

Integrating Unity Catalog with real-time, event-driven platforms like Kafka to enable seamless data governance and support real-time use cases

In this article, I aim to present a real-time data architecture challenge: how to effectively manage data governance in a modern environment. One possible solution we are exploring is leveraging Databricks Unity Catalog, as most of the current organizational data resides in the PySpark Data lake framework. The focus is on integrating Unity Catalog with real-time, event-driven platforms like Kafka to enable seamless data governance and support real-time use cases. This article explores potential use cases, challenges, and approaches to designing such a system."

Unity Catalog and Kafka serve fundamentally different purposes, and whether Unity Catalog can replace Kafka depends on the specific use case.

Here's a breakdown of their roles:

Unity Catalog

- **Primary Purpose:** Unity Catalog is a governance and metadata management solution for data lakes and warehouses. It provides centralized data discovery, access control, auditing, and lineage tracking for data stored across multiple cloud storage solutions like AWS S3, Azure Data Lake, or Google Cloud Storage.
- **Use Cases:**
 - Metadata management.
 - Data access governance.
 - Centralized auditing.
 - Data lineage tracking.

Kafka

- **Primary Purpose:** Kafka is a distributed messaging and event-streaming platform. It's used to build real-time data pipelines and applications that process and react to streams of data in real time.
- **Use Cases:**

- Real-time data streaming and event processing.
 - Decoupling services via publish/subscribe messaging.
 - Log aggregation.
 - Real-time analytics.
-

Can Unity Catalog Replace Kafka?

No, Unity Catalog cannot directly replace Kafka because:

1. Real-Time Messaging vs. Governance:

- Kafka is designed for real-time message processing and streaming, enabling data to flow between systems in near real-time.
- Unity Catalog focuses on governing and organizing static or semi-static datasets in data lakes or warehouses.

2. Event-Driven Architecture:

- Kafka supports event-driven architectures, where services communicate via events and can respond immediately to changes.
- Unity Catalog does not provide event-driven capabilities or real-time streaming functionality.

3. Decoupling Microservices:

- Kafka acts as a messaging layer to decouple microservices and ensure scalability and fault tolerance.
- Unity Catalog is designed for metadata and governance; it doesn't facilitate inter-service communication.

4. Data Velocity:

- Kafka is optimized for high-throughput, low-latency streaming.
 - Unity Catalog is optimized for managing and governing datasets.
-

When to Use Both Together

While Unity Catalog cannot replace Kafka, they can complement each other in a data architecture:

- Kafka for real-time data streaming and event-driven processing.
- Unity Catalog for governing and organizing the data generated or consumed by Kafka.

