DATA+AI SUMMIT
BY databricks

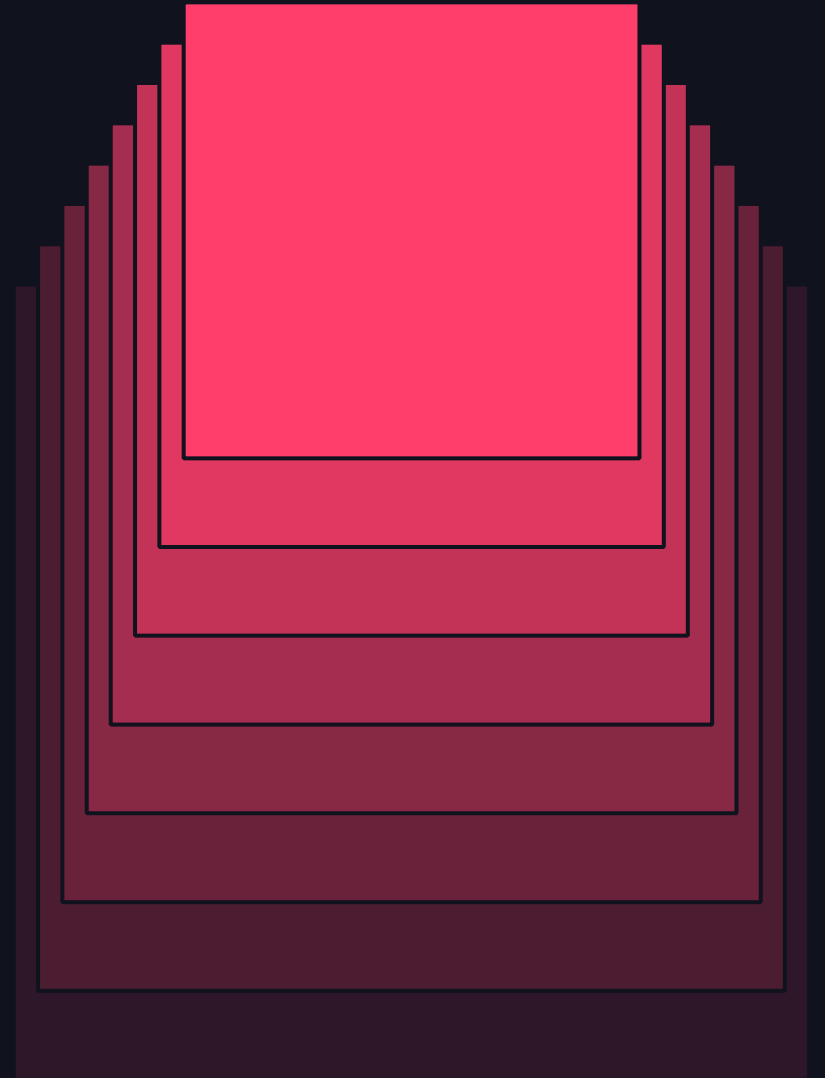# BUILDING HIGH-QUALITY AND TRUSTED DATA PRODUCTS

# WITH DATABRICKS

Karthik Subbarao | Pawarit Laosunthara
June 2024

# Karthikeya Sampa Subbarao

**Databricks EMEA**

❏ Specialist Solutions Architect @ **Databricks**

❏ Experienced Software Engineer & Architect

❏ 10+ years in the Tech industry

❏ Data Architecture, Security & Governance SME



CONNECT

DATA AI SUMMIT

# Pawarit Laosunthara

**Databricks AMER**



❏ Sr. Solutions Architect @ Databricks

❏ Working with Databricks' largest customers
in Logistics, Financial Services, Manufacturing

