# SPOTIFY DATA PIPELINE
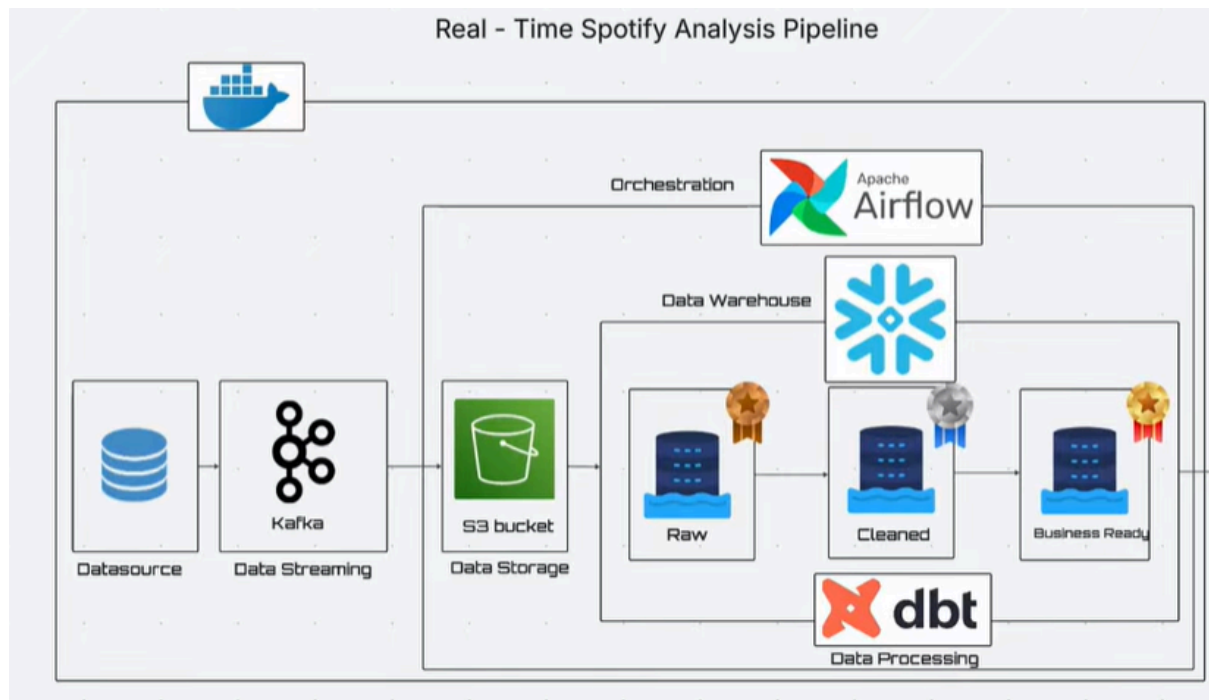
Architecture:



Tech-Stack Used:
1.Kafka
2.Snowflake
3.Airflow
4.Dbt
5.Docker
6.AWS S3

Development
1. Create a dir simulator (mkdir simulator), this will contain the producer code that will produce data in real-time.
2.Create docker dir which will have a docker-compose file having all the containers: zookeeper, kafka, kafdrop(for UI), minio for storage, postgres for initializing airflow & airflow containers. We will also create a .env file.(Make your credentials)
3. Create the requirement.txt file and install all dependencies:
pip install -r requirements.txt  # here -r is for read
4. Do docker compose up after changing dir to docker.
5. Check if kafdrop, airflow and minio are working by clicking the port in docker desktop. Check login creds as well.
6. Now let's come to producer code:

6.1 We set up Kafka config and created a Kafka Producer with those config and a value serializer that will convert the dictionary to json.

6.2 Then we created a song_artist_pair list of dict, & for each artist-song pair we generate a uuid(universally unique identifier) [uuid5: `same song_id each time`]. A **namespace** is just a fixed UUID used as a *seed* to generate other UUIDs in a consistent way.

Key properties of uuid5:

- ✅ **Deterministic** (same input → same UUID every time)

- ❌ **Not random**

- ✅ Ideal for **natural keys** → **surrogate IDs**

uuid4 is used to create "random" uuid

6.3 Then we generated random users, and created a "generate_event" function to generate random records using the above data.

7. **Kafka Consumer**: We will first connect to s3 service and check if a bucket already exists, else we create a new one.