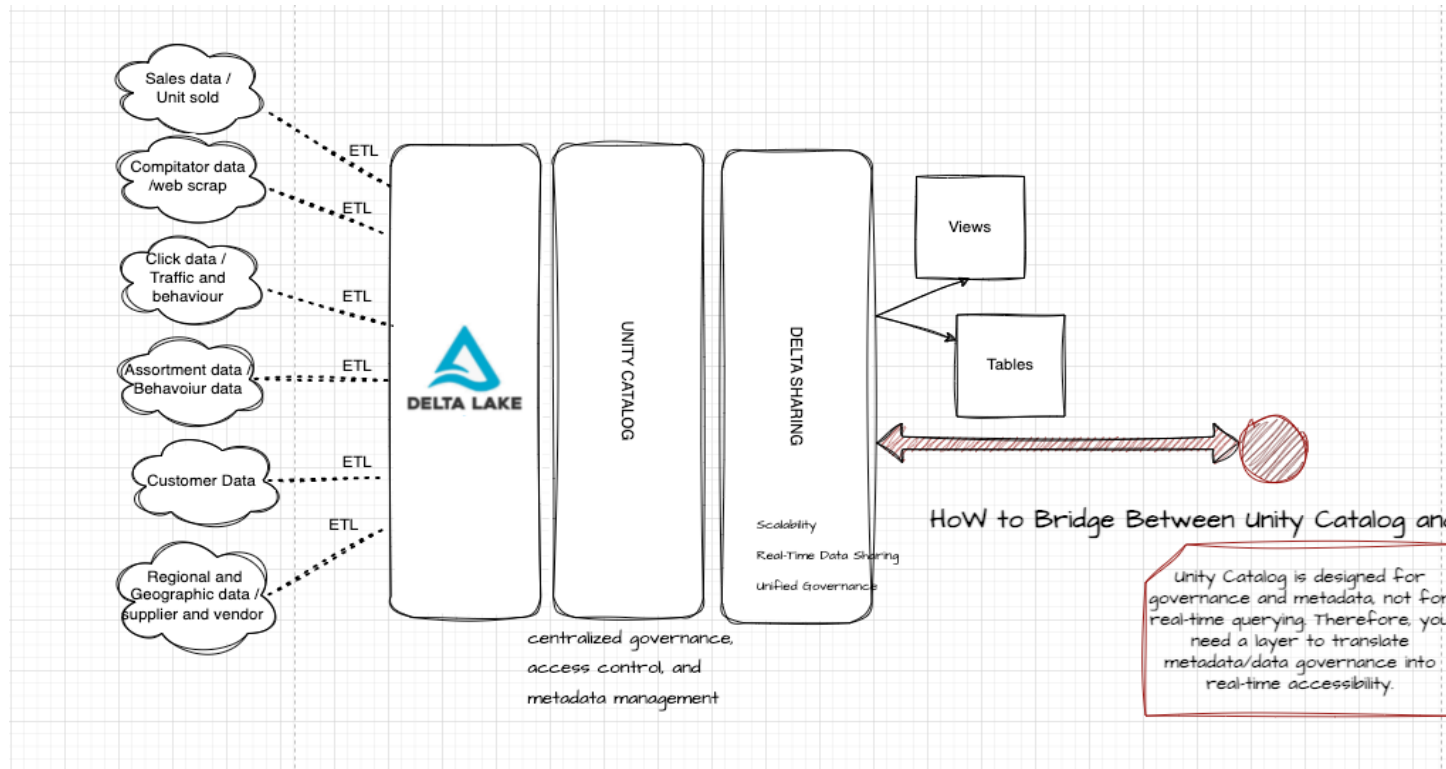
If you're looking to replace Kafka, you'd need another streaming or messaging system, such as:

- AWS Kinesis (streaming)
- Google Pub/Sub (messaging)
- Azure Event Hubs (streaming)

If your use case involves real-time messaging, Unity Catalog is not a suitable replacement for Kafka.

I have a data storage, basically data lake environment, where different kind of data is present, like Sales data /

Unit sold, Click data / Traffic and behavior, Assortment data / Behavior data., Regional and Geographic data / supplier and vendor, So my problem these data, I need a unified source of all these data, which can be fed to, which can be managed in Unity catalog, and then further it can be published to Kafka for making it more reliable for event streaming. What I would say that definitely there is a use case, which exists for Unity catalog inside our current organization, because the data management becomes more easier, data handling becomes more easier, data governance, centralized auditing, data lineage tracking, all these becomes easier for us.



But the problem is that how our real-time system will interact with Unity catalog.

This seemed to be a modern data architecture challenge—how to integrate data governance and management (via Unity Catalog) with real-time event-driven systems (via Kafka). Let's see break down this down into a solution if this makes some sense:

Proposed Architecture for Your Use Case

1. Data Ingestion into a Data Lake:
 - All your data sources (e.g., like **Sales data / Unit sold, Click data / Traffic and behavior, Assortment data / Behavior data., Regional and Geographic data / supplier and vendor**) are ingested into the data lake using ETL/ELT pipelines.
 - Tools like Databricks, Apache NiFi, or Azure Data Factory can streamline this process by ingesting and transforming data into the desired format.
2. Centralized Governance via Unity Catalog:

- Unity Catalog is used to manage, govern, and organize all datasets (structured and unstructured) within your data lake.
- Features you leverage:
 - Centralized data discovery for your analysts and teams.
 - Access control to ensure only authorized systems/services can access data.
 - Data lineage tracking to provide insights into the origin and transformations of data.
 - Auditing for compliance purposes.

