

Mid-Term Assignment

DADS 6005 : Data Streaming and Real-Time Analytics

Chayaphon Sornthananon, 6610422007

Semester 1/2024

Lecturer : Asst.Prof Ekarat Rattagan



Table of Content

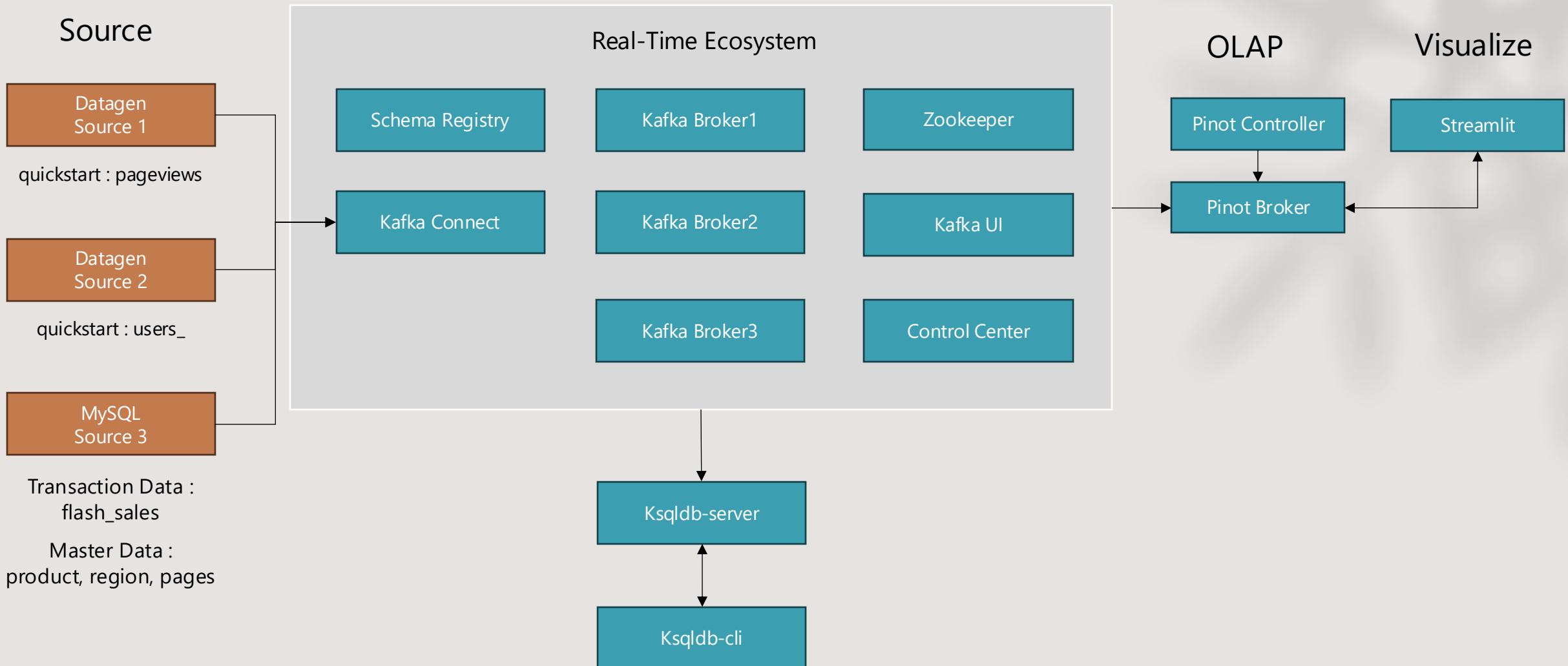
<u>Requirements</u>	P.3
<u>System Diagram</u>	P.4
<u>Code Reference</u>	P.5
<u>Docker</u>	P.6
<u>Role and Responsibility for each services</u>	P.7
<u>a.Data Source</u>	P.8-9
<u>b.Kafka System : Kafka Broker</u>	P.10
<u>b.Kafka System : Kafka Topics and Partition</u>	P.11-12
<u>b.Kafka System : Additional Explanation</u>	P.13
<u>b.Kafka System : Schema-Registry</u>	P.14
<u>b.Kafka System : Kafka Connect</u>	P.15
<u>b.Kafka System : Connectors for data source</u>	P.16
<u>c. KsqlDB</u>	P.17-18
<u>c. KsqlDB : Clean or transform data (topic4)</u>	P.19
<u>c. KsqlDB : Aggregation (join + group by) (topic5)</u>	P.20
<u>c. KsqlDB : Tumbling Window (Topic 6s)</u>	P.21
<u>c. KsqlDB : Hopping Window (Topic 7)</u>	P.22
<u>c. KsqlDB : Session Window (Topic 8)</u>	P.23
<u>c. KsqlDB : Testing the correctness</u>	P.24
<u>d. Apache Pinot</u>	P.25-26
<u>e. Streamlit Dashboard</u>	P.27-29
<u>Appendix</u>	P.30-33

Requirements

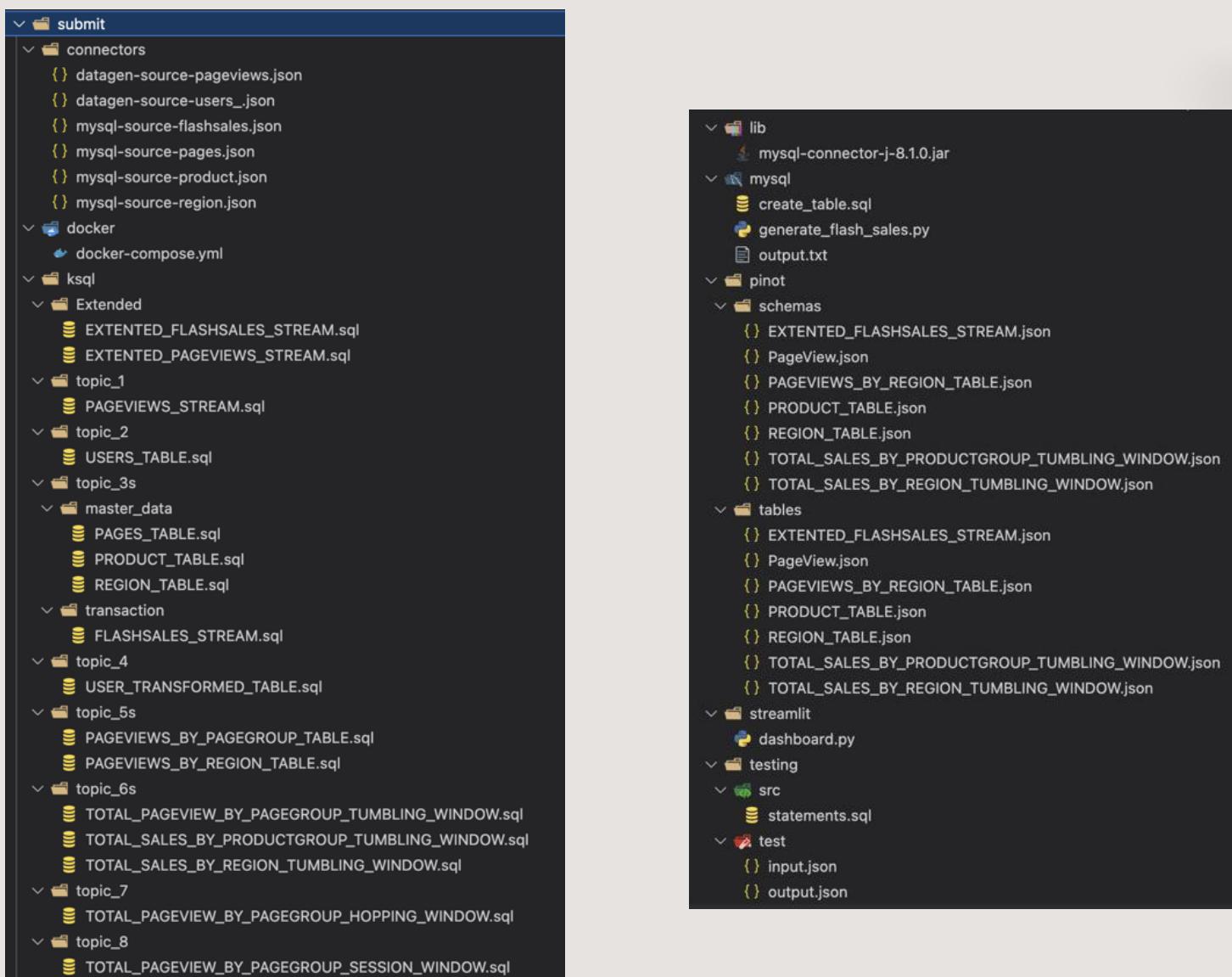
Data streaming and Realtime analytics System Requirement

- a. Data source 3 sources [1]
 - i. Source 1 : PageView (stream datagen) (topic1)
 - ii. Source 2 : Users_ (stream datagen) (topic2)
 - iii. Source 3 : Your design (relational database) (topic3)
- b. Kafka system [2]
 - i. 5 partitions
 - ii. 3 brokers
 - iii. 8 topics
 - iv. Schema Register
 - v. Kafka connect
- c. ksqlDB operation [3]
 - i. Clean or transform data (topic4)
 - ii. Aggregation (join + group by) (topic5)
 - iii. Windows
 - 1. Tumbling (topic6)
 - 2. Hopping (topic7)
 - 3. Session (topic8)
 - iv. Testing the correctness [4]
- d. Apache Pinot [5,6]
 - i. At least 3 queries
- e. Dashboard [7]
 - i. At least 4 panels
 - ii. Bonus: Interactive dashboard

System Diagram



Code Reference



Docker

In this project we use docker technology to help us running each services.