7.1 Initialize a Kafka Consumer and store the events in s3 bucket according to batch size.

7.2 auto_offset_reset="earliest"  means it will continue from where it last stopped

7.3 kafka_group_id - It saves the offset , if we change this whole offset will start again

7.4 json.load will convert json to dict,  take a string as input and returns a dictionary as output.

7.5 json.dumps will convert dict to json, take a dictionary as input and returns a string as output.
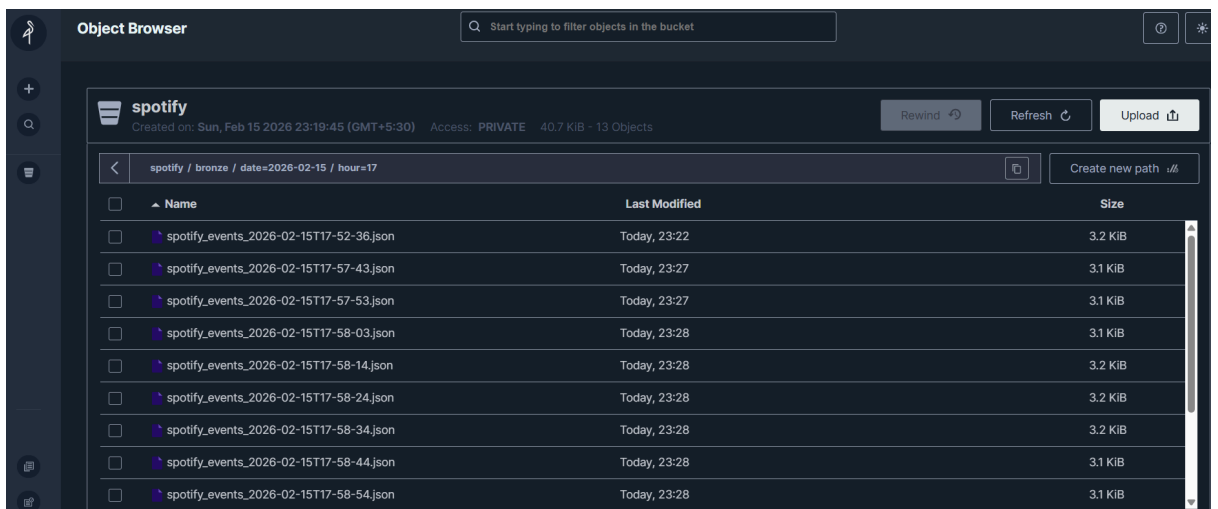


```
Produced event: pause - God's Plan by Drake (user 881dc191-47cf-4642-b4b2-9334c488cd67)
Produced event: skip - Love Story by Taylor Swift (user dc9d0fc7-9cba-480b-a0d8-093db735f4e1)
Produced event: add_to_playlist - God's Plan by Drake (user 881dc191-47cf-4642-b4b2-9334c488cd67)
Produced event: pause - Stronger by Kanye West (user 384d94c7-154e-4602-8d08-8c665b3de7b9)
Produced event: skip - Levitating by Dua Lipa (user dc9d0fc7-9cba-480b-a0d8-093db735f4e1)
Produced event: add_to_playlist - Love Story by Taylor Swift (user f78657c8-14a0-46b8-9443-04f5d3a411ff)
Produced event: pause - Stronger by Kanye West (user 881dc191-47cf-4642-b4b2-9334c488cd67)
Produced event: pause - Stronger by Kanye West (user fd69b2ea-ccd2-4d21-b918-fe873c4db470)
Produced event: skip - Levitating by Dua Lipa (user 19d48e3d-20fc-45e7-8454-04fcb810a7c1)
Produced event: skip - Love Story by Taylor Swift (user c0d8a451-a49a-41d6-b5fd-8235a6a2e18d)
Produced event: play - Stronger by Kanye West (user 71fffa9f-1423-4495-a311-74e985682796)
Produced event: play - Stronger by Kanye West (user c0d8a451-a49a-41d6-b5fd-8235a6a2e18d)
Produced event: skip - Levitating by Dua Lipa (user 19d48e3d-20fc-45e7-8454-04fcb810a7c1)

producer_v1*        0  0
```

```
Bucket spotify already exists.
Listening of events on Kafka topic 'spotify-events'.........
Uploaded 10 events to MinIO: bronze/date=2026-02-15/hour=17/spotify_events_2026-02-15T17-52-36.json
Uploaded 10 events to MinIO: bronze/date=2026-02-15/hour=17/spotify_events_2026-02-15T17-57-43.json
Uploaded 10 events to MinIO: bronze/date=2026-02-15/hour=17/spotify_events_2026-02-15T17-57-53.json
Uploaded 10 events to MinIO: bronze/date=2026-02-15/hour=17/spotify_events_2026-02-15T17-58-03.json
Uploaded 10 events to MinIO: bronze/date=2026-02-15/hour=17/spotify_events_2026-02-15T17-58-14.json
Uploaded 10 events to MinIO: bronze/date=2026-02-15/hour=17/spotify_events_2026-02-15T17-58-24.json
Uploaded 10 events to MinIO: bronze/date=2026-02-15/hour=17/spotify_events_2026-02-15T17-58-34.json
Uploaded 10 events to MinIO: bronze/date=2026-02-15/hour=17/spotify_events_2026-02-15T17-58-44.json
Uploaded 10 events to MinIO: bronze/date=2026-02-15/hour=17/spotify_events_2026-02-15T17-58-54.json
Uploaded 10 events to MinIO: bronze/date=2026-02-15/hour=17/spotify_events_2026-02-15T17-59-04.json
Uploaded 10 events to MinIO: bronze/date=2026-02-15/hour=17/spotify_events_2026-02-15T17-59-14.json
```

