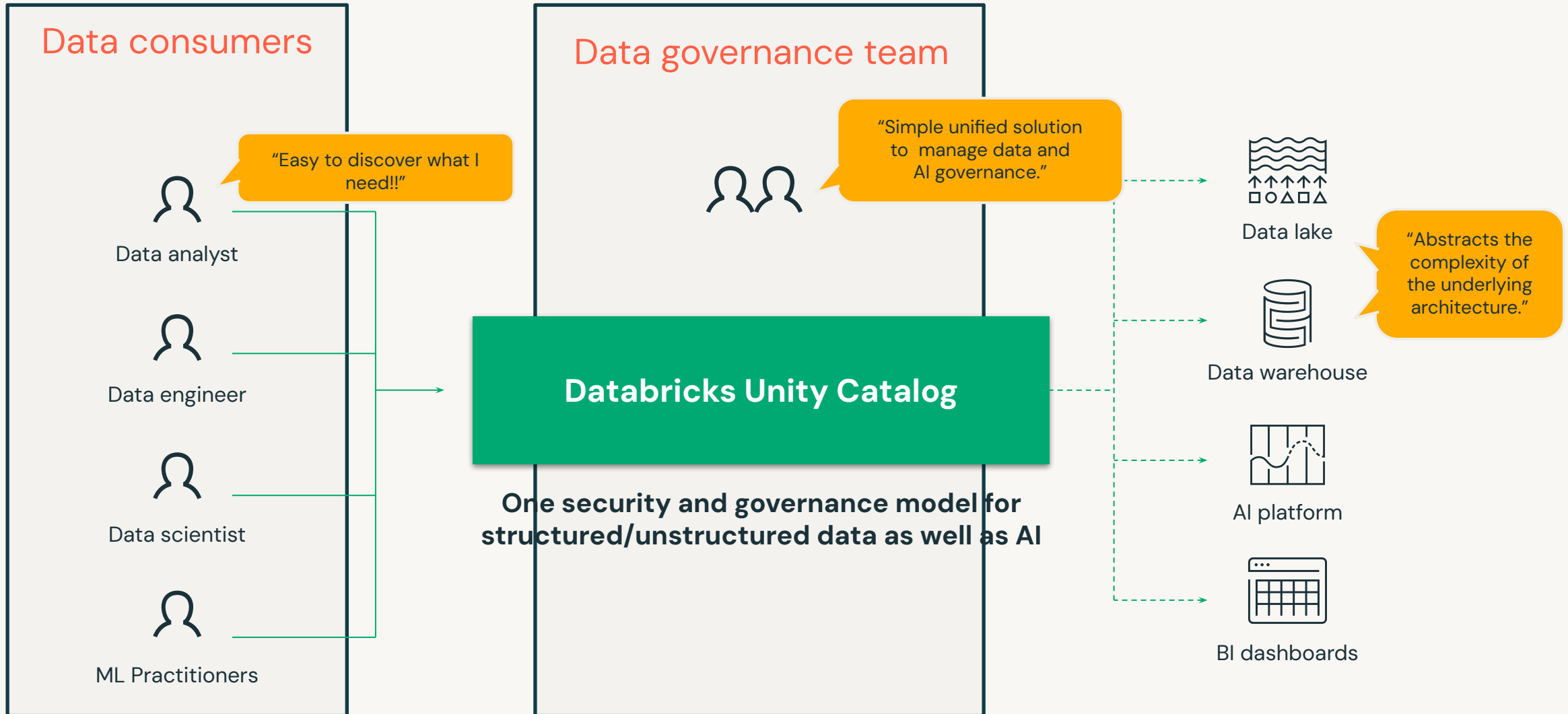Zen Van Gaever, Databricks Data Engineer @ delaware

# 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.

# Data and AI governance is complex

# Databricks unifies governance for data and AI

# Data and AI governance is growing in importance

> Organizations are finally realizing the value of data as an asset that needs to be protected, managed and maintained to increase asset value. "
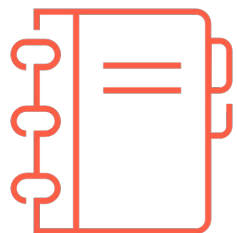>
> —IDC

> Organizations seeing the highest returns from AI have a framework for AI governance to cover every step of the model development process. "
>
> — The State of AI in 2022, McKinsey & Co.