What is Docker

Docker is an **open-source containerization platform** that allows developers to package applications into **containers**.

Containers are lightweight, portable, and isolated environments that include everything needed to run an application, such as code, libraries, and dependencies.

Docker ensures that applications behave consistently across development, testing, and production environments.

Why Use Docker

Docker provides a solution to many challenges in modern application development and deployment:

1.Consistency: Docker ensures that the application runs the same in all environments, eliminating wording "it works on my machine" issues.

2.Portability: Docker containers are platform-independent and can run on any system with Docker installed, including local machines, servers, and cloud platforms.

3.Efficiency: Containers share the host operating system, making them faster and more lightweight compared to traditional virtual machines.

4.Scalability: Applications can scale easily by spinning up or down multiple containers as needed.

5.Ease of Deployment: Docker simplifies the deployment process by packaging applications and dependencies into one container.

This screenshot show all containers run for this project.

Container list										
Containers		Name	State	Quick Actions	Stack	Image	Created	IP Address	Published Ports	Ownership
<input type="checkbox"/>	zookeeper	running	Start Stop Kill Restart Pause Resume Remove	midterm	confluentinc/cp-zookeeper:7.7.1	2024-11-15 20:51:00	172.90.1.100	2181:2181	administrators	
<input type="checkbox"/>	kafka1	running	Start Stop Kill Restart Pause Resume Remove	midterm	confluentinc/cp-kafka:7.7.1	2024-11-15 20:51:00	172.90.1.101	9010:9010 9191:9191	administrators	
<input type="checkbox"/>	kafka2	running	Start Stop Kill Restart Pause Resume Remove	midterm	confluentinc/cp-kafka:7.7.1	2024-11-15 20:51:00	172.90.1.102	9020:9020 9192:9192	administrators	
<input type="checkbox"/>	kafka3	running	Start Stop Kill Restart Pause Resume Remove	midterm	confluentinc/cp-kafka:7.7.1	2024-11-15 20:51:00	172.90.1.103	9030:9030 9193:9193	administrators	
<input type="checkbox"/>	kafka-ul	running	Start Stop Kill Restart Pause Resume Remove	midterm	provectuslabs/kafka-ulatest	2024-11-15 20:51:00	172.90.1.104	8080:8080	administrators	
<input type="checkbox"/>	schema-registry	running	Start Stop Kill Restart Pause Resume Remove	midterm	confluentinc/cp-schema-registry:7.7.1	2024-11-15 20:51:00	172.90.1.105	8081:8081	administrators	
<input type="checkbox"/>	kafka-connect	healthy	Start Stop Kill Restart Pause Resume Remove	midterm	confluentinc/cp-kafka-connect:7.7.1	2024-11-15 20:51:00	172.90.1.106	8083:8083	administrators	
<input type="checkbox"/>	control-center	running	Start Stop Kill Restart Pause Resume Remove	midterm	confluentinc/cp-enterprise-control-center:7.7.1	2024-11-15 20:51:00	172.90.1.107	9021:9021	administrators	
<input type="checkbox"/>	ksqldb-server	running	Start Stop Kill Restart Pause Resume Remove	midterm	confluentinc/cp-ksqldb-server:7.7.1	2024-11-15 20:51:00	172.90.1.108	8088:8088	administrators	
<input type="checkbox"/>	ksqldb-cli	running	Start Stop Kill Restart Pause Resume Remove	midterm	confluentinc/ksqldb-cl:0.29.0	2024-11-15 20:51:00	172.90.1.109	-	administrators	
<input type="checkbox"/>	pinot-controller	running	Start Stop Kill Restart Pause Resume Remove	midterm	apachepinot/pinot:0.12.0	2024-11-15 20:51:00	172.90.1.111	9000:9000	administrators	
<input type="checkbox"/>	pinot-broker	running	Start Stop Kill Restart Pause Resume Remove	midterm	apachepinot/pinot:0.12.0	2024-11-15 20:51:00	172.90.1.112	8099:8099	administrators	
<input type="checkbox"/>	pinot-server	running	Start Stop Kill Restart Pause Resume Remove	midterm	apachepinot/pinot:0.12.0	2024-11-15 20:51:00	172.90.1.113	-	administrators	
<input type="checkbox"/>	midterm-streamlit-1	running	Start Stop Kill Restart Pause Resume Remove	midterm	python:3.9-slim	2024-11-15 20:51:00	172.90.1.115	8501:8501	administrators	
<input type="checkbox"/>	mysql	running	Start Stop Kill Restart Pause Resume Remove	midterm	mysql:8	2024-11-15 20:51:00	172.90.1.116	3306:3306	administrators	

How Docker Works

1.Docker file (Input):

- Developers write a dockerfile or docker-compose, which includes instructions to build a Docker image. The image contains the application code, dependencies, and runtime.
- Docker-compose help to pack all containers within one stack help for dependencies check on running and easy to maintain.

2.Docker Engine (Process):

- The Docker Engine reads the dockerfile to create an **image**, which serves as a blueprint for creating containers.
- When a container is launched from the image, it runs as an isolated, standalone instance.

3.Running Containers (Output):

- The containerized application is executed in an environment isolated from the host system, ensuring stability and consistency.
- Docker also enables networking, data sharing, and scaling through its runtime engine.

Role and Responsibility for each services

1. Zookeeper

- **Purpose:** Manages metadata and coordination for Kafka brokers.
- **Responsibilities:**
 - Handles configuration management and leader election for Kafka.
 - Keeps track of Kafka broker states.

• **Interactions:** Essential for Kafka brokers (kafka1, kafka2, kafka3) to function, as it maintains metadata about partitions and replicas.

2. Kafka Brokers (kafka1, kafka2, kafka3)

- **Purpose:** These brokers form the Kafka cluster.

- **Responsibilities:**
 - Store and serve data streams.
 - Distribute partitioned data across brokers for scalability.
 - Ensure fault tolerance with data replication.

- **Interactions:**
 - Communicate with Zookeeper to get metadata.
 - Work together to replicate data and distribute partitions.
 - Provide the backbone for other services like schema-registry, Kafka-connect, and KsqlDB-server.

3. Schema Registry

- **Purpose:** Centralizes schema management for Kafka topics.

- **Responsibilities:**
 - Stores schemas for Avro, JSON, and Protobuf data.
 - Ensures compatibility between producers and consumers.

- **Interactions:**
 - Works with Kafka brokers to store schemas in specific topics.
 - Enables serialization and deserialization for Kafka-connect and KsqlDB-server.

4. Kafka Connect

- **Purpose:** Facilitates data ingestion and export from Kafka.

- **Responsibilities:**
 - Connects external systems (e.g., databases, files) with Kafka.
 - Provides scalable and fault-tolerant connectors.

- **Interactions:**
 - Uses schema-registry for data serialization.
 - Relies on Kafka brokers for publishing and consuming messages.
 - Interacts with databases like MySQL to pull or push data.

5. KSQLDB (Server and CLI)

- **Purpose:** Enables real-time stream processing using SQL-like queries.

- **Responsibilities:**
 - Processes and transforms streaming data.
 - Creates derived topics or materialized views.

- **Interactions:**
 - Relies on Kafka brokers for streaming data.
 - Uses schema-registry for managing schemas.
 - Works with Kafka-connect for external data integration.

6. Kafka UI

- **Purpose:** Web interface to manage and monitor Kafka clusters.

- **Responsibilities:**
 - Visualizes Kafka topics, brokers, and consumer groups.
 - Simplifies topic management and monitoring.
- **Interactions:**
 - Connects to Kafka brokers for real-time data visualization.
 - Uses schema-registry and kafka-connect for enhanced insights.

7. Control Center

- **Purpose:** Monitoring and managing the Confluent platform.

- **Responsibilities:**
 - Provides a UI to monitor Kafka cluster health, latency, and throughput.
 - Supports managing topics, schemas, and connectors.
- **Interactions:**
 - Connects to Kafka brokers, schema-registry, Kafka-connect, and KsqlDB-server to consolidate management.

8. Pinot Components (Controller, Broker, Server)

- **Purpose:** Real-time OLAP (Online Analytical Processing) system for querying Kafka data.

- **Responsibilities:**
 - **Controller:** Manages Pinot cluster configuration.
 - **Broker:** Handles query routing.
 - **Server:** Stores and serves data segments.
- **Interactions:**
 - Pulls data from Kafka topics for real-time analysis.
 - Uses zookeeper for metadata and coordination.

9. MySQL

- **Purpose:** Relational database used for persistent storage.

- **Responsibilities:**
 - Acts as a data source or sink for Kafka via Kafka-connect.
- **Interactions:**
 - Integrates with Kafka for data ingestion/export.

10. Streamlit Dashboard

- **Purpose:** Frontend for visualizing analytics results.

- **Responsibilities:**
 - Displays real-time analytics and insights from Kafka data.
- **Interactions:**
 - Queries data from Pinot (pinot-broker).
 - Presents data to users through an intuitive UI.

*PS. Kafka UI and Control Connect service as the similar purpose, but control center is paid version (can trial some period) and it also included ui for ksql operation.

a.Data Source

Datagen
Source 1

quickstart : pageviews

Datagen
Source 2

quickstart : users_

MySQL
Source 3

Transaction Data :
flash_sales

Master Data :
product, region, pages

Source 1 : Using Kafka's built-in **data generation connectors** to simulate real-time data with quick start pageviews

Source 2 : Using Kafka's built-in **data generation connectors** to simulate real-time data with quick start user_

Source 3 : Using **mysql source connectors** to periodical load data according to poll interval.

There is 1 transactional event data

- Flash sales is the transaction record for each product sales event.

There are 3 master data

- Product : describe the product name and product category (group) from product_id
- Region : describe the region name from region_id
- Pages : describe the pages name and page category (group) from page_id

For flash sales event I use python code to continue iterate simulate sales transaction within 1 secs per iteration.



a.Data Source

What : Source is where data was created and then we push or poll data to kafka topic in real-time.

