🎧

# Spotify End-To-End Data Pipeline Project

## Business Use Case

A music enthusiast wants to create a **database of top global songs**, which will be automatically updated on a **daily or weekly** schedule by extracting data from the **Spotify API**. The aim is to track popular songs, artists, and playlists over time. To achieve this, a data pipeline is built to automate the process of **extracting**, **transforming**, and **visualizing** the music data using AWS services and Python code.

## Introduction

The **Spotify End-to-End Data Pipeline** project demonstrates the complete process of extracting, transforming, and loading (ETL) Spotify music data. The primary goal of this project is to create an automated pipeline that weekly collects music data, processes it, and stores it in a database for further analysis.
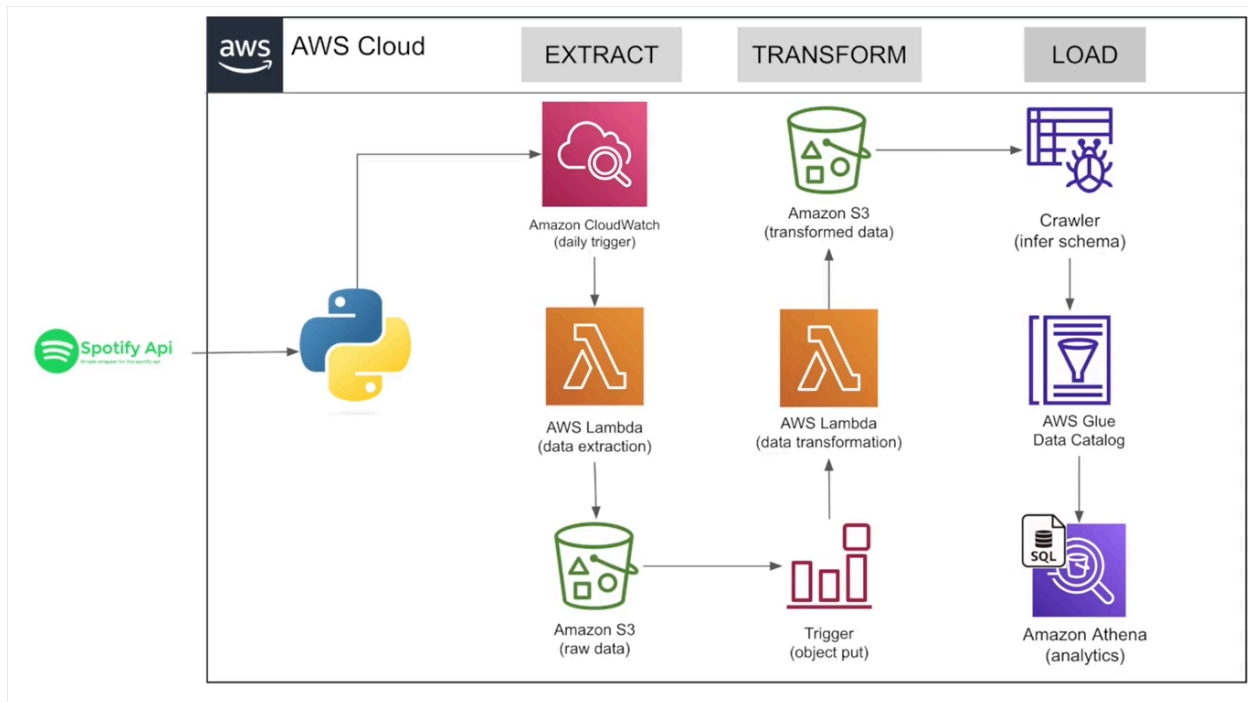
## 💡 Key AWS Services & Technologies:

- **Amazon CloudWatch**: Triggers the pipeline on a scheduled basis (daily or weekly).

- **AWS Lambda**: Runs Python code for data extraction and transformation.

- **Amazon S3**: Stores both raw and transformed data.

