# AWS-Based Spotify End-to-End Analysis Project

## Project Overview

This project leverages various AWS services to create a data pipeline for analyzing Spotify data. Data is ingested into an S3 bucket, processed through AWS Glue, queried using AWS Athena, and visualized in Amazon QuickSight. This document outlines the steps taken to implement this solution, the technologies used, and the architecture followed.

### Technologies Used

- **Amazon S3**: For storage of raw and processed data.
- **AWS Glue**: For ETL (Extract, Transform, Load) and cataloging data.
- **AWS Athena**: For querying the processed data.
- **Amazon QuickSight**: For data visualization and dashboard creation.

---

## Step-by-Step Process

### 1. Create an IAM User

To ensure secure access to the AWS services, we first create an IAM (Identity and Access Management) user.

**Steps:**

- Log in to the AWS Management Console.
- Navigate to **IAM** and click on **Users**.
- Create a new user, providing necessary credentials.
- Assign the required roles and policies to allow access to the S3, Glue, Athena, and QuickSight services.
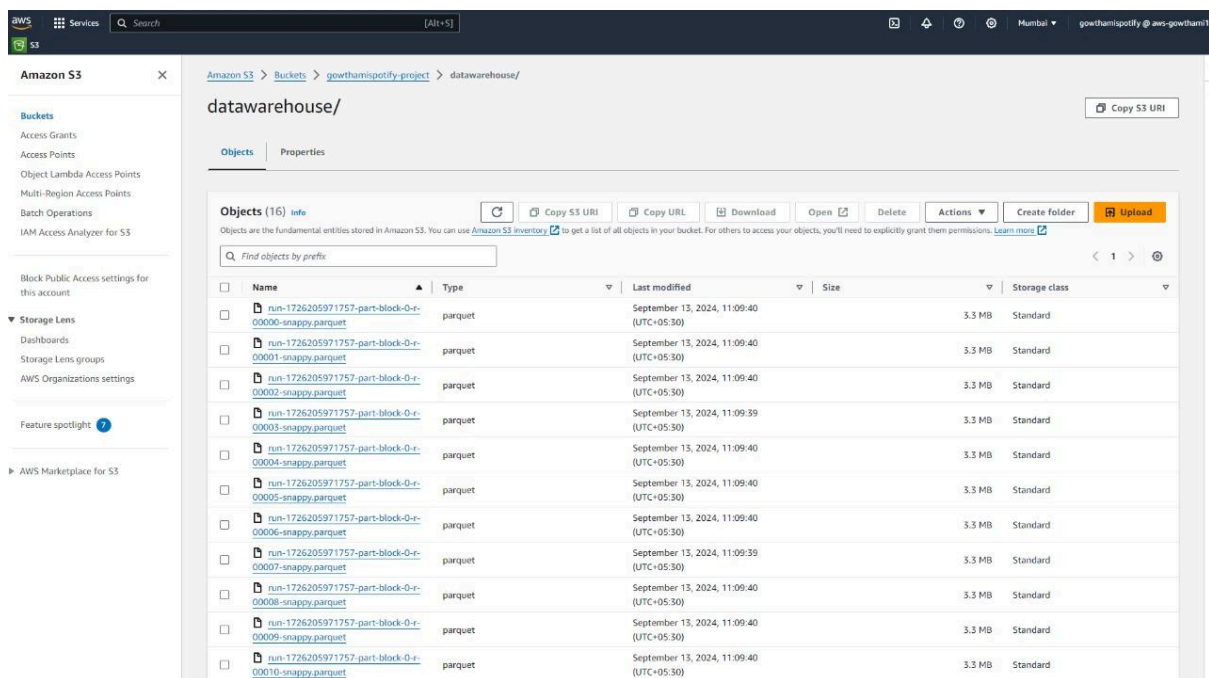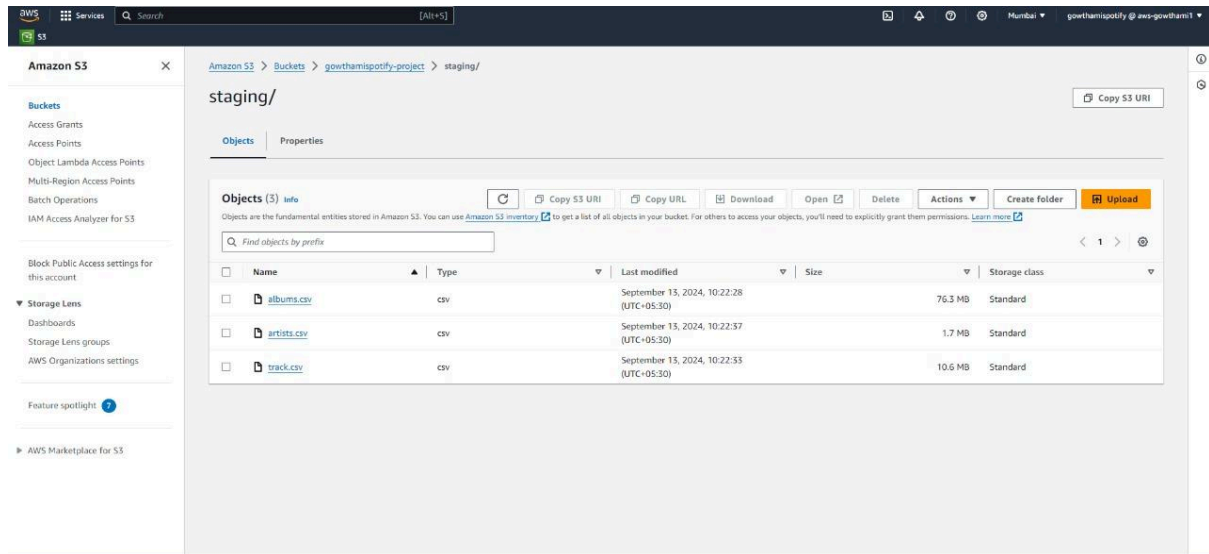
### 2. Create an S3 Bucket