Why : The data can either be push by producer or poll from source directly by kafka connect.

How :

Input

- Source 1 & 2 **quick start data generation**
- Source 3 Mysql Relational Database

Process

- Source 1 & 2 Kafka's built-in **quick start data generation** connector generate and serialize data and operate on the **push** mechanism to continuously simulate and stream real-time data to Kafka topic 1 & 2 accordingly.
- Source 3 use **mysql source connectors** operates on the **poll** mechanism with periodical pulling data according to poll interval configuration. (in this example topic3 for flash_sales table) and extra 3 topics for master data

Output

- Serialize data located in kafka partition under its associate topic.

b.Kafka System : Kafka Broker

Brokers							
Uptime			Partitions				
Broker Count	Active Controller	Version	Online 522 of 522	URP 0	In Sync Replicas 746 of 746	Out Of Sync Replicas 0	
Broker ID	Disk usage	Partitions skew ⓘ	Leaders	Leader skew	Online partitions	Port	Host
1 ✓	1.01 GB, 248 segment(s)	-0.30%	174	-	248	9092	kafka1
2	1.01 GB, 246 segment(s)	-1.10%	172	-1.10%	246	9092	kafka2
3	1.01 GB, 252 segment(s)	1.30%	176	1.10%	252	9092	kafka3

According to the requirements we set-up 3 brokers

What : Broker is the main server where Kafka operate and store each message under its associated topic.

Why : Multiple brokers will help on load balancing (scalable) and fault-tolerant with replica (copy) data to other brokers if config replica > 1 in topic.

How: Leader and Follower Brokers

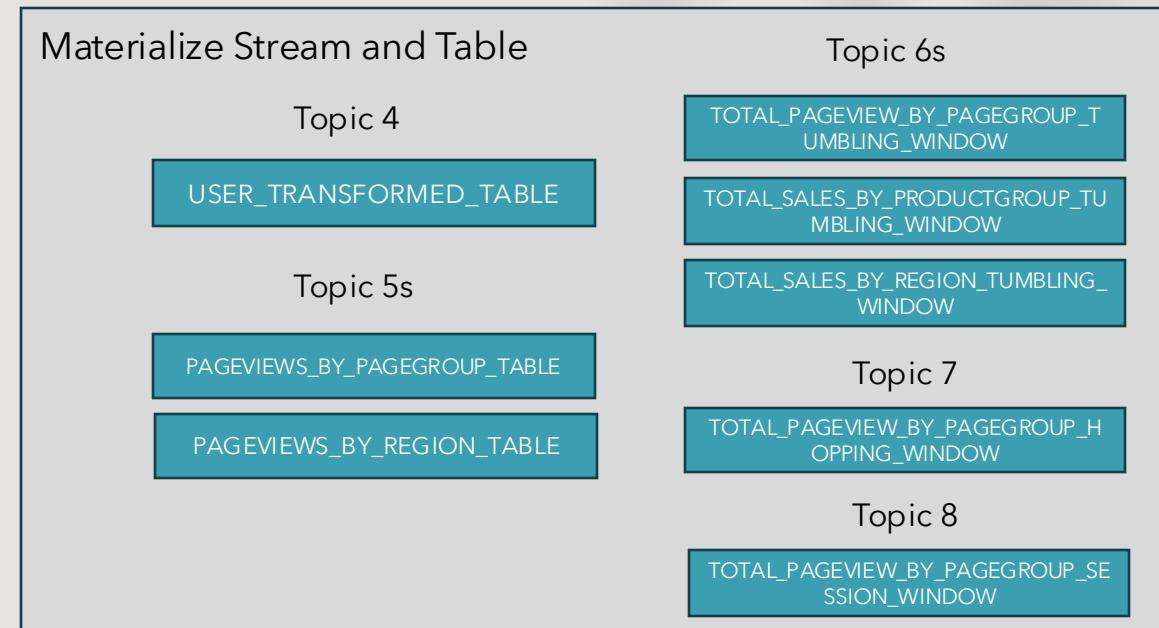
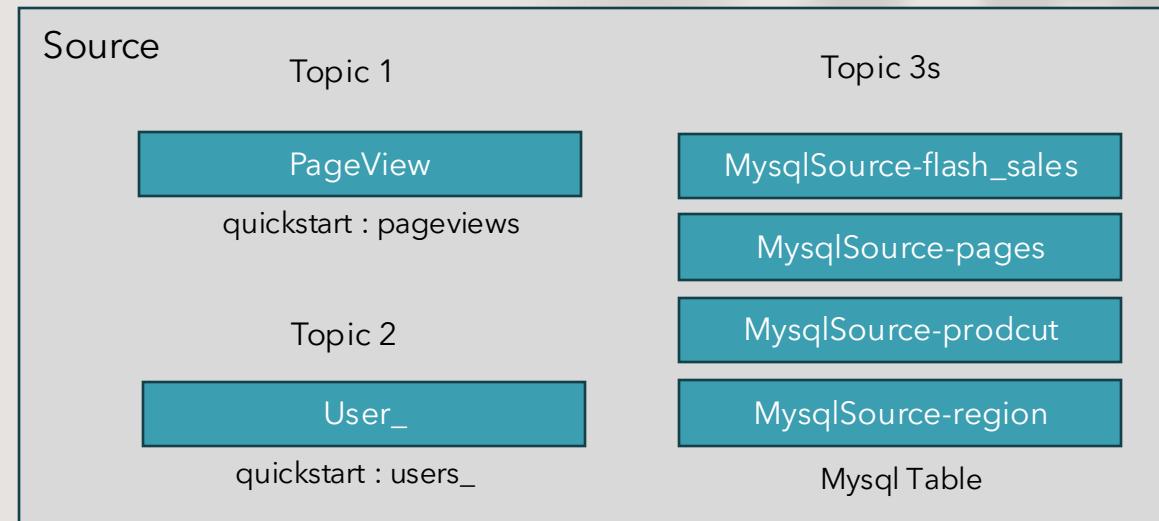
For each partition, **one broker is assigned as the leader**, while the others are followers. The leader broker handles all client requests (reads and writes) for its partition, while follower brokers replicate the data to stay synchronized.

- If a leader broker fails, one of the followers takes over as the new leader, ensuring uninterrupted data access.

b.Kafka System : Kafka Topics and Partition

I have created 16 topics, each topic has 5 partitions, and important topic (event base) has 2 replica and none importance (master data) has 1 replica.

Topics					
	Topic Name	Partitions	Out of sync replicas	Replication Factor	Number of messages
■	EXTENDED_FLASHSALES_STREAM	5	0	2	23544
■	EXTENDED_PAGEVIEWS_STREAM	5	0	2	29369
■	MysqlSource-flash_sales	5	0	2	34282
■	MysqlSource-pages	5	0	1	53640
■	MysqlSource-product	5	0	1	53550
■	MysqlSource-region	5	0	1	5355
■	PAGEVIEWS_BY_PAGEGROUP_TABLE	5	0	2	70159
■	PAGEVIEWS_BY_REGION_TABLE	5	0	2	29146
■	PageView	5	0	2	71313
■	TOTAL_PAGEVIEW_BY_PAGEGROUP_HOPPING_WINDOW	5	0	2	52946
■	TOTAL_PAGEVIEW_BY_PAGEGROUP_SESSION_WINDOW	5	0	2	52915
■	TOTAL_PAGEVIEW_BY_PAGEGROUP_TUMBLING_WINDOW	5	0	2	29051
■	TOTAL_SALES_BY_PRODUCTGROUP_TUMBLING_WINDOW	5	0	2	23473
■	TOTAL_SALES_BY_REGION_TUMBLING_WINDOW	5	0	2	21979
■	USER_TRANSFORMED_TABLE	5	0	1	13670
■	Users_	5	0	2	23702



b.Kafka System : Kafka Topics and Partition

What :

Topic is a category or feed name where records (messages) are sent by producers and read by consumers.

- It's similar to a queue in other messaging systems but with key differences that allow Kafka to handle distributed and durable message storage.

Partition is a subset of a topic and represents a unit of parallelism.

- Each topic is divided into one or more partitions, which are distributed across Kafka brokers.

Why :

Kafka topics and partitions are essential for Kafka's scalability, performance, and reliability.

Topic:

- **Purpose:** Organizes data logically, ensuring related events are grouped together.
- **Importance:** Makes it easy for producers and consumers to handle specific data streams. For example:
 - Topic PageView: Handles user view page on website events.
 - Topic User_: Handles user registration events.

Partition:

- **Purpose:** Distributes data across brokers for scalability and increase performance of read and write parallelly.
- **Importance:**
 - **Scalability:** Partitions enable Kafka to handle large volumes of data by distributing it across multiple brokers.
 - **Parallelism:** Multiple consumers can process partitions independently, improving performance.
 - **Fault Tolerance:** Kafka replicates partitions across brokers to ensure data durability and availability.

How :

Input (Data) : Producer writes data to a topic -> Partition assignment

Process :

- Each partition stores messages in an append-only log with a unique offset (sequence number) for each message.
- Kafka replicates each partition across multiple brokers for fault tolerance.
- Consumers subscribe to a topic and pull messages from partitions.

Output :

- Messages are stored in the topic partitions and replicated for durability.
- Each consumer reads from one or more partitions, enabling distributed processing of the data stream.
- Kafka retains messages for a configurable period (e.g., 7 days), allowing consumers to replay or reprocess data if needed.

b.Kafka System : Additional Explanation

In this setup, I have configured **3 brokers** to provide **load balancing** and **fault tolerance**.