- **AWS Glue Crawler**: Automatically infers schema and catalogs the data for querying.
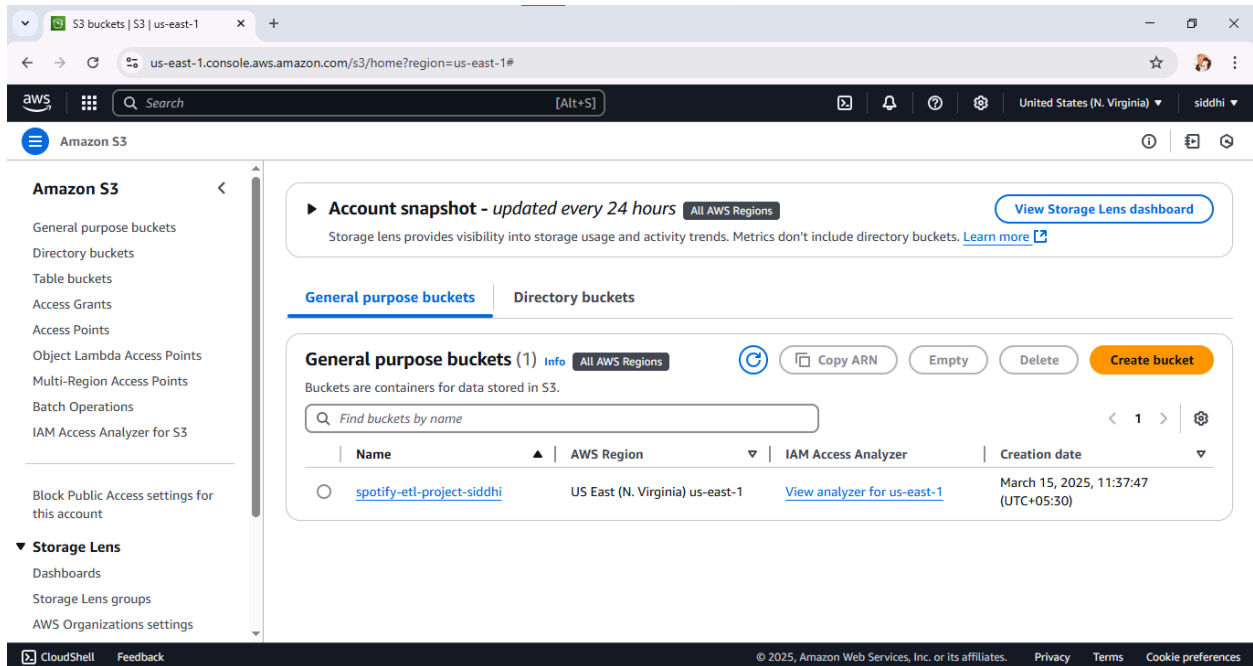
- **Amazon Athena**: Executes interactive SQL queries directly on data in Amazon S3.

- **Spotify API**: Source of music data (playlists, tracks, artists).
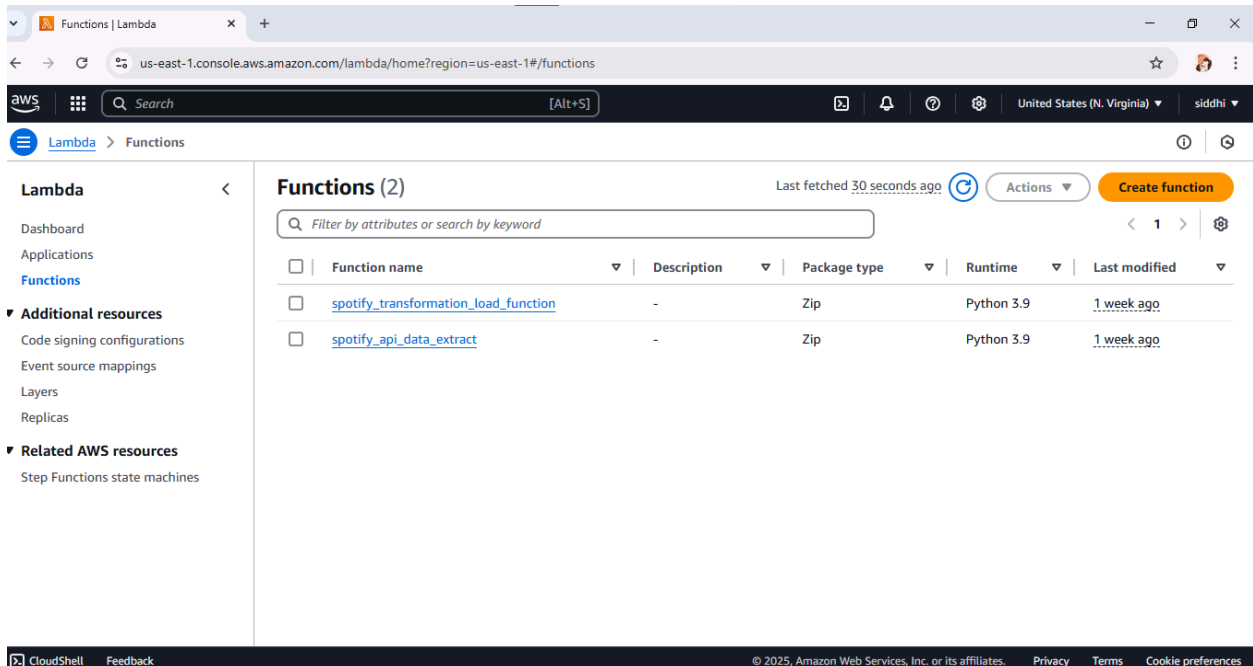
# ARCHITECTURE  DIAGRAM :