> AI is now an enterprise essential, and as such, AI governance will join cybersecurity and compliance as a board-level topic. "
>
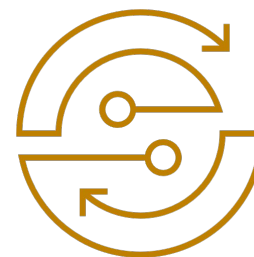> —Forrester, 2023 AI Predictions report

# Key elements of governance
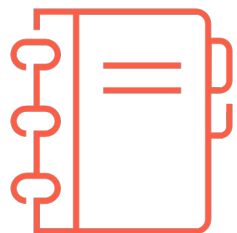
**Secure and auditable**
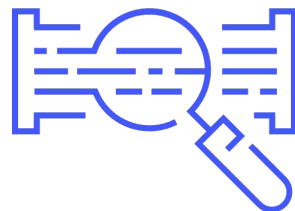
**Discoverable**

**Usable and performant**

**Accurate and high quality**

# How do we get there?

**Centralized access control for data & AI**

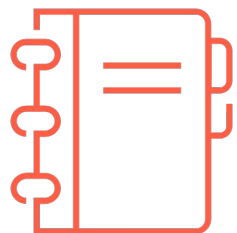**Built-in data search and discovery**

**Performance & scalability**

**Lineage & data quality**

# How do we get there?

**Centralized access control for data & AI**
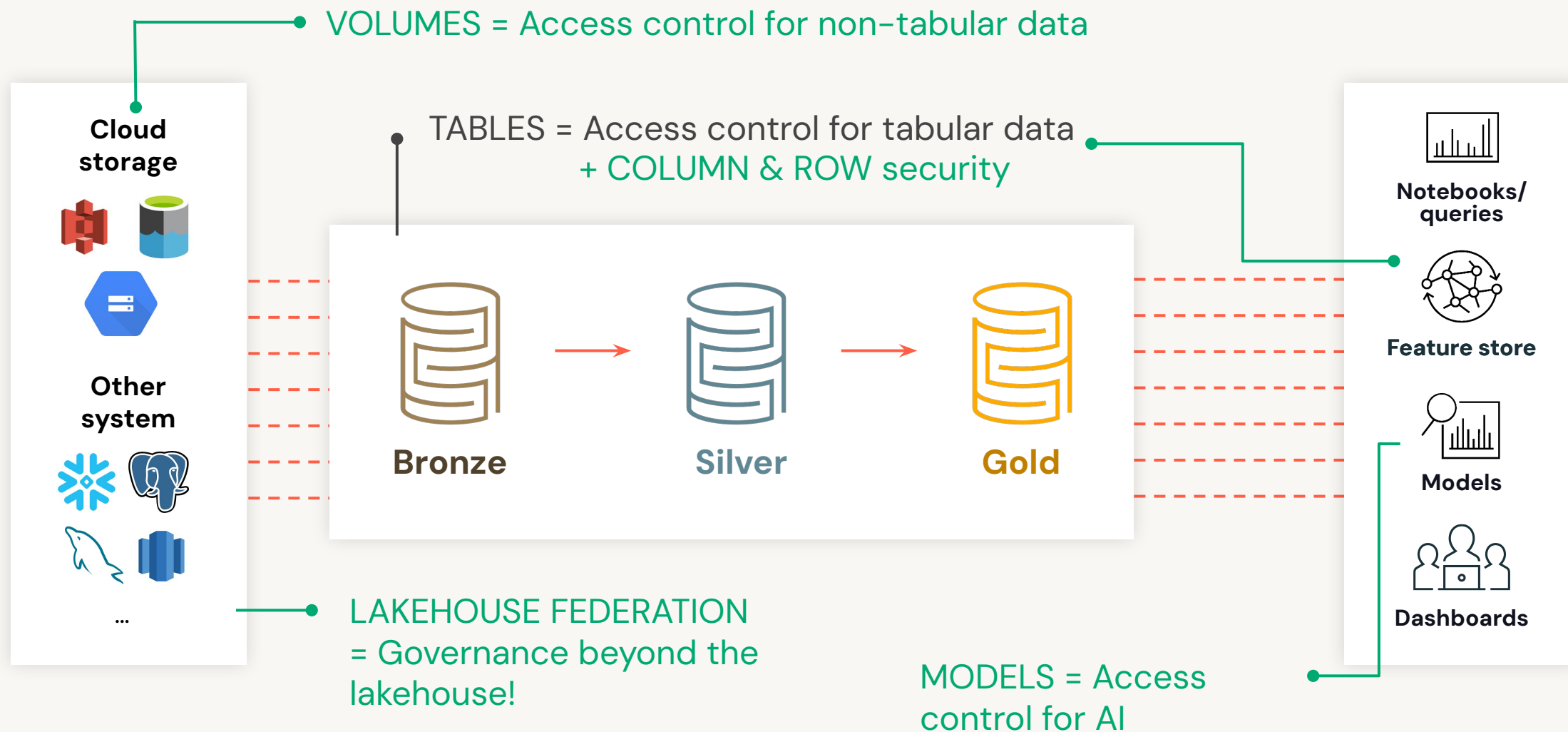
**Built-in data search and discovery**
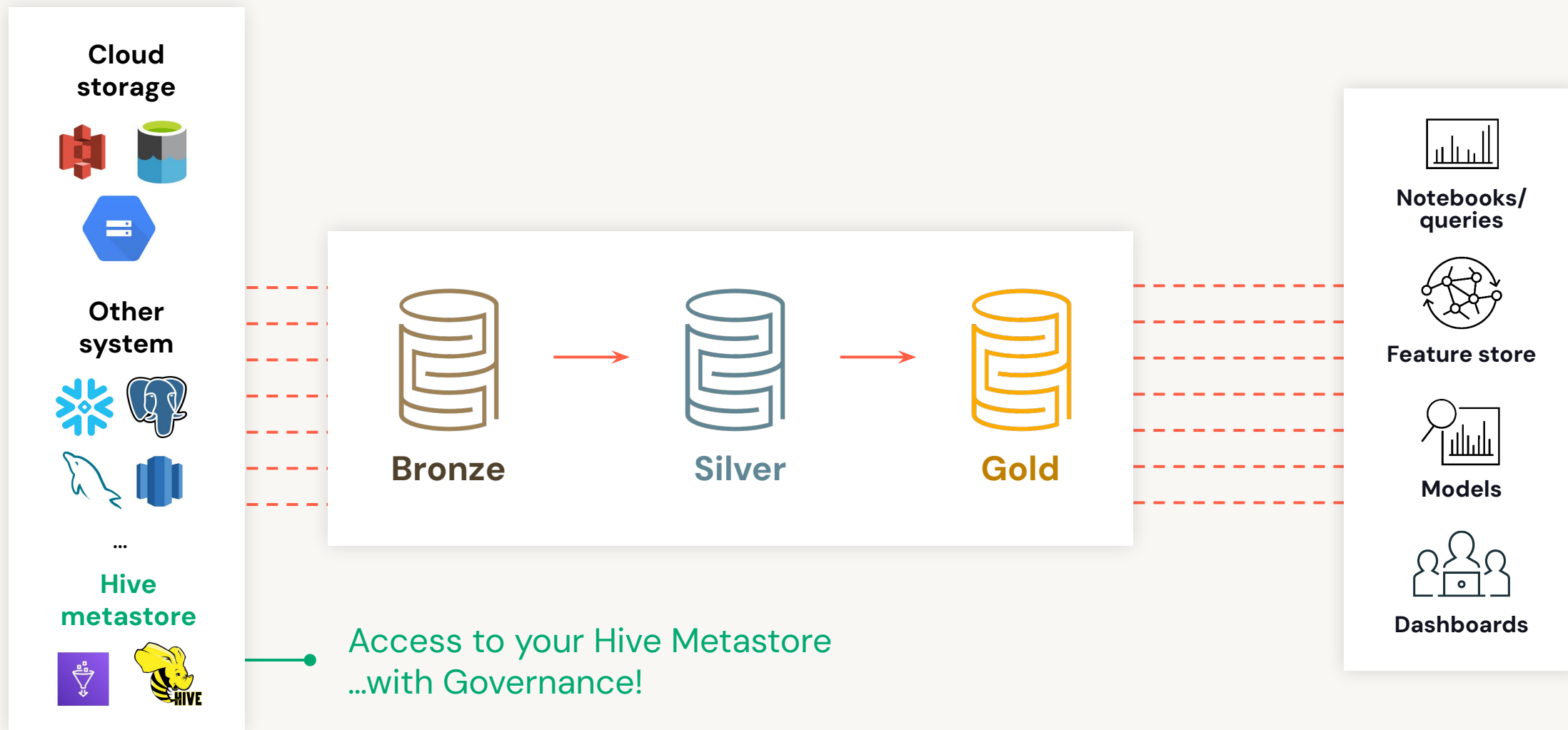
**Performance & scalability**