1. Partitions

- I have created topic with **5 partitions**. This partitioning improves **performance** by spreading the data across multiple brokers, which allows parallel read and write operations.
- By having 5 partitions, the data storage is split into 5 segments, which can be distributed across different brokers. This division increases throughput for both reading and writing data, as different clients can access separate partitions simultaneously.

2. Replicas

- **Replication** is a mechanism to achieve fault tolerance in Kafka. Each partition has a specified number of replicas to ensure data availability in case of broker failures.

In my setup:

- **Event transactions (important topics)** have set **replication factor of 2**. This means each partition in this topic has two copies on different brokers. If one broker fails, another broker with a replica can still serve the data, ensuring reliability.
- **Master data (less critical topics)** have a **replication factor of 1**. This means there is only one copy of each partition, reducing storage requirements at the cost of lower redundancy.

• Disadvantages of Replicas:

While replicas provide data redundancy, having more replicas consumes additional disk space and can impact performance, as each write operation must be synchronized across all replicas, increasing the I/O load.

3. Leader and Follower Brokers

For each partition, **one broker is assigned as the leader**, while the others are followers. The leader broker handles all client requests (reads and writes) for its partition, while follower brokers replicate the data to stay synchronized.

- If a leader broker fails, one of the followers takes over as the new leader, ensuring uninterrupted data access.

4. Cleanup Policy and Retention

In this setup, each topic is configured with a **delete** cleanup policy to manage data storage effectively. This policy removes data that exceeds the defined retention period.

- **Master Data:** The retention period for master data (less critical topics) is set to **1 day**. This helps reduce storage usage by retaining only the most recent data, as older records are less likely to be needed.
- **Transaction Data:** For transaction data (critical topics), the retention period is set to **7 days**. This ensures that important data remains available for a longer duration to support critical operations before being deleted.

b.Kafka System : Schema-Registry

What : **Schema Registry** is a central component that manages and stores the schemas used for serializing and deserializing data in Kafka. It is commonly used with formats like **Avro**, **Protobuf**, or **JSON Schema**.

- It provides a repository to store schemas, enabling producers and consumers to agree on the data structure (schema) for communication.
- Each schema is versioned, allowing you to track changes and ensure compatibility between producers and consumers.

Why : The Schema Registry plays a critical role in maintaining **data integrity, compatibility, and evolution** in distributed systems.

1. Ensures Data Consistency:

- Producers must write data that adheres to a defined schema.
- Consumers can safely read and deserialize data, knowing the schema structure.

2. Manages Schema Evolution:

- As schemas change over time (e.g., adding fields), the Schema Registry ensures changes are compatible (backward or forward compatibility).

3. Simplifies Consumer-Producer Communication:

- Producers and consumers don't need to embed schemas in every message. Instead, they refer to a schema stored in the Schema Registry using a schema ID.

4. Reduces Data Overhead:

- Instead of embedding the full schema with each message, only the schema ID is sent, reducing message size and bandwidth usage.

5. Prevents Errors:

- Validates messages against the schema before they are sent to Kafka, avoiding corrupt or malformed data.

How :

Input:

- A producer sends page view event with the schema { "viewtime": "long", "userid": "string", "pageid": "string" } to Kafka.
- This schema is registered with the Schema Registry as Schema ID 2.

Process:

- The producer serializes the event using Schema ID 2 and sends it to a Kafka topic.
- The consumer reads the event and retrieves Schema ID 2 from the Schema Registry to deserialize it.

Output:

- The consumer processes the deserialized message: { "viewtime": 1, "userid": "User_2", "pageid": "Page_53" }.

Subject	Id	Type	Version	Compatibility
EXTENDED_FLASHSALES_STREAM-value	16	AVRO	1	BACKWARD
EXTENDED_PAGEVIEWS_STREAM-value	9	AVRO	1	BACKWARD
MysqlSource-flash_sales-value	3	AVRO	1	BACKWARD
MysqlSource-pages-value	4	AVRO	1	BACKWARD
MysqlSource-product-value	5	AVRO	1	BACKWARD
MysqlSource-region-value	6	AVRO	1	BACKWARD
PAGEVIEWS_BY_PAGEGROUP_TABLE-value	20	AVRO	1	BACKWARD
PAGEVIEWS_BY_REGION_TABLE-value	22	AVRO	1	BACKWARD
PageView-value	2	AVRO	1	BACKWARD
TOTAL_PAGEVIEW_BY_PAGEGROUP_HOPPING_WINDOW-value	35	AVRO	1	BACKWARD
TOTAL_PAGEVIEW_BY_PAGEGROUP_SESSION_WINDOW-value	35	AVRO	1	BACKWARD
TOTAL_PAGEVIEW_BY_PAGEGROUP_TUMBLING_WINDOW-value	26	AVRO	1	BACKWARD
TOTAL_SALES_BY_PRODUCTGROUP_TUMBLING_WINDOW-value	29	AVRO	1	BACKWARD
TOTAL_SALES_BY_REGION_TUMBLING_WINDOW-value	36	AVRO	3	BACKWARD
USER_TRANSFORMED_TABLE-value	7	AVRO	1	BACKWARD
Users_-value	1	AVRO	1	BACKWARD

Offset	Partition	Timestamp
0	4	11/15/2024, 12:24:08
Key	Value	Headers
	{"viewtime": 1, "userid": "User_2", "pageid": "Page_53"}	


```
"type": "record",
"name": "pageviews",
"namespace": "ksql",
"fields": [
  {
    "name": "viewtime",
    "type": "long"
  },
  {
    "name": "userid",
    "type": "string"
  },
  {
    "name": "pageid",
    "type": "string"
  }
],
"connect.name": "ksql.pageviews"
```

b.Kafka System : Kafka Connect

What : Kafka Connect acts as a framework to integrate Kafka brokers with external systems. It uses connectors that can be configured as **source connectors** to ingest data into Kafka or **sink connectors** to export data from Kafka to external systems. Kafka Connect typically uses a **polling mechanism** for data movement.

Why : To store and consume data in real-time with low latency we use Kafka services, with external source we require connector who act as a messenger between broker and source/sink.

How : Input and process : (in my configuration)

- **Datagen as Source (pageviews and users):**

- Using Kafka's built-in **data generation connectors** to simulate real-time data for pageviews and users.
- These connectors generate data at configurable time intervals and send it directly to Kafka topics.

- **MySQL as a Source**

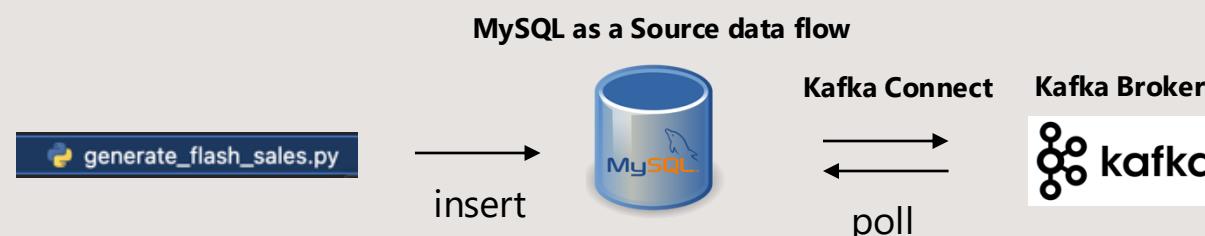
- **Flash Sales Transactions:**

- The flash sales transaction data is set up with an **incremental mode** using the sale_id column. This mode checks for new entries based on the sale_id. If a new sale_id is detected, it is loaded into the corresponding Kafka topic (MysqlSource-flash_sales).
 - The **poll interval** for this load is configured at **1000 milliseconds** (1 second), meaning Kafka will check MySQL every second for new sales transactions.

- **Master Data (product, pages, and region):**

- For static master data, I've set for **bulk mode**. In bulk mode, Kafka performs a full load of the data at each poll interval, rather than checking for individual updates.
 - The **poll interval** is set to **60000 milliseconds** (1 minute), meaning every minute, Kafka fetches the latest data from MySQL in bulk and updates the topic. In real-life master data will be flush and load on daily interval.

Output : Data arrived serialized and stored in Kafka partition under its associated topic, ready for consumer to query.



b.Kafka System : Connectors for data source

Connector name	Status	Category	Plugin name	Topics	Number of tasks
mysql-source-pages	Running	Source	JdbcSourceConn...	--	1
datagen-source-users_	Running	Source	DatagenConnector	--	1
mysql-source-region	Running	Source	JdbcSourceConn...	--	1
mysql-source-product	Running	Source	JdbcSourceConn...	--	1
mysql-source-flashesales	Running	Source	JdbcSourceConn...	--	1
datagen-source-pageviews	Running	Source	DatagenConnector	--	1

datagen-source-pageviews

- simulate page view event every 1 secs

datagen-source-users_

- simulate users_ registration event every 3 secs

mysql-source-product

Bulk load data from mysql table pages every 60 secs

```
"table.whitelist": "product",
"mode": "bulk",
"topic.prefix": "MysqlSource-",
"poll.interval.ms": "60000",
```

mysql-source-flashesales

Incremental load data from mysql table flash_sales every 1 secs

```
"table.whitelist": "flash_sales",
"mode": "incrementing",
"incrementing.column.name": "sale_id",
"topic.prefix": "MysqlSource-",
"poll.interval.ms": "1000",
```

mysql-source-pages

Bulk load data from mysql table pages every 60 secs

```
"table.whitelist": "pages",
"mode": "bulk",
"topic.prefix": "MysqlSource-",
"poll.interval.ms": "60000",
```

mysql-source-region

Bulk load data from mysql table region every 60 secs

```
table.whitelist": "region",
"mode": "bulk",
"topic.prefix": "MysqlSource-",
"poll.interval.ms": "60000",
```

c. KsqlDB



What : Ksqldb is an event streaming database built on top of Apache Kafka. It allows us to work with real-time data streams using SQL-like queries. With ksqlDB, we can process, transform, and analyze streaming data in a declarative and user-friendly way, without needing to write complex code.