❏ Previous roles
    ❏ Tech Lead at Thoughtworks
    ❏ Data Scientist at Airbus

**CONNECT**

DATA+AI SUMMIT

# Agenda

- Data Products & Lifecycle

- Data Contracts & Governance

- Publishing & Discovery

- Demo

- Interoperability

- Takeaways

# 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
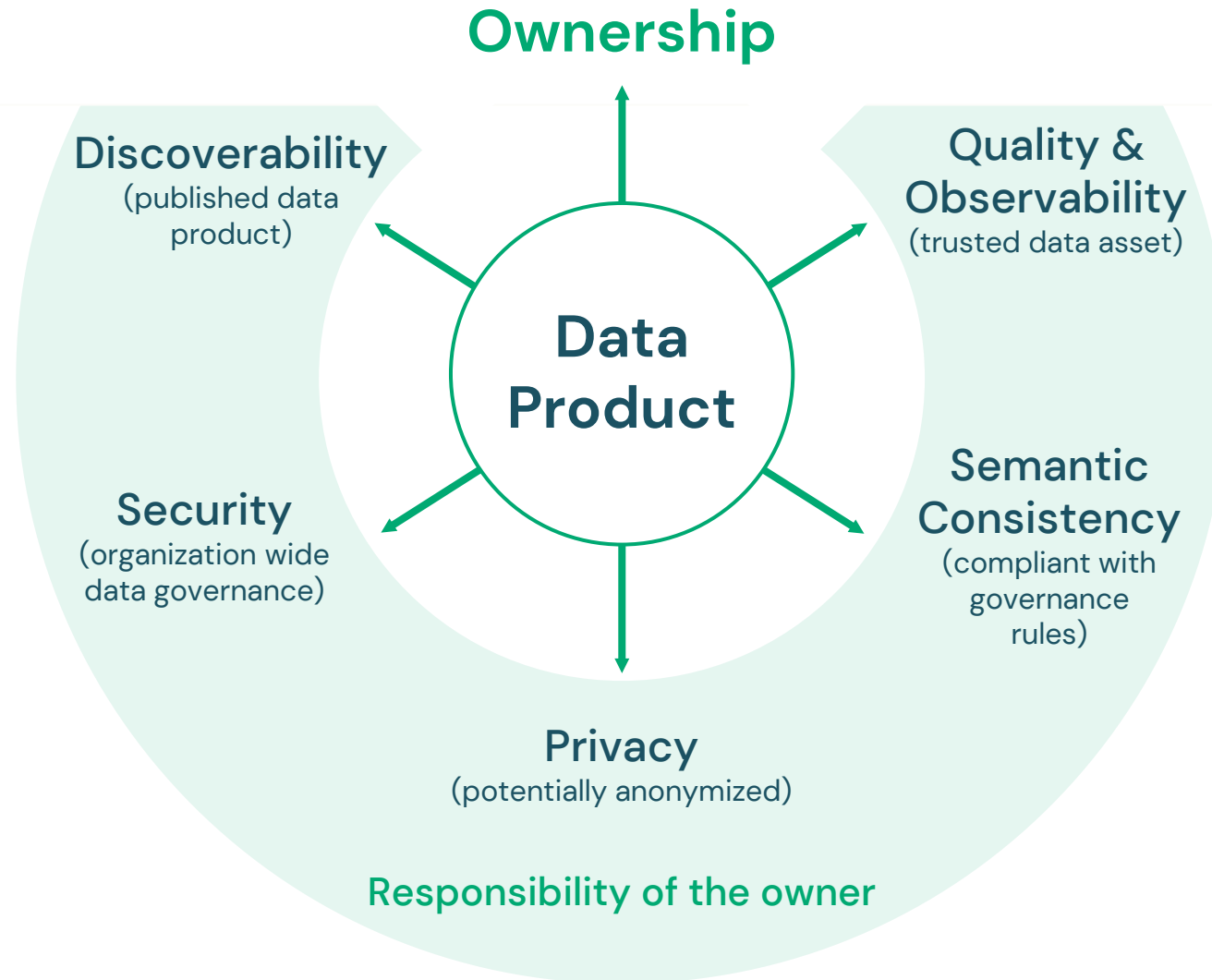
# Data Products

# Data and "product thinking"

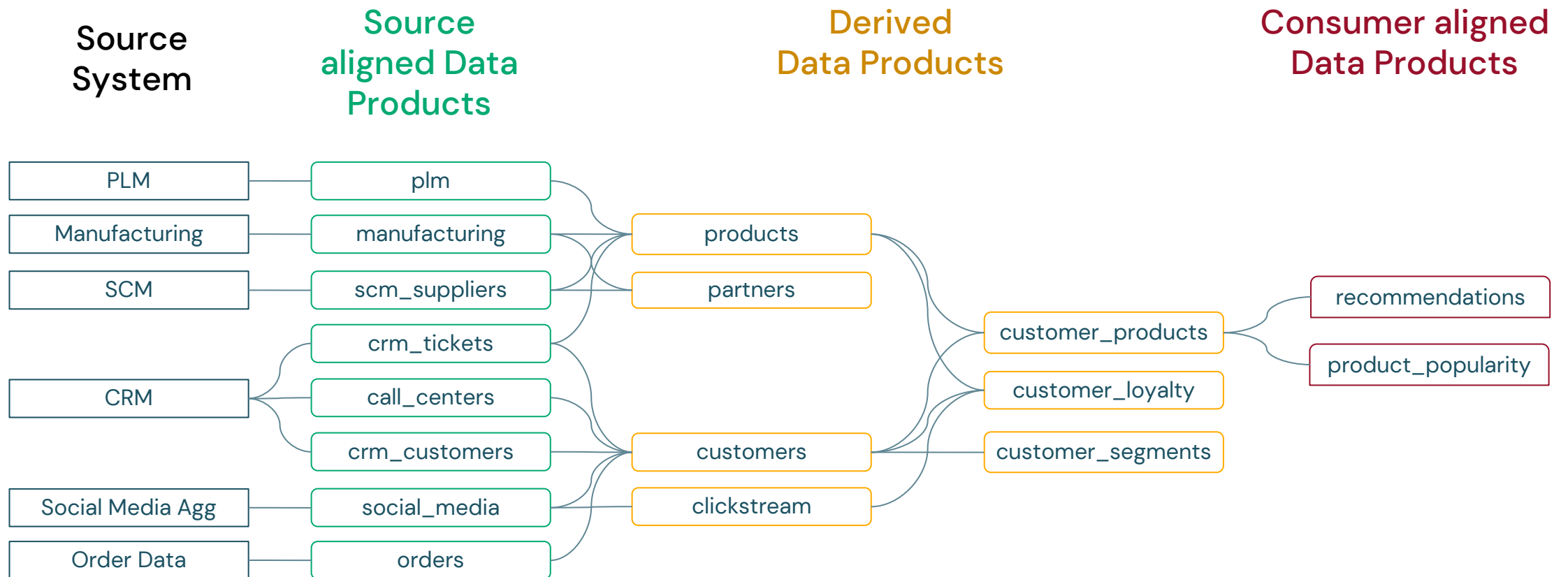To publish data as data products, "product thinking" needs to be applied

Data products should:

- have an owner  and is built for specific audiences

- follow a defined product life cycle

- be defined and described by data contracts

- be published following an agreed governance process

# Adding Data Product attributes to the concept



**Ownership**

**Discoverability**
(published data product)

**Quality & Observability**
(trusted data asset)

**Data Product**

**Security**
(organization wide data governance)

**Semantic Consistency**
(compliant with governance rules)

**Privacy**
(potentially anonymized)

**Responsibility of the owner**

# Data product hierarchy



| Source System | Source aligned Data Products | Derived Data Products | Consumer aligned Data Products |
|---|---|---|---|
| PLM | plm | products | recommendations |
| Manufacturing | manufacturing | partners | product_popularity |
| SCM | scm_suppliers | customer_products | |
| | crm_tickets | customer_loyalty | |
| CRM | call_centers | customers | |
| | crm_customers | customer_segments | |
| Social Media Agg | social_media | clickstream | |
| Order Data | orders | | |

9

# Data Product Lifecycle

# Data Products

## Typical Lifecycle

# Data Products

## Typical Challenges

How can I satisfy different consumption modes?