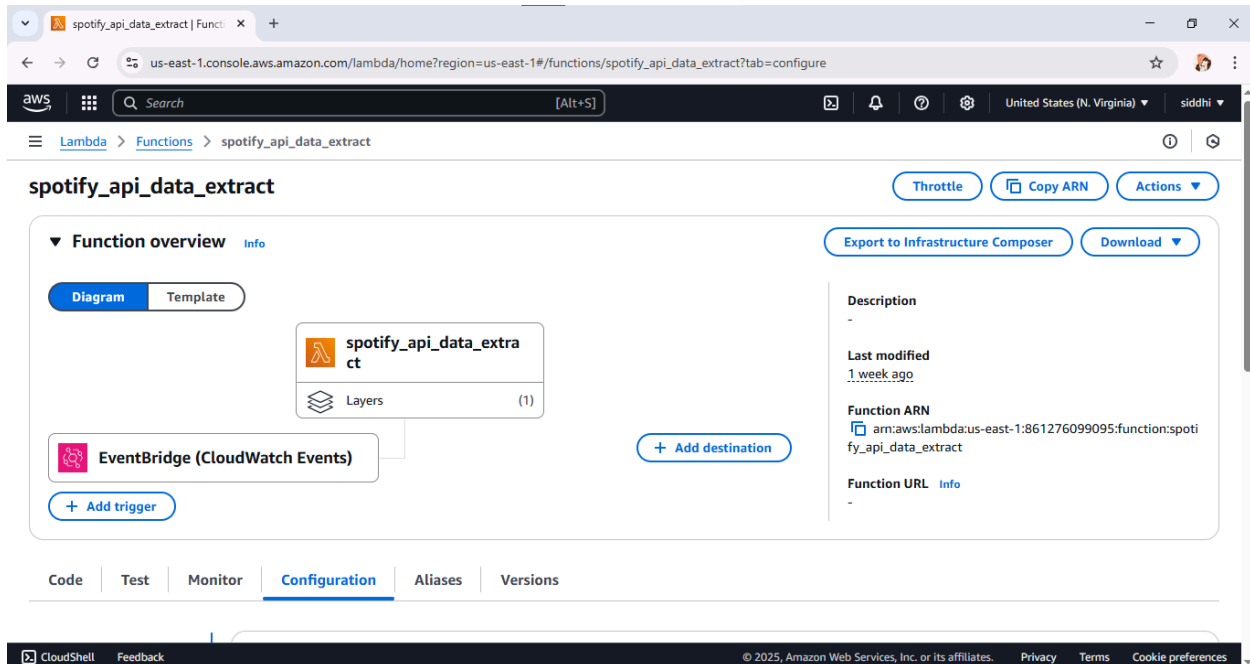# 🛠️DEMONSTRATION OF THE PROJECT AND STEPS TO BUILD   ETL PIPELINE :