**Object Browser**

Q Start typing to filter objects in the bucket

**spotify**
Created on: Sun, Feb 15 2026 23:19:45 (GMT+5:30)  Access: **PRIVATE**  40.7 KiB - 13 Objects

Rewind | Refresh | Upload

spotify / bronze / date=2026-02-15 / hour=17    Create new path

| Name | Last Modified | Size |
|---|---|---|
| spotify_events_2026-02-15T17-52-36.json | Today, 23:22 | 3.2 KiB |
| spotify_events_2026-02-15T17-57-43.json | Today, 23:27 | 3.1 KiB |
| spotify_events_2026-02-15T17-57-53.json | Today, 23:27 | 3.1 KiB |
| spotify_events_2026-02-15T17-58-03.json | Today, 23:28 | 3.1 KiB |
| spotify_events_2026-02-15T17-58-14.json | Today, 23:28 | 3.2 KiB |
| spotify_events_2026-02-15T17-58-24.json | Today, 23:28 | 3.2 KiB |
| spotify_events_2026-02-15T17-58-34.json | Today, 23:28 | 3.2 KiB |
| spotify_events_2026-02-15T17-58-44.json | Today, 23:28 | 3.1 KiB |
| spotify_events_2026-02-15T17-58-54.json | Today, 23:28 | 3.1 KiB |