(tables, files, ML $models, streams$)

?

**Consumption + Value Creation**

feedback ⟶

Inception → **Design** → **Creation** → **Publishing** → **Operation + Governance** → Retirement

feedback

information

**Design**

?

What do I need to build?

**Creation**

?

How do we ensure fresh/reliable data?

How will different personas collaborate?
$Analysts, Engineers, Data\ Scientists$

**Publishing**

?

Where does the data came from?

How can I monitor usage/operations?

**Operation + Governance**

?

# Data Products

## Mapping the Databricks Lakehouse to the data product lifecycle

**Data and AI Governance**

📖 Unity Catalog

| Discovery | Lineage |
| Insights | Quality |

**Data Warehousing**
Dashboards

**Mosaic AI**
ml*flow*   Notebooks

**Consumption + Value Creation**

feedback ⟍

feedback   information

| Inception | → | Design | → | Creation | → | Publishing | → | Operation + Governance | → | Retirement |

Owner

Team

**Docs repo**

data contract   data product spec

**Orchestration**
Workflows   Delta Live Tables

**Data Warehousing**
Databricks SQL
Lakeview

**Mosaic AI**
ml*flow*   Notebooks
Feature Store

**Data and AI Governance**

📖 Unity Catalog

Access Control   Data Explorer   Auditing   System tables

Lineage   Lakehouse Monitoring

# Data Contract and Governance

# Data Contract

| Data description |
|---|
| name, owner, description, source systems, … |

| Data SLAs |
|---|
| last update, expiration dates, retention time, usage restrictions, code of conduct, re-sharing conditions, … |

| Data schema |
|---|
| tables, columns, anonymization and encryption info, … |

| Security |
|---|
| who is allowed to use the data product |

| Data quality |
|---|
| applied quality checks, quality metrics, … |

| Explanatory add-ons (optional) |
|---|
| notebook, dashboard, sample code, … |

# Bitol – Linux Foundation AI & Data sandbox project

Example: Open Data Contract Standard (ODCS)

# Data Product Certification

## Example process to achieve standards and consistency



**(2)** Assess Feedback Approve

Governance Team

**Certified data products** have a "stamp of approval" from the governance team (golden data products)

**(1)** Propose data contract

**(3)** Contract approval

**(6)** Understand usage (via data contract)

Domain 1

**(4)**

Publish

Catalog / Marketplace

Discover **(5)**

Domain 2

**(7)** Use data product

Cloud storage

# Independent and certified data products

## Balance between centralization and autonomy

high quality data products
Semantically consistent
easy to integrate with
other data products

**Certified Data Products**

Data Product

**Independent Data Products**

Data Product

good quality
…but no guarantee that they
can be combined easily

agree on rules and policies
for certified data products

**Governance Team**

approve and govern
data contracts

domains can publish data products
for use by adjacent teams
*as they see best*

**Autonomous Data Domain**

# Data Product
# Publishing and Discovery

# Publishing use cases for data products

Share them internally / externally and with or without restrictions

```
Data Products
```

**External Sharing**

Restricted

Unrestricted

- Databricks Marketplace (request for approval)
- Private Exchange
- Databricks Cleanrooms

Databricks Marketplace (instant access)

**Internal Sharing**

Restricted

Unrestricted

Same cloud region:
- Unity Catalog

Cross-cloud / cross-region:
- Private Exchange (request for approval)

Same cloud region:
- Unity Catalog

Cross-cloud / cross-region:
- Private Exchange (instant access)