**Lineage & data quality**

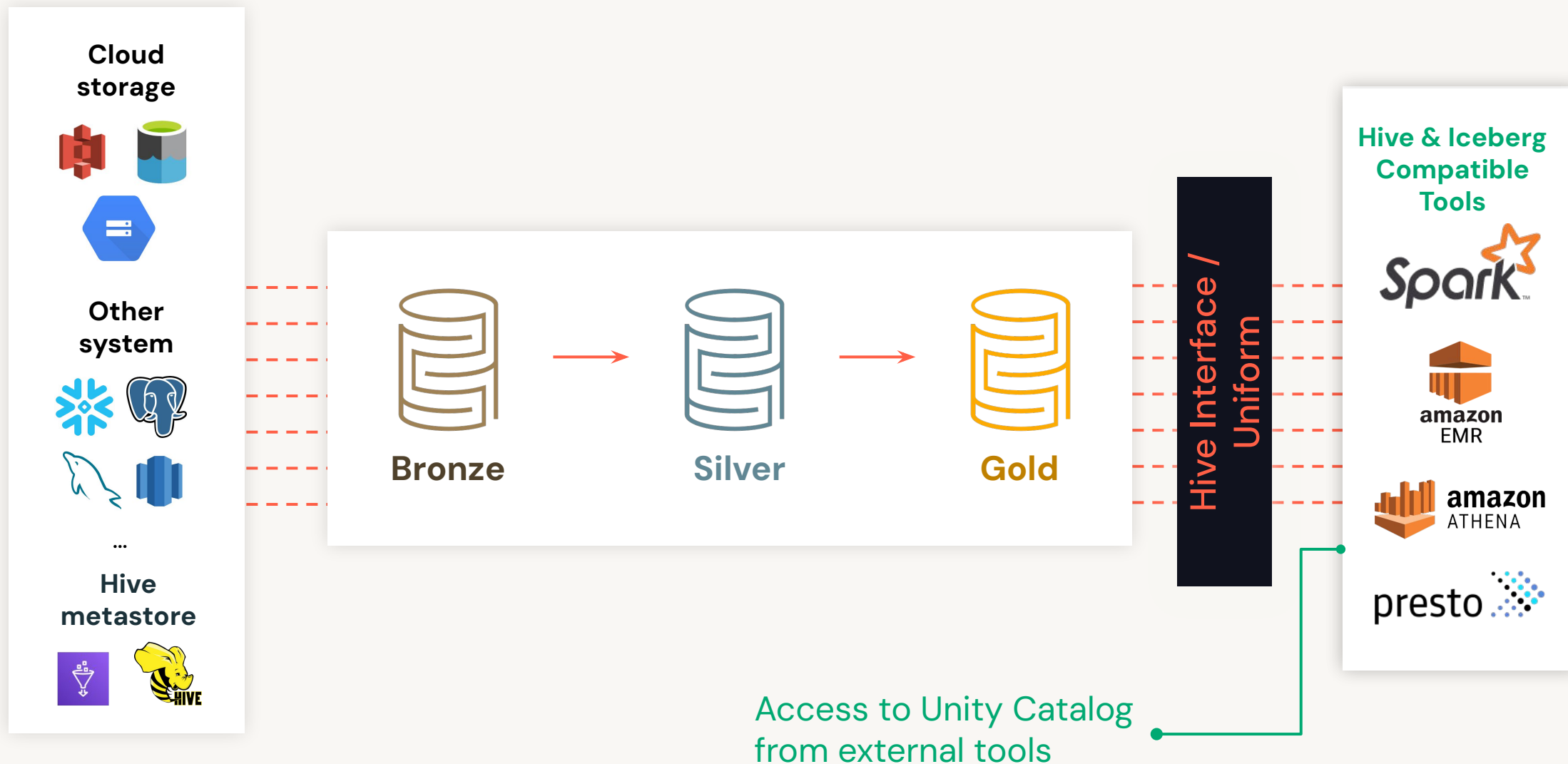# Access control for data + AI assets

VOLUMES = Access control for non-tabular data

Cloud storage

TABLES = Access control for tabular data
+ COLUMN & ROW security

Notebooks/ queries

Bronze → Silver → Gold

Feature store

Other system

...

LAKEHOUSE FEDERATION = Governance beyond the lakehouse!

Models

Dashboards

MODELS = Access control for AI

# Hive metastore federation `Coming Soon!`



**Cloud storage**

**Other system**

...

**Hive metastore**

**Bronze** → **Silver** → **Gold**

Access to your Hive Metastore
...with Governance!

**Notebooks/queries**

**Feature store**

**Models**

**Dashboards**

# Hive metastore interface & Uniform <span>Private Preview</span>

**Cloud storage**

**Other system**

...

**Hive metastore**

**Bronze** → **Silver** → **Gold**

Hive Interface / Uniform

**Hive & Iceberg Compatible Tools**

Spark

amazon EMR

amazon ATHENA

presto

Access to Unity Catalog from external tools

11

# Clusters



**Shared Cluster**

Support concurrent users/workloads with full isolation guarantees, fine grained security



**Single User Cluster**

Support advanced use cases such as GPUs, distributed ML, RDDs

# UC Shared clusters

## Develop and deploy your workloads on shared clusters!



**Interactive development**

**Scheduled workflows**

Shared Cluster

Support concurrent users/workloads with full isolation guarantees

**Recently launched (DBR 13+)**
- Python cluster libraries (13.1)
- Scalar Python/Pandas UDFs (13.2)
- Single node ML & MLflow
- Structured Streaming
- Network connectivity from Python

**Private previews:**
- Scala
- Init scripts / Jars
- MLR (light, single node ML)

# Fine-grained access control for single-user clusters

## Coming Soon!



**Query views and masked tables securely from single-user clusters!**

View and masked table access:

- Data filtered via secure, serverless filtering service,
- Filtered results are streamed back to the single user cluster

# Demo: Access Control for the Lakehouse

# How do we get there?

**Centralized access control for data & AI**

**Built-in data search and discovery**

**Performance & scalability**

**Lineage & data quality**

# Describe and tag data

# Next Step : LakehouseIQ

# Find and discover assets



Search for any lakehouse entity!