Why : Materialize View

- **Streams** are continuous, appending new rows as events occur, suitable for real-time processing of events. (e.g., 2,3.1, e)
- **Tables** capture the current state, providing the latest data per primary key, which is ideal for static or master data that updates over time. (e.g., 1,3,5)

Windowed Aggregations Table:

- The setup also includes **windowed tables** (e.g., 6,7,8), which aggregate data over specific time windows. These tables are useful for time-based analyses, like tracking session activities or events within a particular timeframe (tumbling , hopping and session windows)

c. KsqlDB

- **Stream as a Materialized View**

A **stream** in ksqlDB represents a continuous, append-only flow of data. While it is not inherently a “materialized view”, when query a stream, ksqlDB processes the data in real-time and can materialize the results for a query.

Characteristics:

1. **Unbounded Data:** Streams represent an unending flow of events.
2. **Not State-Based:** Streams are focused on capturing raw, unaggregated data.
3. **Transient Materialization:** When you run a query on a stream, ksqlDB processes the events in real-time and temporarily materializes the results for the duration of the query.
4. **No Persistent State:** Streams do not store state or maintain results over time unless explicitly written to a Kafka topic or derived into a table.

- **Table as a Materialized View**

A **table** in ksqlDB is explicitly a materialized view. It represents the latest state of data, where incoming records are aggregated or updated based on a key.

Characteristics:

1. **State-Based:** Tables maintain a persistent, stateful view of the latest data.
2. **Materialized by Default:** Tables store their state (e.g., aggregations, counts) in an internal **RocksDB store** for fast lookups and updates.
3. **Supports Updates:** Tables can update their state when new data arrives, similar to how a database table works.
4. **Queryable State:** Tables allow you to execute point-in-time queries against the current state.

How :

Input: A producer pushes data into a Kafka topic, which serves as the source for ksqlDB queries.

Process: ksqlDB queries process data from the topic and may create materialized views for stateful operations. If the query outputs transformed, cleaned, or aggregated data, ksqlDB writes the results to a **new Kafka topic** for further use.

Output: Queries produce results either as **materialized views** (internally stored for efficient lookups) or as **new Kafka topics**, depending on the query type and configuration.

c. KsqlDB : Clean or transform data (topic4)

```
CREATE TABLE USER_TRANSFORMED_TABLE AS
SELECT
userid,
TIMESTAMPTOSTRING(registertime, 'yyyy-MM-dd HH:mm:ss') AS registertime,
gender,
regionid,
interests[1] AS interests_1,
interests[2] AS interests_2,
contactinfo['phone'] AS phone,
contactinfo['city'] AS city,
contactinfo['state'] AS state,
contactinfo['zipcode'] AS zipcode
FROM USERS_TABLE
```

What : Sometime our data came with unmeaningful or not ready to use we may need to perform clean or transformation data like joining key to master table or covert value to be easy for visualize.

Process :

This materialize process transformation from USER_TABLE and put to USER_TRANSFORMED_TABLE

- Transform registertime from unixtimestamp to string format with human readable.
- Transform Interests with array of 2 key value pair to each individual column interest1 and interest2.
- Transform contactinfo by mapping it separate columns

input

Value	Header	Key
1		{
2		"registertime": 1509337823511,
3		"userid": "User_5",
4		"regionid": "Region_7",
5		"gender": "OTHER",
6		"interests": [
7		"News",
8		"Travel"
9],
10		"contactinfo": {
11		"zipcode": "94403",
12		"state": "CA",
13		"phone": "6503349999",
14		"city": "San Mateo"
15		}
16		}

Output

{	"REGISTERTIME": {	"string": "2017-08-07 06:39:59"
,	"GENDER": {	"string": "FEMALE"
,	"REGIONID": {	"string": "Region_7"
,	"INTERESTS_1": {	"string": "News"
,	"INTERESTS_2": {	"string": "Movies"
,	"PHONE": {	"string": "6503349999"
,	"CITY": {	"string": "San Mateo"
,	"STATE": {	"string": "CA"
,	"ZIPCODE": {	"string": "94403"
}		}

Userid	Registertime	Regionid	Gender	Interests	ContactInfo
User_2	1518474649324	Region_6	FEMALE	["Game","News"]	{"zipcode":"94403","phone":"650334..."}
User_7	1506229160658	Region_8	MALE	["News","Travel"]	{"zipcode":"94301","phone":"6503889..."}
User_9	1513798275186	Region_6	MALE	["News","Travel"]	{"zipcode":"94301","phone":"6503889..."}

Userid	Registertime	Gender	Regionid	Interests_1	Interests_2	Phone	City	State	Zipcode
User_5	2017-08-28 19:47...	MALE	Region_3	News	Travel	4083366881	San Jose	CA	95112
User_6	2017-05-04 18:0...	MALE	Region_4	Game	Sport	4083366881	San Jose	CA	95112

c. KsqlDB : Aggregation (join + group by) (topic5)

```
CREATE TABLE PAGEVIEWS_BY_PAGEGROUP_TABLE AS
SELECT
    pt.page_group as key,
    CAST(pt.page_group AS STRING) AS page_group,
    COUNT(*) AS count_view
FROM pageviews_stream pvs
LEFT JOIN pages_table pt ON pt.page_id = pvs.pageid
GROUP BY pt.page_group
EMIT CHANGES;
```

What : Sometime our data came with lower granularity (transaction level) and we may need to perform aggregation to be easy for understand big picture of data and describable as per each dimension.

Process :

This materialize process join pageviews_stream with pages_table to get page_group attribute and aggregate by count(*) count view per each page_group. And put to PAGEVIEWS_BY_PAGEGROUP_TABLE

Input

1 select * from PAGEVIEWS_STREAM EMIT CHANGES;

● Add query properties
auto.offset.reset = Latest

Running... Stop

Data structure STREAM
Total messages
Messages/sec
Total message bytes
Message fields

VIEWTIME	USERID	PAGEID
1865601	User_8	Page_44
1865591	User_9	Page_15
1865581	User_1	Page_72
1865571	User_7	Page_10

The above statistics are computed for the

Output

1 select * from PAGEVIEWS_BY_PAGEGROUP_TABLE EMIT CHANGES;

● Add query properties
auto.offset.reset = Latest

Running... Stop

Data structure TABLE
Total messages
Messages/sec
Total message bytes
Message fields

KEY	PAGE_GROUP	COUNT_VIEW
Dry Grocery	Dry Grocery	40997
Furniture	Furniture	41408
Fashion	Fashion	41276

c. KsqlDB : Tumbling Window (Topic 6s)

A **Tumbling Window** is a type of time window that divides time into fixed intervals with no overlap between the windows. Each event belongs to exactly one window, making it easy to calculate metrics such as total sales for a specific time period.

```
CREATE TABLE TOTAL_SALES_BY_PRODUCTGROUP_TUMBLING_WINDOW AS
SELECT
    product_group AS key,
    CAST(product_group AS STRING) AS product_group,
    SUM(sale_price * quantity) AS total_sales,
    TIMESTAMPTOSTRING(WINDOWSTART, 'yyyy-MM-dd HH:mm:ss', 'UTC') AS window_start,
    TIMESTAMPTOSTRING(WINDOWEND, 'yyyy-MM-dd HH:mm:ss', 'UTC') AS window_end
FROM EXTENDED_FLASHSALES_STREAM
WINDOW TUMBLING (SIZE 1 HOUR)
GROUP BY product_group
EMIT CHANGES;
```

In this example, ksqlDB calculates the total sales for each product_group every hour. Each window is distinct and does not overlap with others.

Scenario:

- **Window 1:** From 9:00 to 10:00 — calculates total sales in this period only.
- **Window 2:** From 10:00 to 11:00 — calculates total sales for this interval only.
- **Window 3:** From 11:00 to 12:00 — calculates total sales for this interval only.

If a sale occurs at 10:15, it will be included in **Window 2 (10:00 to 11:00)**. Each window displays the total sales for that specific period.

c. KsqlDB : Hopping Window (Topic 7)

A **Hopping Window** is a type of time window that allows for overlapping intervals. It is defined by a **size** (duration of the window) and an **advance interval** (how frequently the window starts). Events can belong to multiple windows due to the overlap.

```
CREATE TABLE TOTAL_SALES_BY_PRODUCTGROUP_HOPPING_WINDOW AS
SELECT
    product_group AS key,
    CAST(product_group AS STRING) AS product_group,
    SUM(sale_price * quantity) AS total_sales,
    TIMESTAMPTOSTRING(WINDOWSTART, 'yyyy-MM-dd HH:mm:ss', 'UTC') AS window_start,
    TIMESTAMPTOSTRING(WINDOWEND, 'yyyy-MM-dd HH:mm:ss', 'UTC') AS window_end
FROM EXTENDED_FLASHSALES_STREAM
WINDOW HOPPING (SIZE 1 HOUR, ADVANCE BY 30 MINUTES)
GROUP BY product_group
EMIT CHANGES;
```

In this example, ksqlDB calculates the total sales for each product_group using a 1-hour window that advances every 30 minutes. Each window overlaps with the previous and next window.

Scenario:

- **Window 1:** From 9:00 to 10:00 — calculates total sales in this period.
- **Window 2:** From 9:30 to 10:30 — calculates total sales for this overlapping interval.
- **Window 3:** From 10:00 to 11:00 — calculates total sales for this interval.
- **Window 4:** From 10:30 to 11:30 — calculates total sales for this interval.

If a sale occurs at 9:45:

- It will be included in **Window 1 (9:00 to 10:00)** and **Window 2 (9:30 to 10:30)** because of the overlap.