# Data Contract with Unity Catalog metadata



Ownership

Tags derived from data contract, e.g. PII

Markdown allows hyperlinks to link back to external data contract*

Summaries and important details

Lineage

Code definition

Associated collateral

Stored externally

* consider using notebooks in the platform

# Data Contract with Databricks Marketplace

## Share summaries and important details and link full data contract
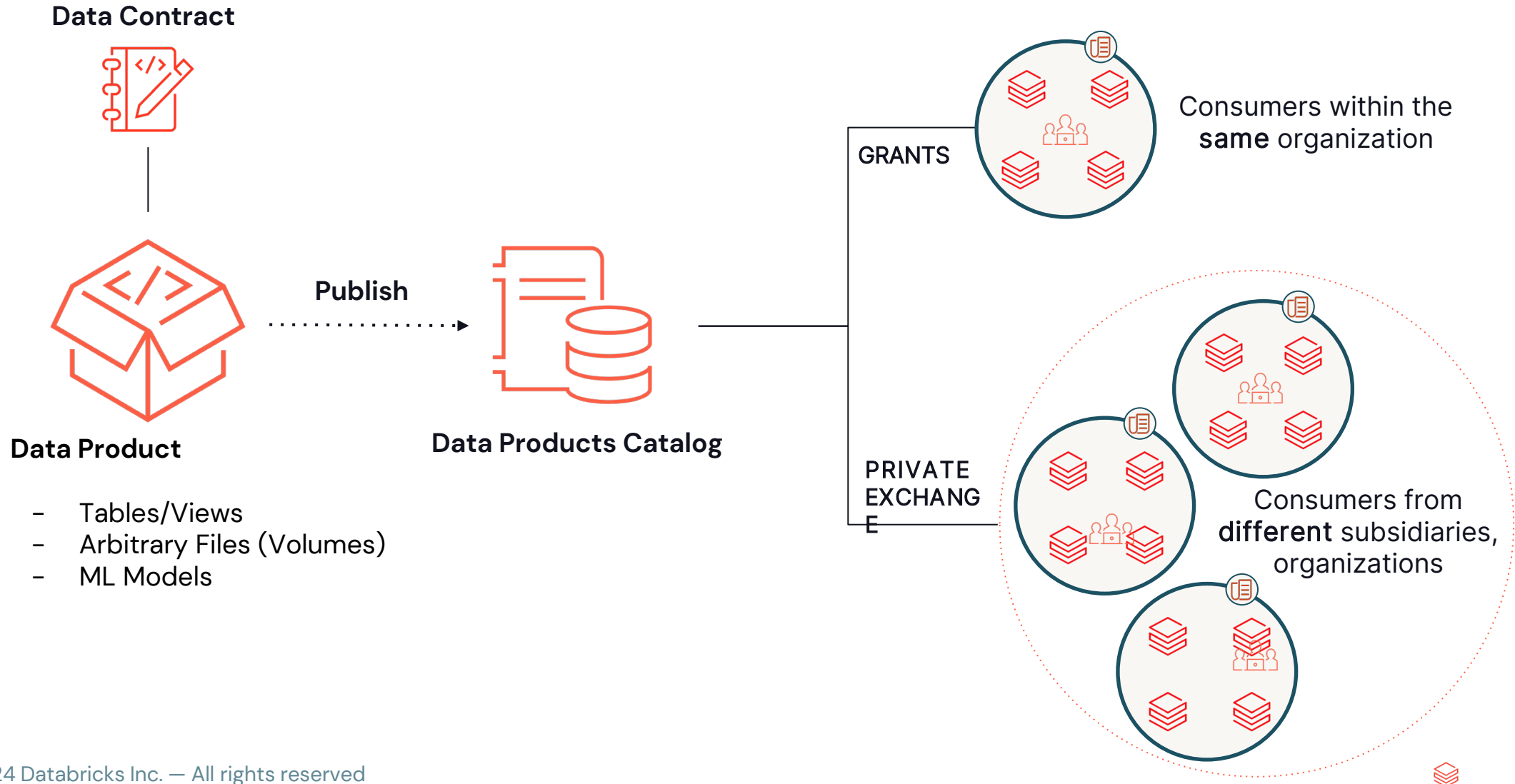