Now we have data in s3 bucket

8. Move data from S3 to Snowflake using Airflow:

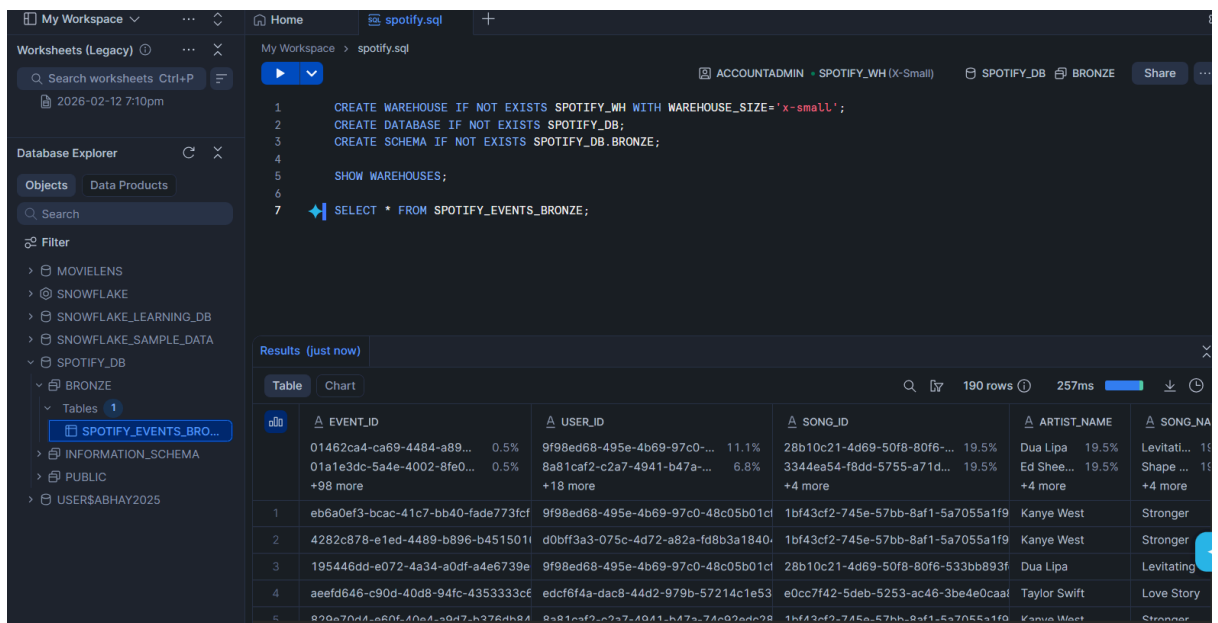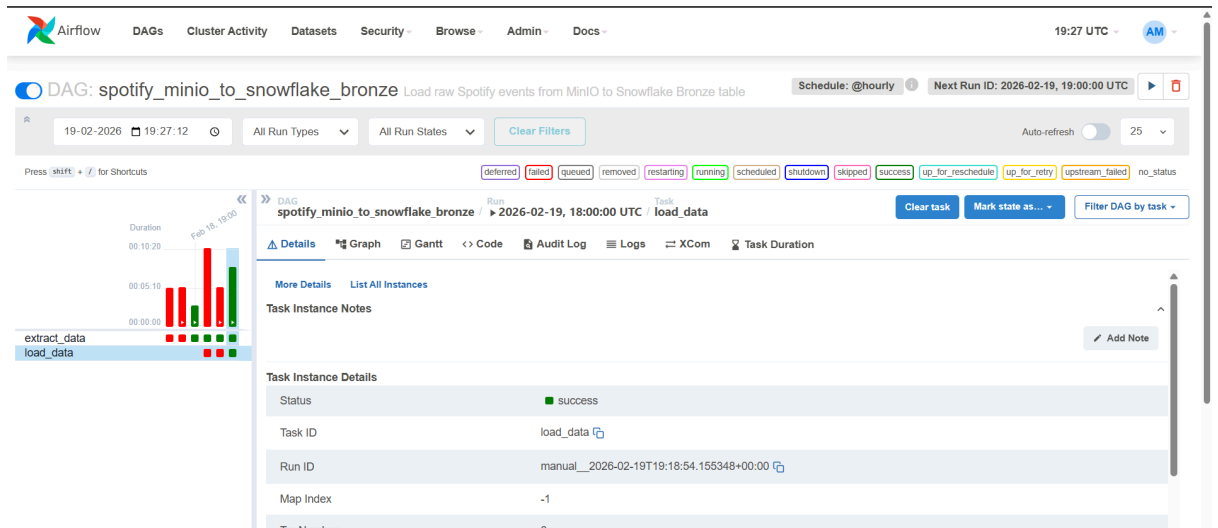8.1 Create a DAG (spotify_minio_to_snowflake_bronze)

8.2 Two tasks : extract_task to extract data from s3, and load_task to load data in snowflake bronze layer.

8.3 Extract task will extract json files data and store it into temp file on local path

8.4 load_to_snowflake will load raw data in the SF bronze layer for further processing.

**context : Airflow context dictionary that contains metadata and information about a running Directed Acyclic Graph (DAG), its tasks, and the Airflow environment.

Below is the successful run of Airflow DAG which loaded data into snowflake DB bronze layer table:

9. DBT:

    9.1 Create a new folder spotify_dbt

    9.2 Initialize dbt: dbt init and enter the details of snowflake account , warehouse, database, default schema(we used Transform schema created as default), threads, it will create a profile.

    9.3 Create sources.yml file containing metadata for silver layer source.

    9.4 Create silver dir & spotify_sql table to transform data from bronze layer and store in silver layer.

    9.5 Similarly create models for the gold layer.

    9.6 Run >> dbt run inside spotify_dbt dir

Note: Please make sure to update .env files with correct username, password and account of snowflake before running.

Thanks for reading the documentation.