Each window displays the total sales for its respective overlapping period.

c. KsqlDB : Session Window (Topic 8)

A **Session Window** groups events based on activity rather than fixed intervals. The window remains open as long as events continue to arrive within a specified **inactivity gap**. Once no events are received for the duration of the gap, the session closes.

```
CREATE TABLE TOTAL_SALES_BY_PRODUCTGROUP_SESSION_WINDOW AS
SELECT
    product_group AS key,
    CAST(product_group AS STRING) AS product_group,
    SUM(sale_price * quantity) AS total_sales,
    TIMESTAMPTOSTRING(WINDOWSTART, 'yyyy-MM-dd HH:mm:ss', 'UTC') AS window_start,
    TIMESTAMPTOSTRING(WINDOWEND, 'yyyy-MM-dd HH:mm:ss', 'UTC') AS window_end
FROM EXTENDED_FLASHSALES_STREAM
WINDOW SESSION (15 MINUTES)
GROUP BY product_group
EMIT CHANGES;
```

In this example, ksqlDB calculates the total sales for each product_group within session windows. A session starts when an event occurs and remains active as long as new events arrive within a 15-minute gap. If no events are received within the gap, the session closes.

Scenario:

- **Session 1:** Starts at 9:00 with sales at 9:05, 9:10, and 9:12.
 - If no further sales occur within the next 15 minutes, the session ends at 9:12.
- **Session 2:** Starts at 9:30 with sales at 9:32 and 9:40.
 - Ends if no sales occur by 9:55.

If a sale occurs at 9:10 and another at 9:12:

- Both events belong to **Session 1** because the gap between them is less than 15 minutes.

If another sale occurs at 9:30:

- A **new session (Session 2)** starts at 9:30, as the gap from the last sale in Session 1 is greater than 15 minutes.

c. KsqlDB : Testing the correctness

statements.sql

```
CREATE STREAM USER_DATA (
    registertime BIGINT,
    userid VARCHAR,
    regionid VARCHAR,
    gender VARCHAR,
    interests ARRAY<VARCHAR>,
    contactinfo STRUCT<
        zipcode VARCHAR,
        state VARCHAR,
        phone VARCHAR,
        city VARCHAR>
) WITH (
    KAFKA_TOPIC = 'User_',
    VALUE_FORMAT = 'AVRO',
    TIMESTAMP = 'registertime'
);

CREATE TABLE USERS_PER_REGION AS
SELECT regionid,
    COUNT(*) AS USER_COUNT
FROM USER_DATA
GROUP BY regionid
EMIT CHANGES;
```

output.json

```
{
    "outputs": [
        {
            "topic": "USERS_PER_REGION",
            "key": "Region_1",
            "value": {"USER_COUNT": 1},
            "timestamp": 1500023204411
        },
        {
            "topic": "USERS_PER_REGION",
            "key": "Region_2",
            "value": {"USER_COUNT": 1},
            "timestamp": 1500023204412
        },
        {
            "topic": "USERS_PER_REGION",
            "key": "Region_1",
            "value": {"USER_COUNT": 2},
            "timestamp": 1500023204413
        }
    ]
}
```

Input.json

```
{
    "inputs": [
        {
            "topic": "User_",
            "key": null,
            "value": {
                "registertime": 1500023204411,
                "userid": "User_1",
                "regionid": "Region_1",
                "gender": "MALE",
                "interests": ["Sports", "Music"],
                "contactinfo": {
                    "zipcode": "12345",
                    "state": "NY",
                    "phone": "1234567890",
                    "city": "New York"
                }
            }
        },
        {
            "topic": "User_",
            "key": null,
            "value": {
                "registertime": 1500023204412,
                "userid": "User_2",
                "regionid": "Region_2",
                "gender": "FEMALE",
                "interests": ["Movies", "Travel"],
                "contactinfo": {
                    "zipcode": "54321",
                    "state": "CA",
                    "phone": "0987654321",
                    "city": "Los Angeles"
                }
            }
        }
    ]
}
```

```
(base) ek@ek-docker:~/dads6005/midterm/volumes$ docker exec ksqldb-cli ksql-test-runner -i /opt/app/test/input.json -s /opt/app/src/statements.sql -o /opt/app/test/output.json
    >>> Test passed!
(base) ek@ek-docker:~/dads6005/midterm/volumes$
```

What : ksql-test-runner

A tool designed for testing ksqlDB queries without requiring a live Kafka cluster. It takes static input, runs the SQL, and validates the output.

Tests ksqlDB queries (statements.sql) against mock input (input.json) and expected output (output.json).

Why : Ensures queries work as intended, preventing issues in production and enabling CI/CD integration.

How : 1. Takes input data. 2. Runs SQL queries. 3. Compares results to expected output.

"Test passed!" if results match; error output otherwise.

d. Apache Pinot

The screenshot shows the Apache Pinot Query Console interface. At the top right, it displays "CLUSTER NAME PinotCluster". The main area is divided into several sections:

- SQL EDITOR:** Contains the query: `1 select * from EXTENDED_FLASHSALES_STREAM limit 10`. Below the editor are checkboxes for "Tracing", "Use V2 Engine", and a "Timeout (in Milliseconds)" input field. A "RUN QUERY" button is located to the right.
- QUERY RESPONSE STATS:** Displays performance metrics for the last query execution. The table includes columns: timeUsedMs, numDocsScanned, totalDocs, numServersQueried, numServersResponded, numSegmentsQueried, numSegmentsProcessed, numSegmentsMatched, numConsumingSegmentsQueried, numEntriesScannedInFilter, and numEntriesScannedPos. The data row shown is: 6, 10, 29772, 1, 1, 5, 1, 1, 5, 0, 120.
- EXTENDED_FLASHSALES_STREAM SCHEMA:** Shows the schema for the EXTENDED_FLASHSALES_STREAM table. It includes columns: CUSTOMER_ID, DISCOUNT, FS_PRODUCT_ID, PRODUCT_GROUP, PRODUCT_NAME, PROVINCE_ID, QUANTITY, REGION_NAME, SALES_AMOUNT, SALE_DATETIME, SALE_ID, and SALE_PRICE. The data rows listed are:

CUSTOMER_ID	DISCOUNT	FS_PRODUCT_ID	PRODUCT_GROUP	PRODUCT_NAME	PROVINCE_ID	QUANTITY	REGION_NAME	SALES_AMOUNT	SALE_DATETIME	Sale_ID	Sale_Price
20131155	19	P98	Other	Fitness Equipment	41	8	Eastern	1582.4	1731653249000	13615	197.8
85201879	35.86	P77	Dry Grocery	Tea & Coffee	75	6	Central	661.62	1731653175000	13618	110.27
40070700	9.41	P45	Fashion	Casual Wear	21	7	Central	2395.68	1731653361000	13619	342.24
58156650	46.03	P17	Electronics	Portable Chargers	26	6	Eastern	551.4	1731653164000	13620	91.9
57766604	43.16	P14	Electronics	Digital Cameras	46	3	Central	366.78	1731653273000	13621	122.26
64172175	37.3	P25	Electronics	Virtual Reality Headsets	66	7	Central	787.78	1731653262000	13622	112.54
96783588	16.92	P88	Dry Grocery	Flour & Baking Essentials	41	9	Central	2084.58	1731653286000	13624	231.62
84199343	34.73	P51	Furniture	Bedroom Furniture	33	4	Central	472.6	1731653328000	13629	118.15
- ROWS PER PAGE:** Set to 10, showing 1-10 of the results.

d. Apache Pinot

What : Serving as online alnalytic processing (OLAP) layer which able to serve huge traffic and load for data reading with good performance.

Why : because real-time data processing and dashboard need a low latency compute and frequency, so Pinot is built for serving this.

How : When create real-time table in Pinot server it will consume data from Kafka topics (from each brokers) as real-time and keep data within it for a better performance on query.