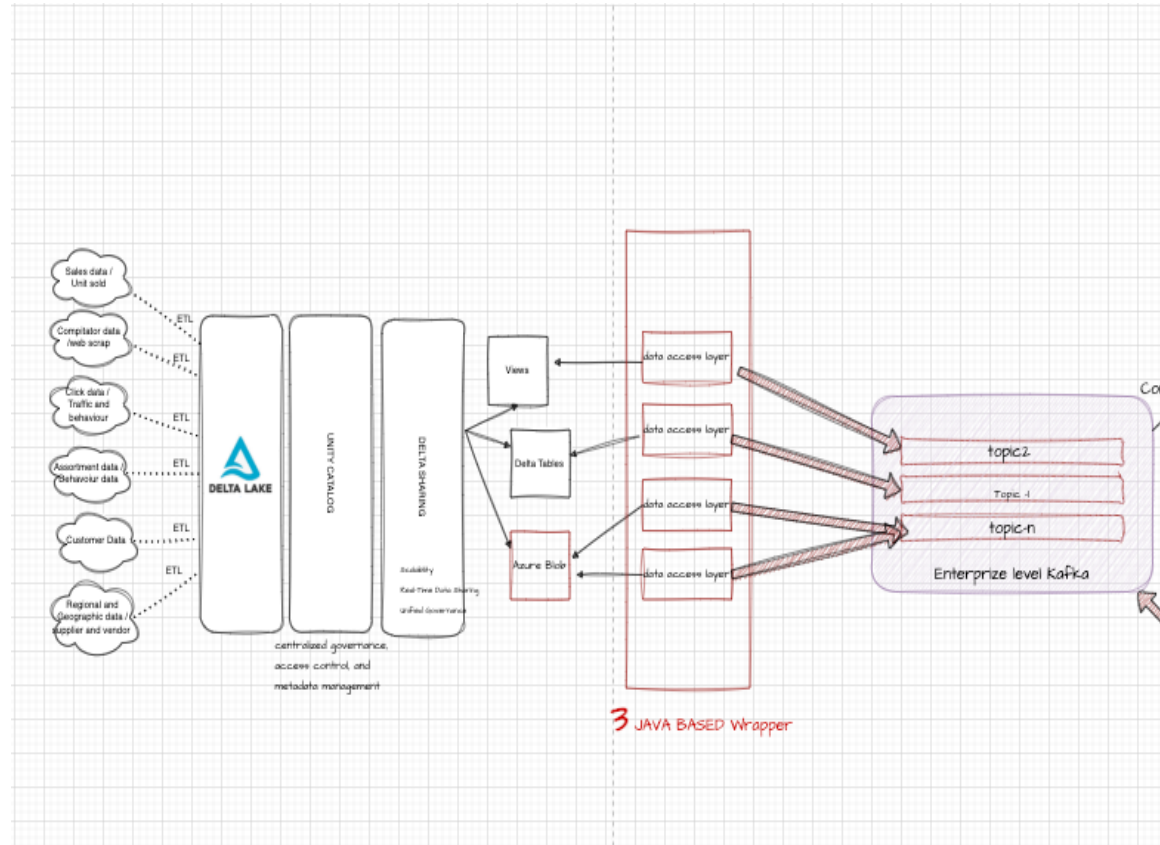
3. Real-Time Data System Integration with Kafka:

- Publish curated datasets to Kafka topics for real-time event-driven use cases:
 - When new or updated data becomes available in the lake (managed by Unity Catalog), publish an event to Kafka.
 - Downstream consumers can subscribe to these Kafka topics for real-time data processing or application updates.
- Kafka ensures:
 - Scalability for real-time processing.
 - Reliability and fault tolerance for event-driven systems.

4. Bridge Between Unity Catalog and Real-Time Systems:

- How Real-Time Systems Interact with Unity Catalog:
 - Unity Catalog is designed for governance and metadata, not for real-time querying. Therefore, you need a layer to translate metadata/data governance into real-time accessibility.
 - Implement a data access layer or metadata API:
 - Expose APIs or services that query Unity Catalog for metadata about the datasets (e.g., schema, access permissions).

- Combine metadata with direct data access paths (e.g., S3, Azure Blob, Delta Tables) to retrieve or stream relevant data in real-time.
- This access layer can feed data to Kafka for event-driven applications or directly to consumers if needed.



5. Integration Points Between Unity Catalog and Kafka:

- **Use change data capture (CDC) pipelines to detect updates in data managed by Unity Catalog and push events to Kafka. Tools like Debezium or Delta Live Tables can assist in this.**
- Data transformations or aggregations managed in Unity Catalog can be processed in batch mode or micro-batches (e.g., using Apache Spark) and then published to Kafka for real-time consumption.