# Demo

# Demo

## What you'll see

➔ How consumers can discover data products

➔ How a data contract can look like in Databricks

➔ How to monitor the quality of the data products

# Demo

**Data Contract**

**Data Product**

- Tables/Views
- Arbitrary Files (Volumes)
- ML Models

**Publish**

**Data Products Catalog**

GRANTS

Consumers within the **same** organization

PRIVATE EXCHANGE

Consumers from **different** subsidiaries, organizations
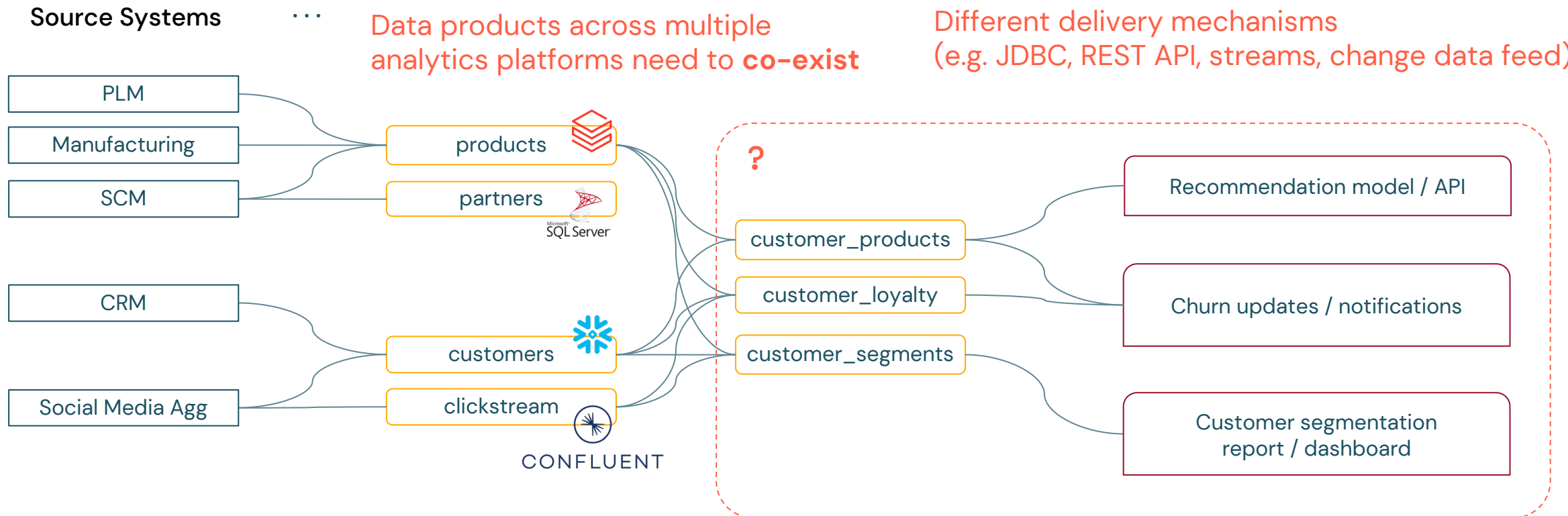
# Interoperability