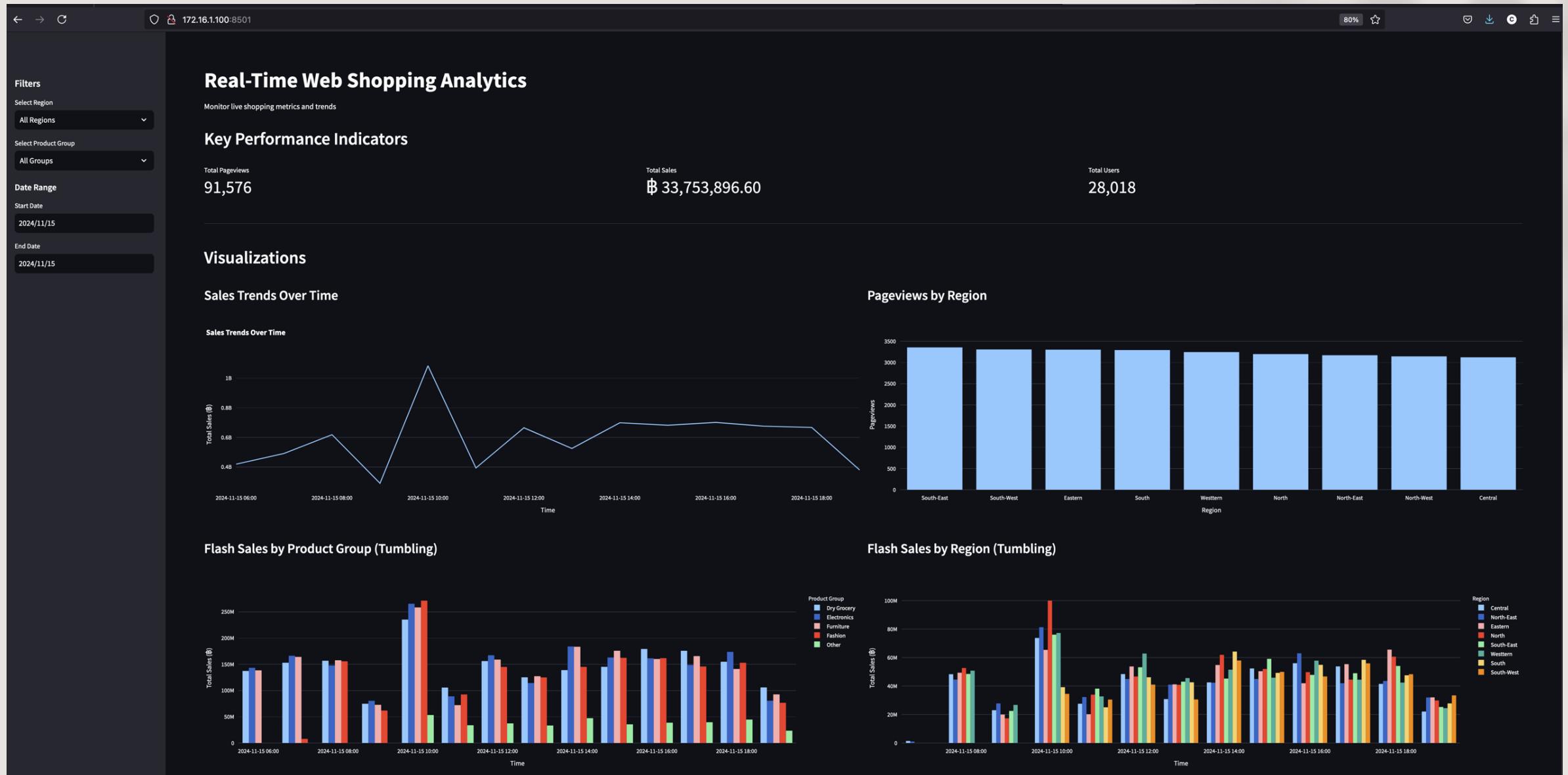
By creating real-time table in Pinot, we must have topic create in Kafka first and then create Schema align with Kafka schema. The schema name should be the same name as table name in case of adding table by Pinot Controller (UI).

I have create real time table with flushing old data every 24 hours, and each segment size as of 100 MB.

All of my table use AVRO format which require Schema-Registry server for identify the message schema for a correct deserialize message.

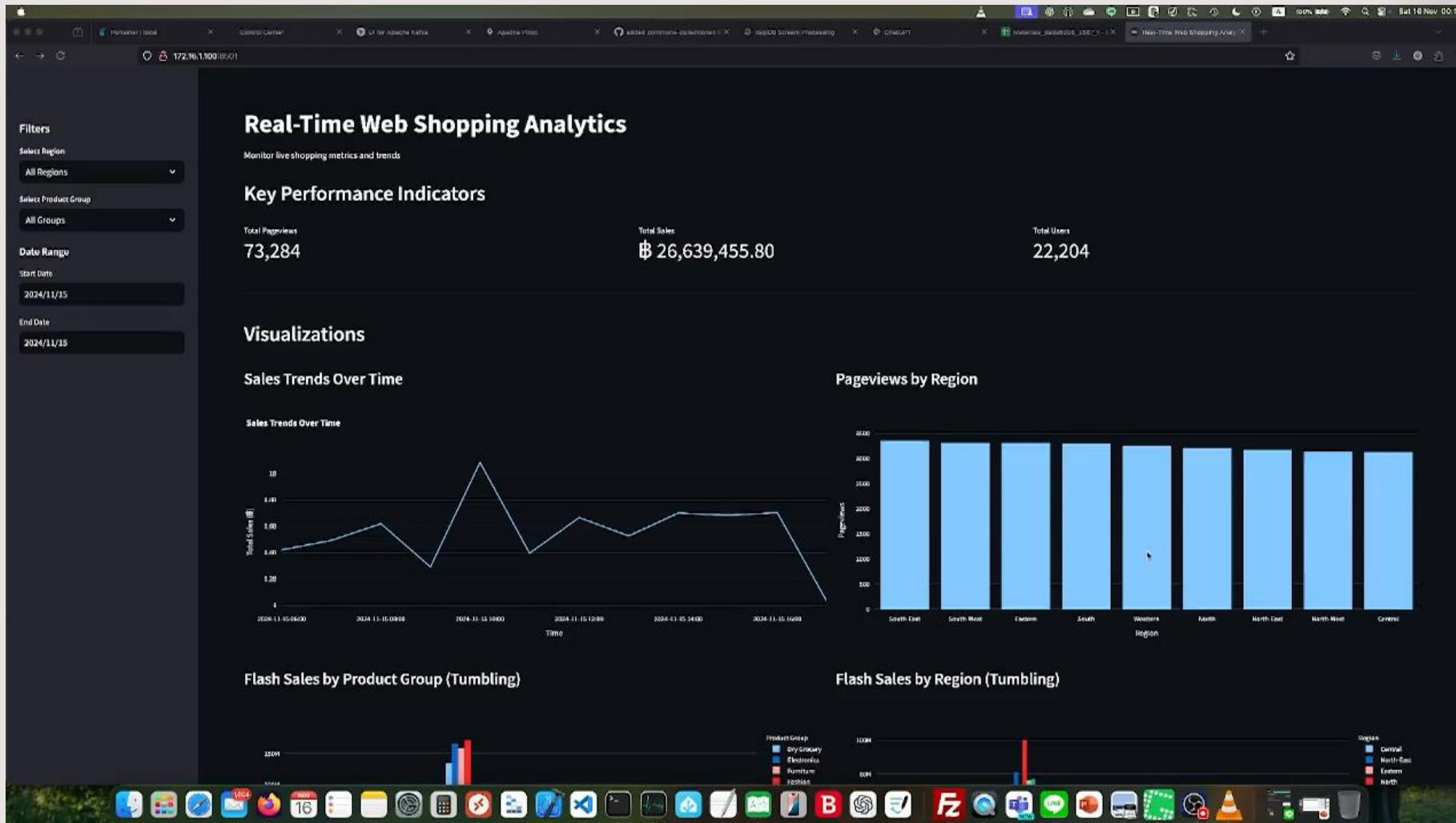
```
"tableIndexConfig": {  
    "streamConfigs": {  
        "streamType": "kafka",  
        "stream.kafka.topic.name": "EXTENDED_FLASHSALES_STREAM",  
        "stream.kafka.broker.list": "kafka1:9092,kafka2:9092,kafka3:9092",  
        "stream.kafka.consumer.type": "lowlevel",  
        "stream.kafka.consumer.prop.auto.offset.reset": "smallest",  
        "stream.kafka.consumer.factory.class.name": "org.apache.pinot.plugin.stream.kafka20.KafkaConsumerFactory",  
        "stream.kafka.decoder.class.name": "org.apache.pinot.plugin.inputformat.avro.confluent.KafkaConfluentSchemaRegistryAvroMessageDecoder",  
        "realtime.segment.flush.threshold.rows": "0",  
        "realtime.segment.flush.threshold.time": "24h",  
        "realtime.segment.flush.threshold.segment.size": "100M",  
        "stream.kafka.decoder.prop.schema.registry.rest.url": "http://schema-registry:8081",  
        "stream.kafka.decoder.prop.format": "AVRO"  
    },
```

e. Streamlit Dashboard



Quick Demo with Interactive Dashboard

quick_demo_interactive_real-time_dashboard.mp4



e. Streamlit Dashboard

What : Streamlit use for data visualization (dashboard) in real-time.

Why : Because our kafka plattform is for real-time data processing so we need dashboard wher we can visualize our data in real-time as well.

How : The dashbaord written by python with library streamlit, plotly, pandas, pinotdb ; it comsume data from Apache Pinot brokers (OLAP).

We createa real-time table in Pinot which read data from Kafka and ingestion to Pinot (data will be keep in Pinot server) which serving as online alnalytic processing (OLAP), by using library pinotdb we can connect and query data from pinot with sql statement and load into pandas data frame then plot visualization with plotly. And all of this process is executing and operation by python with streamlit library.

There 2 parts in dashboard

- Key Performace Indicator is show the metrics of Total Page Views, Total Sales and Total Users.
- Visualization chart there are 4 Charts,
 - Overall Sales Trend
 - PageView By Region (Topic 5 Agrregate Count PageView, Group by Region)
 - Flash Sales By Product Group (Topic 6 Tumbling Windows)
 - Flash Sales By Region (Topic 6 Tumbling Windows)
- There are filter on the left side to see data in specific condition (Date, Product Group, Region Name)

The dashboard is set refresh with every 1 second.