S3 is used to store the raw Spotify data and processed data in different layers.

**Steps:**

- In the AWS Management Console, navigate to **S3**.
- Create a new bucket for storing data.
- Inside the bucket, create two folders (or "layers"):

- **Staging**: Where raw data will be initially uploaded.
- **Data Warehouse**: Where transformed and processed data will be stored after the ETL process.
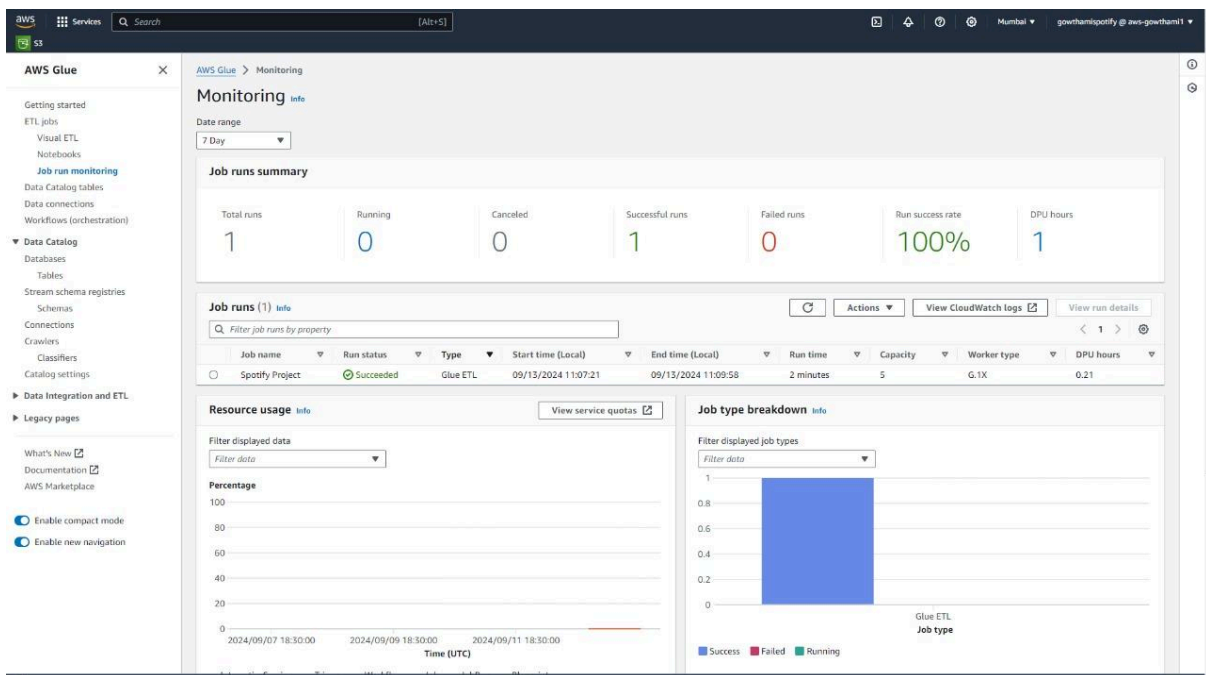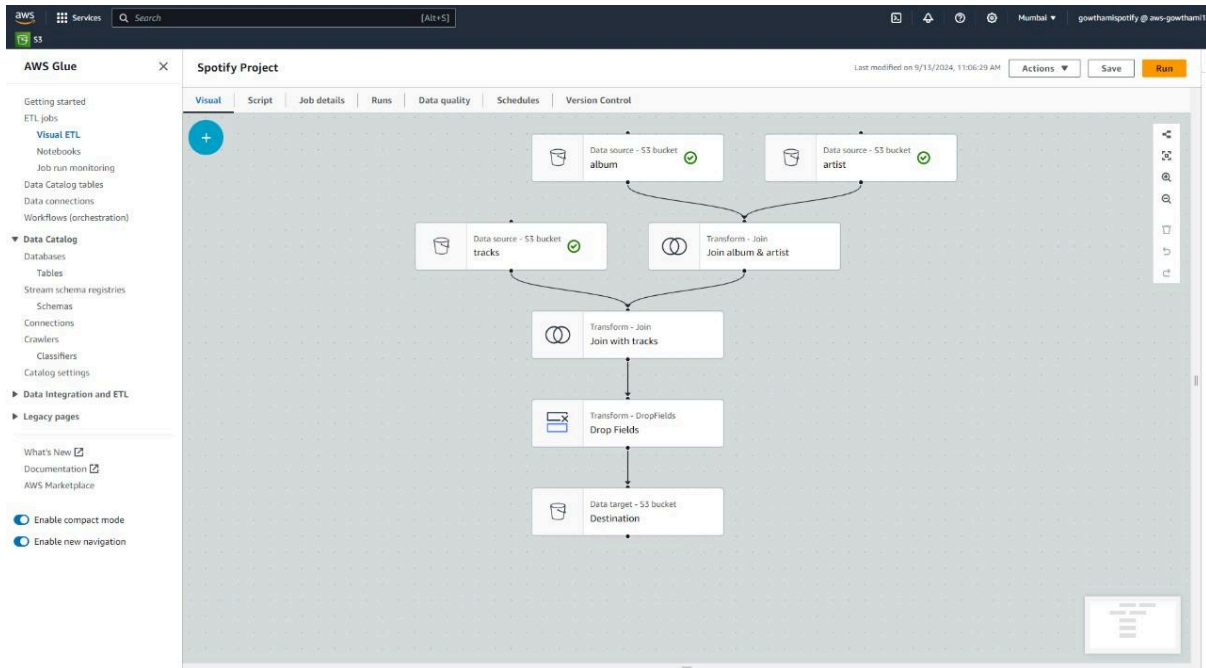




## 3. Create a Data Pipeline Using AWS Glue

Set up an ETL pipeline that transfers and transforms data from the Staging layer to the Data Warehouse layer using AWS Glue.

**Steps:**

- Go to **AWS Glue** in the console.
- Create a new **Glue Job** for the ETL process.

- Configure the job to read from the **Staging** S3 folder, apply transformations (if any), and write the processed data to the **Data Warehouse** folder in S3.





## 4. Set Up a Glue Crawler

AWS Glue Crawler is used to automatically catalog the data in the S3 Data Warehouse layer.

**Steps:**

- In the AWS Glue Console, create a new **Crawler**.

- Specify the **S3 Data Warehouse** folder as the data source.
- Run the crawler, which will create a **Data Catalog** and a **Database** in AWS Glue.







## 5. Set Up Athena and Query Editor

AWS Athena allows you to run SQL queries on the data stored in S3 using the Glue Data Catalog.

**Steps:**

- In the AWS Management Console, navigate to **Athena**.
- Set up the Query Editor by specifying the S3 bucket where query results will be stored.
- Use SQL to query the data stored in the Glue Catalog and S3 Data Warehouse.

## 6. Visualize Data in Amazon QuickSight

Amazon QuickSight is used to create visualizations and dashboards based on the processed data.

**Steps:**

- In **Amazon QuickSight**, connect to your data source by selecting Athena as the database.
- Create a new database from the Athena queries.
- Build visualizations and dashboards to display insights from the Spotify data.