# Reality of Enterprise Data Platforms

Data products come in many shapes and forms

Challenge #1: Ecosystem

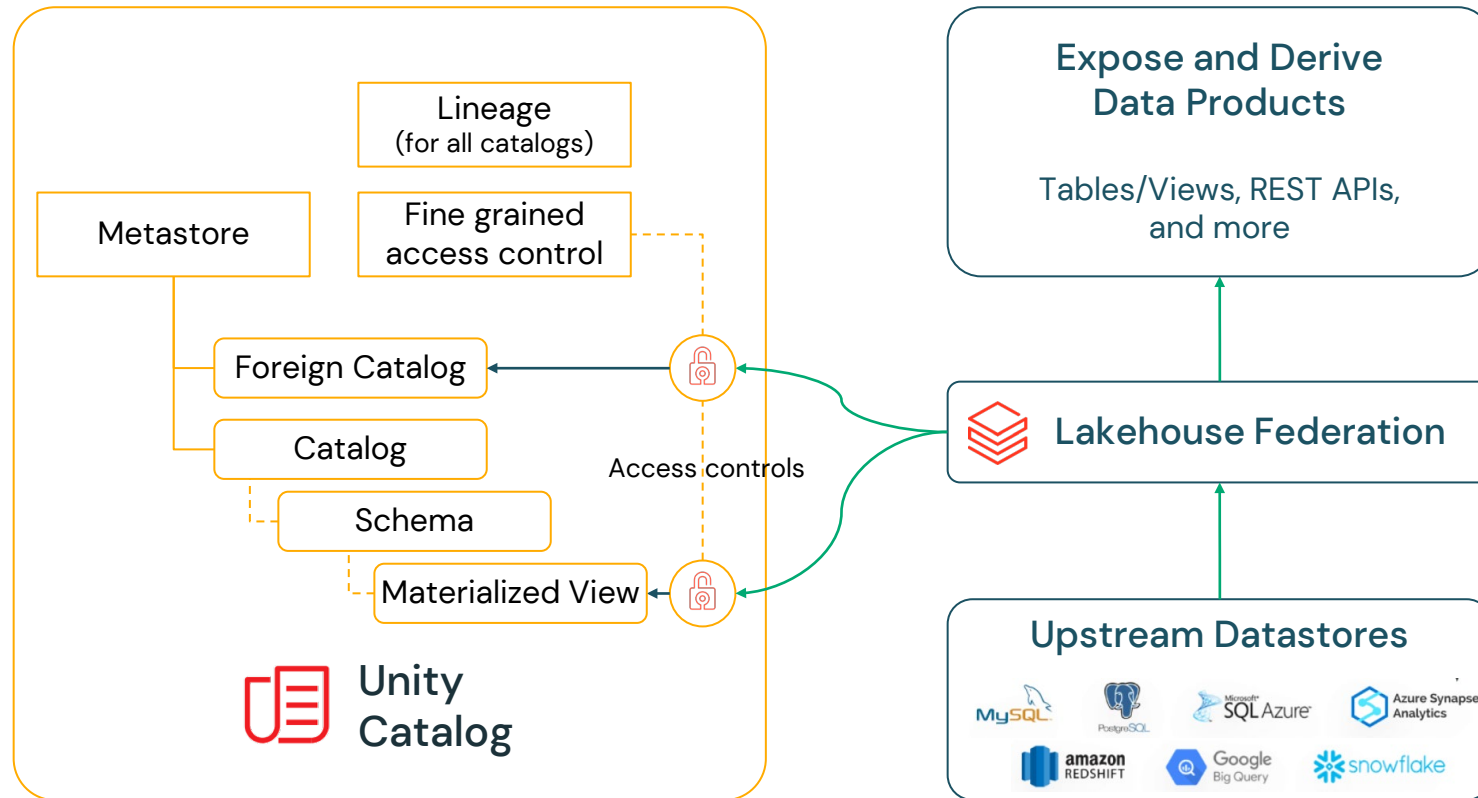Data products across multiple
analytics platforms need to **co-exist**