Appendix

Sample Screenshot Topics for Source 1 & 2

Topics / PageView

Overview Messages Consumers Settings Statistics

Seek Type: Offset | Partitions: All items are selected | Key Serde: String | Value Serde: SchemaRegistry | Clear all | Submit

Search: + Add Filters

	Offset	Partition	Timestamp	Key Preview	Value Preview
<input checked="" type="checkbox"/>	0	4	11/15/2024, 12:24:08	1	{"viewtime":1,"userid":"User_2","pageid":"Page_53"}
				Key Value Headers	Timestamp: 11/15/2024, 12:24:08 Timestamp type: CREATE_TIMESTAMP
				<pre>{"viewtime": 1, "userid": "User_2", "pageid": "Page_53"}</pre>	Key Serde: String Size: 1 Bytes
					Value Serde: SchemaRegistry Size: 21 Bytes
<input checked="" type="checkbox"/>	0	3	11/15/2024, 12:24:09	11	{"viewtime":11,"userid":"User_8","pageid":"Page_8..."}
<input checked="" type="checkbox"/>	0	0	11/15/2024, 12:24:09	21	{"viewtime":21,"userid":"User_7","pageid":"Page_1..."}

Topics / Users

Overview Messages Consumers Settings Statistics

Seek Type: Offset | Partitions: All items are selected. | Key Serde: String | Value Serde: SchemaRegistry | Clear all | Submit

Search: + Add Filters

Offset	Partition	Timestamp	Key Preview	Value Preview
0	0	11/15/2024, 12:23:59	User_2	{"registertime":1513629010348,"userid":"User_2"...}

Key Value Headers

Timestamp: 11/15/2024, 12:23:59
Timestamp type: CREATE_TIME

Key Serde: String
Size: 6 Bytes

Value Serde: SchemaRegistry
Size: 103 Bytes

```
{ "registertime": 1513629010348, "userid": "User_2", "regionid": "Region_6", "gender": "MALE", "interests": [ "News", "Movies" ], "contactinfo": { "zipcode": "94403", "state": "CA", "phone": "6503349999", "city": "San Mateo" } }
```

Offset	Partition	Timestamp	Key Preview	Value Preview
1	0	11/15/2024, 12:24:00	User_2	{"registertime":1500023204411,"userid":"User_2"...}

Sample Screenshot Topics for Source 3

Topics / MysqlSource-flash_sales

Overview Messages Consumers Settings Statistics

Seek Type: Offset | Partitions: All items are selected. | Key Serde: String | Value Serde: SchemaRegistry | Clear all | Submit

Search: + Add Filters

DONE 0 ms 6 KB

Offset	Partition	Timestamp	Key	Preview
9017	2	11/16/2024, 03:27:58	44367	{"sale_id":44367,"customer_id":("long":60816779)...
8920	1	11/16/2024, 03:27:56	44366	{"sale_id":44366,"customer_id":("long":95212771)...
9016	2	11/16/2024, 03:27:54	44364	{"sale_id":44364,"customer_id":("long":215321721...}

Key Value Headers

```
{"sale_id": 44367, "customer_id": { "long": 60816779 }, "province_id": { "int": 42 }, "product_id": { "string": "P27" }, "region_id": { "string": "Region_6" }, "discount": { "bytes": 14.63 }, "original_price": { "bytes": 261.42 }, "sale_price": { "bytes": 223.17 }, "quantity": { "int": 10 }, "sale_datetime": { "long": "2024-11-1ST20:26:16Z" }}
```

Topics / MysqlSource-product

Overview Messages Consumers Settings Statistics

Seek Type: Offset | Partitions: All items are selected. | Key Serde: String | Value Serde: SchemaRegistry | Clear all | Submit

Search: + Add Filters

DONE 0 ms

Offset	Partition	Timestamp	Key	Preview
29523	0	11/16/2024, 11:33:55	P99	{"product_id": "P99","product_group": ("string": "Ot...")
20129	4	11/16/2024, 11:33:55	P89	{"product_id": "P89","product_group": ("string": "Dry...")

Key Value Headers

```
{"product_id": "P99", "product_group": { "string": "Other" }, "product_name": { "string": "Travel Accessories" }}
```

Topics / MysqlSource-region

Overview Messages Consumers Settings Statistics

Seek Type: Offset | Partitions: All items are selected. | Key Serde: String | Value Serde: SchemaRegistry | Clear all | Submit

Search: + Add Filters

DONE 0 ms

Offset	Partition	Timestamp	Key	Preview
1341	1	11/16/2024, 11:33:55	Region_5	{"region_id": "Region_5","region_name": ("string": "S...")}
1341	4	11/16/2024, 11:33:55	Region_9	{"region_id": "Region_9","region_name": ("string": "N...")}

Key Value Headers

```
{"region_id": "Region_5", "region_name": { "string": "South-East" }}
```

Topics / MysqlSource-pages

Overview Messages Consumers Settings Statistics

Seek Type: Offset | Partitions: All items are selected. | Key Serde: String | Value Serde: SchemaRegistry | Clear all | Submit

Search: + Add Filters

DONE 0 ms

Offset	Partition	Timestamp	Key	Preview
28181	0	11/16/2024, 11:32:55	Page_97	{"page_id": "Page_97","page_group": "Other","page..."}

Key Value Headers

```
{"page_id": "Page_97", "page_group": "Other", "page_name": "Garden & Plants"}
```

Appendix 1 : Output of Mysql Table

```
[mysql]> show tables;
+-----+
| Tables_in_dads6005 |
+-----+
| flash_sales
| pages
| product
| region
+-----+
4 rows in set (0.02 sec)
```

```
[mysql]> select * from region;
+-----+-----+
| region_id | region_name |
+-----+-----+
| Region_1 | Central
| Region_2 | North
| Region_3 | North-East
| Region_4 | Eastern
| Region_5 | South-East
| Region_6 | South
| Region_7 | South-West
| Region_8 | Western
| Region_9 | North-West
+-----+-----+
9 rows in set (0.00 sec)
```

```
[mysql]> select * from product limit 30;
+-----+-----+-----+
| product_id | product_group | product_name |
+-----+-----+-----+
| P10 | Electronics | Smartphone Accessories
| P11 | Electronics | Laptops & Notebooks
| P12 | Electronics | Headphones & Earbuds
| P13 | Electronics | Smart Home Devices
| P14 | Electronics | Digital Cameras
| P15 | Electronics | Gaming Consoles
| P16 | Electronics | Television & Accessories
| P17 | Electronics | Portable Chargers
| P18 | Electronics | Bluetooth Speakers
| P19 | Electronics | Wearable Technology
| P20 | Electronics | Computer Components
| P21 | Electronics | Networking Devices
| P22 | Electronics | Drones & Quadcopters
| P23 | Electronics | Smart Watches
| P24 | Electronics | Tablet Accessories
| P25 | Electronics | Virtual Reality Headsets
| P26 | Electronics | Home Audio Systems
| P27 | Electronics | Printers & Scanners
| P28 | Electronics | Projectors
| P29 | Electronics | e-Readers
| P30 | Fashion | Men's Clothing
| P31 | Fashion | Women's Clothing
| P32 | Fashion | Kids' Wear
| P33 | Fashion | Footwear
| P34 | Fashion | Bags & Luggage
| P35 | Fashion | Jewelry
| P36 | Fashion | Watches
| P37 | Fashion | Sunglasses
| P38 | Fashion | Hats & Caps
| P39 | Fashion | Belts
+-----+-----+-----+
30 rows in set (0.00 sec)
```

```
[mysql]> select * from pages limit 30;
+-----+-----+-----+
| page_id | page_group | page_name |
+-----+-----+-----+
| Page_10 | Electronics | Smartphone Accessories
| Page_11 | Electronics | Laptops & Notebooks
| Page_12 | Electronics | Headphones & Earbuds
| Page_13 | Electronics | Smart Home Devices
| Page_14 | Electronics | Digital Cameras
| Page_15 | Electronics | Gaming Consoles
| Page_16 | Electronics | Television & Accessories
| Page_17 | Electronics | Portable Chargers
| Page_18 | Electronics | Bluetooth Speakers
| Page_19 | Electronics | Wearable Technology
| Page_20 | Electronics | Computer Components
| Page_21 | Electronics | Networking Devices
| Page_22 | Electronics | Drones & Quadcopters
| Page_23 | Electronics | Smart Watches
| Page_24 | Electronics | Tablet Accessories
| Page_25 | Electronics | Virtual Reality Headsets
| Page_26 | Electronics | Home Audio Systems
| Page_27 | Electronics | Printers & Scanners
| Page_28 | Electronics | Projectors
| Page_29 | Electronics | e-Readers
| Page_30 | Fashion | Men's Clothing
| Page_31 | Fashion | Women's Clothing
| Page_32 | Fashion | Kids' Wear
| Page_33 | Fashion | Footwear
| Page_34 | Fashion | Bags & Luggage
| Page_35 | Fashion | Jewelry
| Page_36 | Fashion | Watches
| Page_37 | Fashion | Sunglasses
| Page_38 | Fashion | Hats & Caps
| Page_39 | Fashion | Belts
+-----+-----+-----+
30 rows in set (0.00 sec)
```

```
[mysql]
[mysql]> select * from flash_sales order by sale_id desc limit 10;
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| sale_id | customer_id | province_id | product_id | region_id | discount | original_price | sale_price | quantity | sale_datetime |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 44367 | 60816779 | 42 | P27 | Region_6 | 14.63 | 261.42 | 223.17 | 10 | 2024-11-15 20:26:16
| 44366 | 95212771 | 64 | P11 | Region_4 | 48.18 | 486.89 | 252.31 | 2 | 2024-11-15 20:24:19
| 44365 | 74770033 | 58 | P16 | Region_8 | 21.11 | 216.88 | 171.10 | 8 | 2024-11-15 20:23:38
| 44364 | 21532172 | 76 | P56 | Region_5 | 20.85 | 386.04 | 305.55 | 4 | 2024-11-15 20:26:53
| 44363 | 88396172 | 57 | P65 | Region_1 | 45.65 | 320.77 | 174.34 | 1 | 2024-11-15 20:24:04
| 44362 | 68333674 | 23 | P47 | Region_8 | 37.33 | 353.94 | 221.81 | 6 | 2024-11-15 20:25:16
| 44361 | 57175108 | 14 | P66 | Region_6 | 7.61 | 351.79 | 325.02 | 5 | 2024-11-15 20:24:03
| 44360 | 46896886 | 33 | P28 | Region_8 | 38.11 | 320.62 | 198.43 | 1 | 2024-11-15 20:27:29
| 44359 | 81445923 | 32 | P25 | Region_8 | 16.23 | 160.04 | 134.07 | 6 | 2024-11-15 20:23:48
| 44358 | 69216433 | 20 | P28 | Region_2 | 22.76 | 458.38 | 354.05 | 3 | 2024-11-15 20:23:10
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
10 rows in set (0.01 sec)
```