1. Navigate to the **S3 service** in the AWS Console and create a bucket. I created a bucket named **spotify-etl-project-siddhi.**
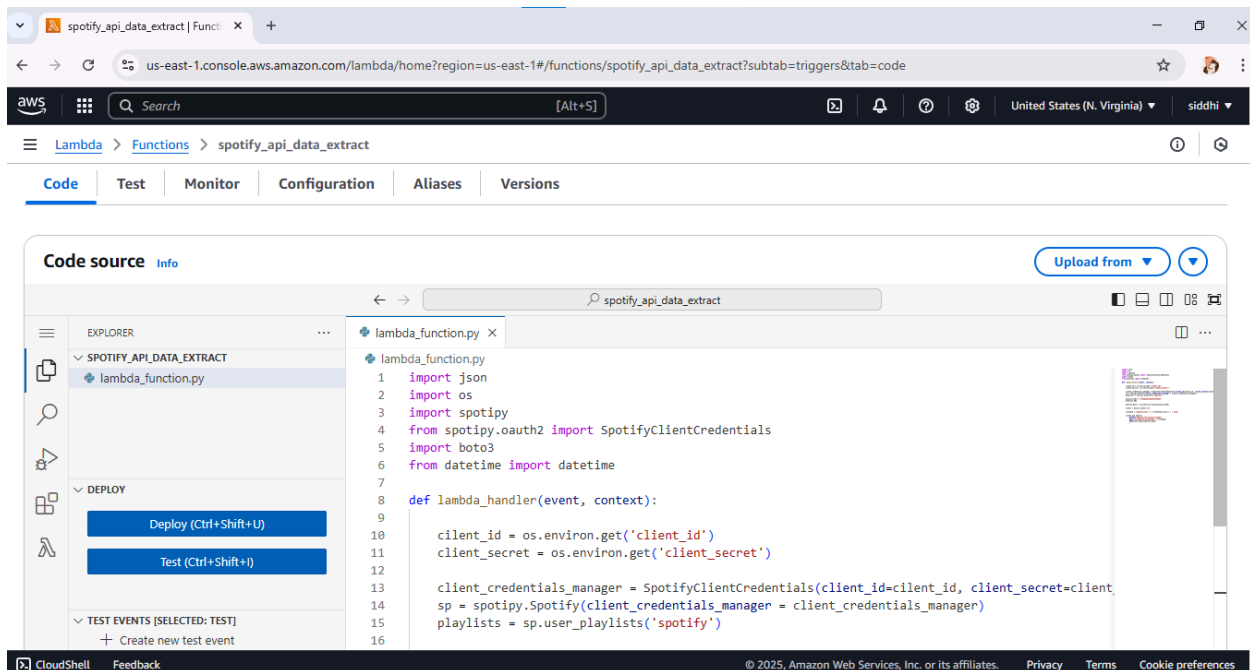
2. Then, I created two Lambda functions named **spotify_api_data_extract** and **spotify_transformation_load_function**.

## 3. I applied a
**CloudWatch alarm** to the **spotify_api_data_extract** function to **trigger it weekly** and collect **updated data**.



## Here is the function code :

4. I added an **S3 trigger** to the **spotify_transformation_load_function**, so it **automatically triggers** whenever the **spotify_api_data_extract** function runs.



**Here is the code :**

5. I created three folders in my bucket:

1. **raw_data** → Contains two subfolders:

   - **to_processed:** Stores data fetched from the API.

   - **processed:** Holds the processed data in **JSON format**.

2. The data in the **processed** folder is transformed by the **Lambda function** and stored separately in the **transformed_data** folder as:

   - **album_data.csv**

   - **artist_data.csv**

   - **songs_data.csv**

3. **athena_query_results** → Stores the **Athena query results**.

6. I created an **AWS Glue database** named **spotify-db**, which helps in:

- **Cataloging and organizing data** stored in S3.

- **Defining table schemas** for the transformed data.

- **Querying the data with Athena** by creating metadata tables.

- **Automating ETL workflows** for data transformation and analysis.

I am using
**Amazon Athena**, connected to **AWS Glue**, to run queries on the **spotify-db** database for data analysis and retrieval.

# Conclusion

The **Spotify ETL project** successfully automates the extraction, transformation, and loading of Spotify data using **AWS services**. The data is stored in **S3**, processed with **Lambda functions**, cataloged with **Glue**, and queried using **Athena**. This pipeline ensures efficient data processing and enables seamless analysis for gaining insights from Spotify's music data.