Challenge #2: Consumption Modes

Different delivery mechanisms
(e.g. JDBC, REST API, streams, change data feed)
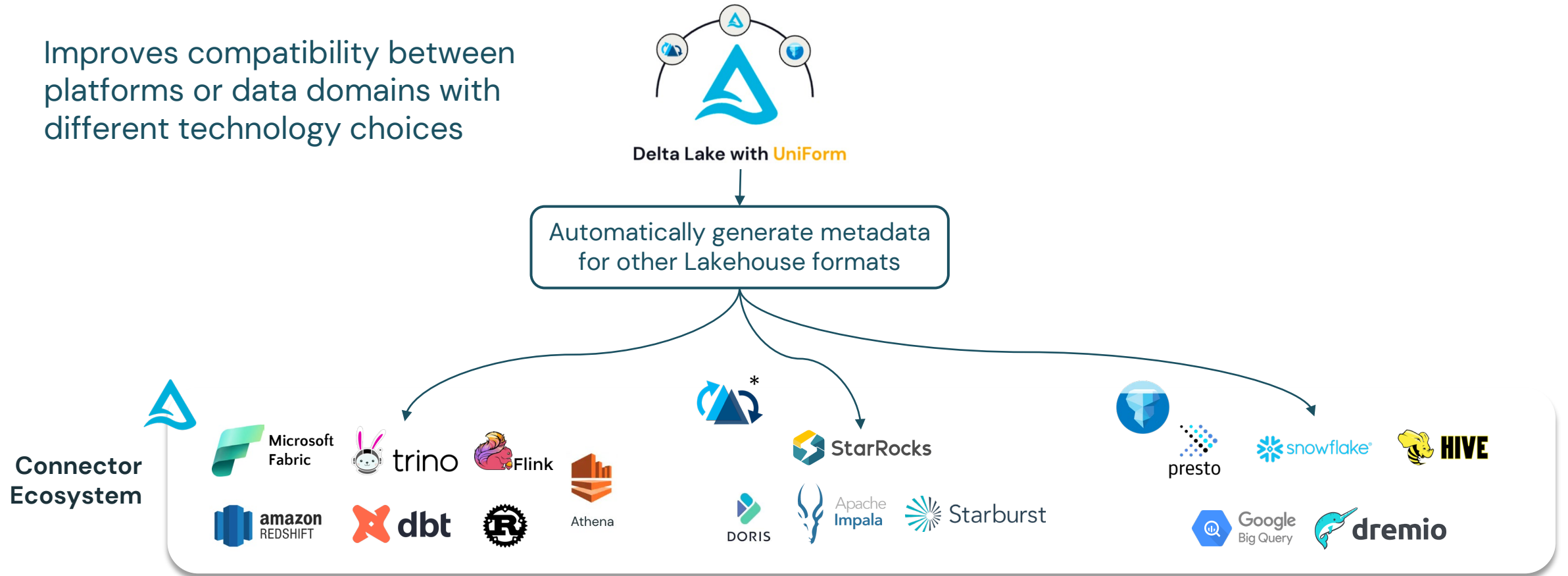
# Build on top of data products from other datastores

Bridge siloes with Lakehouse Federation

# Interoperability with other Lakehouse formats

## UniForm – Universal Format

Improves compatibility between platforms or data domains with different technology choices



**Delta Lake with UniForm**

Automatically generate metadata for other Lakehouse formats

**Connector Ecosystem**

Microsoft Fabric · trino · Flink · Athena · amazon REDSHIFT · dbt · R

StarRocks · DORIS · Apache Impala · Starburst

presto · snowflake · HIVE · Google Big Query · dremio

\* Hudi support in roadmap

# Takeaways

# Building high-quality data products with Databricks

## Data Products should be:

1. ### Discoverable
   - Take advantage of Unity Catalog's AI-assisted documentation, search

2. ### Reliable and Transparent
   - Prevent quality issues with Delta Live Tables
   - Monitor with Lakehouse Monitoring

3. ### Well governed (access controls, PII, auditable)
   - Leverage Unity Catalog's row/column-level security + tagging
   - Activate **system tables** for auditability, lineage

# Learn more at the summit!

**Databricks Events App**

## Tells us what you think

- We kindly request your valuable feedback on this session.

- Please take a moment to rate and share your thoughts about it.

- You can conveniently provide your feedback and rating through the **Mobile App**.

## What to do next?

- Discover more related sessions in the mobile app!

- Visit the Demo Booth: Experience innovation firsthand!

- More Activities: Engage and connect further at the Databricks Zone!

## Get trained and certified

- Visit the Learning Hub Experience at Moscone West, 2nd Floor!

- Take complimentary certification at the event; come by the Certified Lounge

- Visit our Databricks Learning website for more training, courses and workshops!
databricks.com/learn

# Thank You