Data browser in the notebook!

**Frequent notebooks/queries/dashboards:** How do I get started?

**Frequent Joins:** How can I enrich?

**Frequent users:** Who are the experts?

# Data insights

# How do we get there?
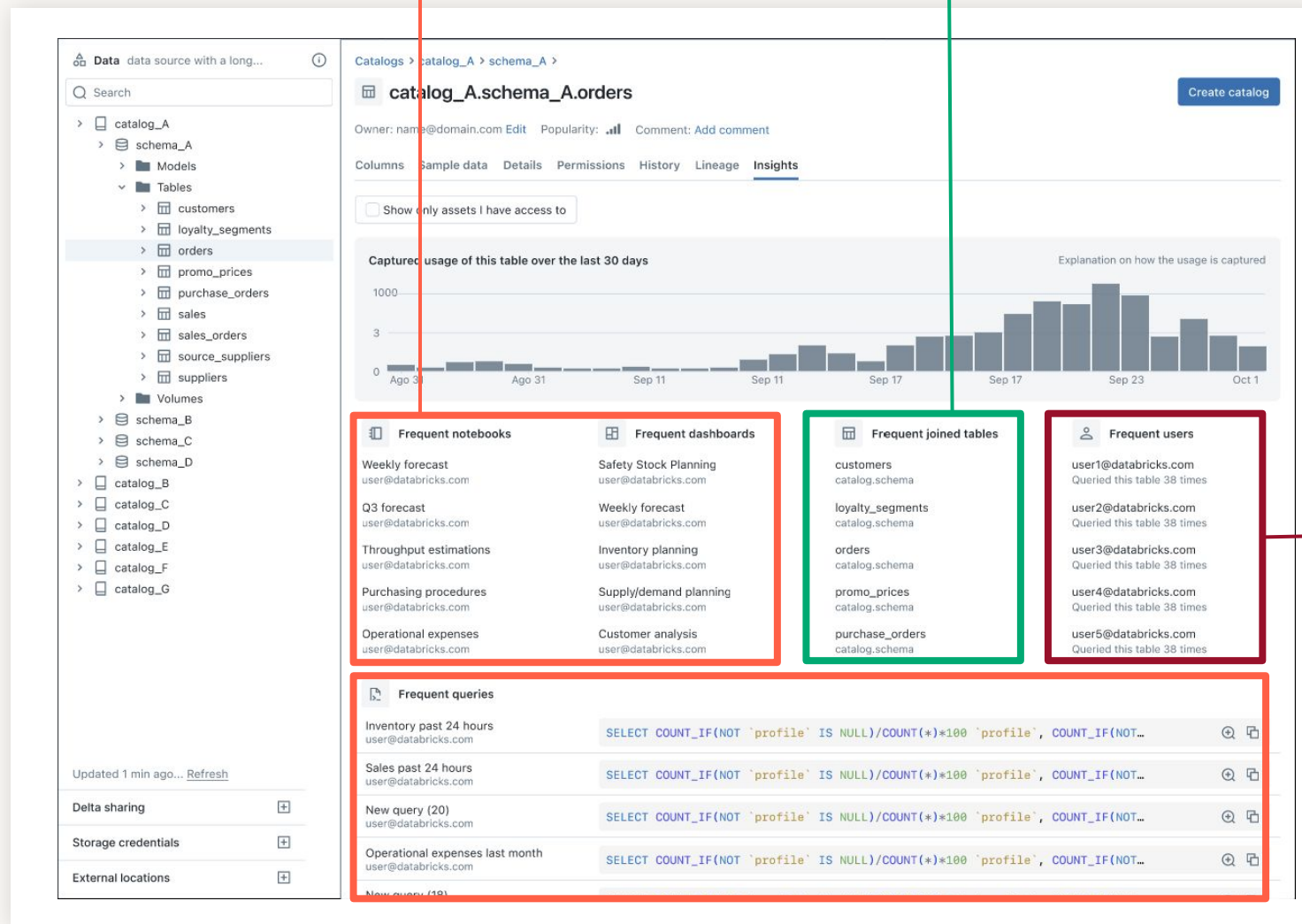
**Centralized access control for data & AI**

**Built-in data search and discovery**

**Performance & scalability**

**Lineage & data quality**

# Predictive optimization

**Challenge**

Tables can can be optimized for faster price-performance, but…

- Which tables?
- How often?
- What optimizations?

## ANALYZE TABLE

March 27, 2023

**Applies to:** ✅ Databricks SQL ✅ Databricks Runtime

The `ANALYZE TABLE` statement collects statistics about one specific ta schema, that are to be used by the query optimizer to find a better qu

## Syntax

```
ANALYZE TABLE table_name [ PARTITION clause ]
    COMPUTE STATISTICS [ NOSCAN | FOR COLUMNS col1 [, ...] | FOR

ANALYZE TABLES [ { FROM | IN } schema_name ] COMPUTE STATISTICS [
```

## Compa
## Delta La

April 18, 2023

See OPTIMIZE.

Delta Lake on Databri
coalesce small files in

## Syntax exa

You trigger compaction by running the `OPTIMIZE` command:

**SQL**   Python   Scala

```
OPTIMIZE delta.`/data/events`
```

# Predictive Optimization

**Solution**

- AI model determines which tables to optimize
  - Bin-packing
  - Liquid clustering
  - Vacuuming
  - Table statistics
  - Intelligent storage tiering
- Databricks automatically performs maintenance
  - Using hyper-optimized serverless compute

# Predictive optimization <span style="background-color:#6dc2a0">Public Preview</span>

**Results:**



Legend: ■ Query Performance  ■ Cost

Categories: No Optimization | Daily, Manual Optimization | Predictive Optimization

**Fewer knobs**

No need to:

- Determine which tables to optimize
- Run and manage the clusters

# How do we get there?

**Centralized access control for data & AI**

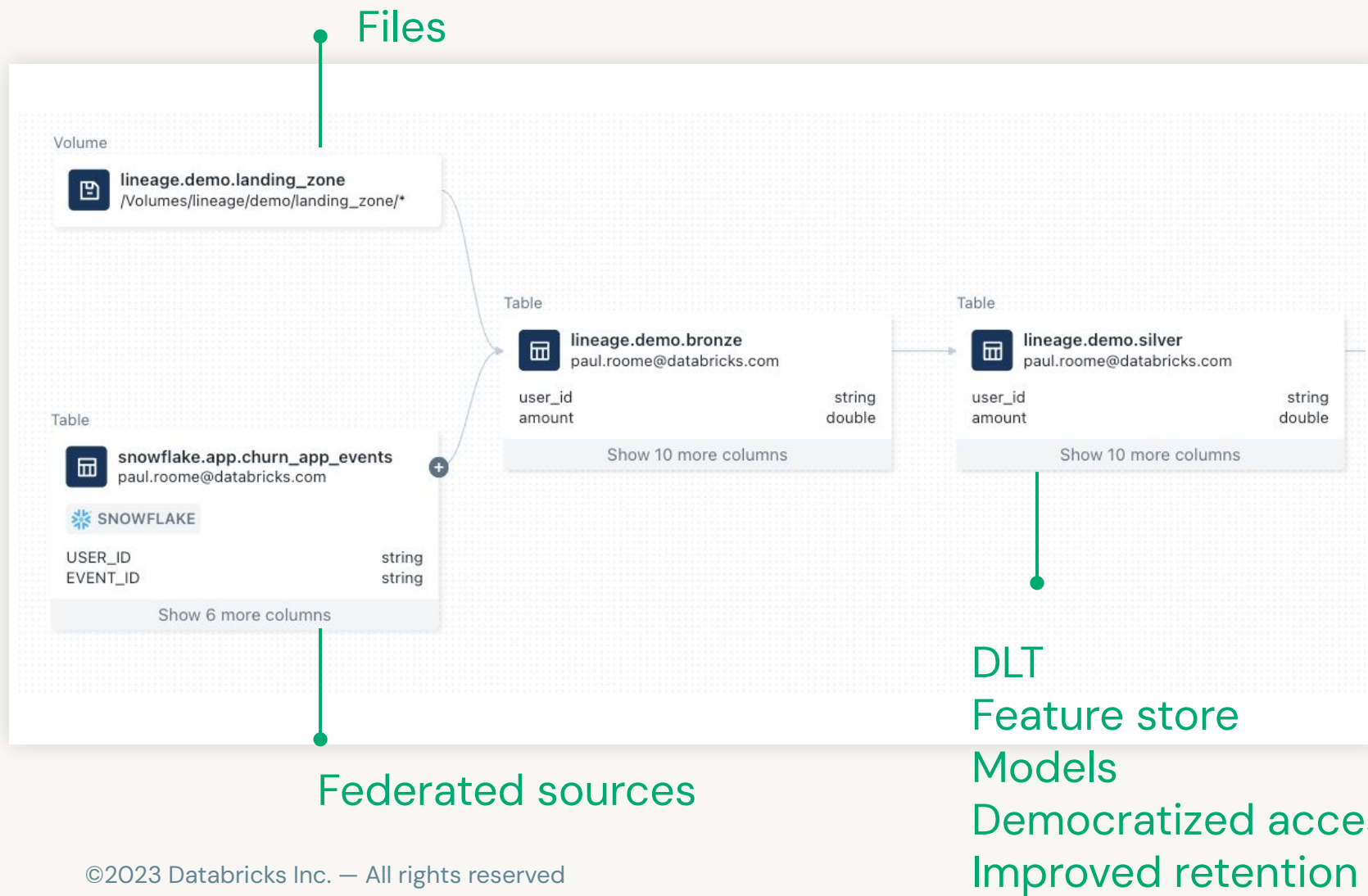**Built-in data search and discovery**

**Performance & scalability**

**Lineage & data quality**

# Lineage for the Lakehouse



**Files**

**Federated sources**

**DLT**
**Feature store**
**Models**
**Democratized access (BROWSE)**
**Improved retention logic**

**Automatic Data Documentation**

**Impact analysis**

**Root cause analysis**

26

# Lakehouse monitoring

**Data quality alerting**

**Auto classification of sensitive data**

**Model performance and drift monitoring**

**Lakehouse monitoring**

**The first AI-powered monitoring service for both Data and ML**

# Lineage + Lakehouse Monitoring = 💖

**Root cause analysis**
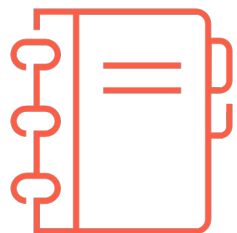


**Issue Detected by Lakehouse Monitoring**

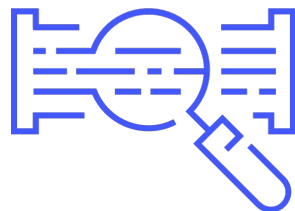# Lineage + Lakehouse Monitoring = 💖

**Proactive alerting**



Quality issue identified

Downstream impact

At Risk

# Demo: Lineage for the Lakehouse

# What we've seen...

**Secure and auditable**

**Discoverable**

**Usable and performant**

**Accurate and high quality**

# Even more on the horizon...



**Auto Tag + Mask**



**LakehouseIQ
Augmented Search +
Discovery**



**Governance Portal**