# MILESTONE 5

# PROJECT 2

# GROUP – 9

# TOPIC:

# SPOTIFY MUSIC STREAMING TREND ANALYSIS

# MEMBERS:

KAVYA PATI BANDLA: Patibandla.ka@northeastern.edu

LIKHITHA GUTTAPALLI: guttapalli.l@northeastern.edu

# SPOTIFY MUSIC TREND ANALYTICS

The end-to-end ETL pipeline built for the Spotify music trend analytics is explained in this report. The pipeline processes raw data from the S3 bucket, transforms it into dimension and fact tables using Amazon glue, and then Amazon Athena is used for analytics

Dataset Information:

- Source: Kaggle
- Dataset link: Spotify Popularity Prediction-ML Practice
- Format:
- Song_data: contains detailed information regarding the audio features of the song
  - Contains 18836 rows and 10 columns of different data types
- Song_info: the information regarding the song characteristics such as the artist, album, playlist etc.
  - Contains 11836 rows and 4 columns of different data types
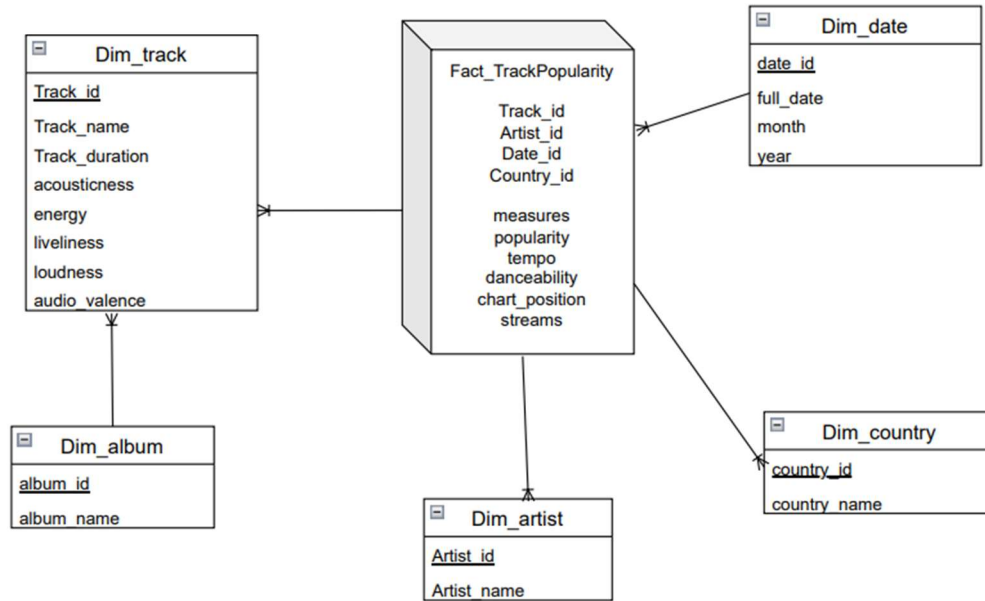- Generated the song_history.json file for details regarding the number of streams and chart positions.


Dimension Tables:

- Dim_track: Contains information regarding the track such as track_id, track name, duration, energy, danceability, liveliness..
- Dim_artist: Contains the artist information such as the artist_id and artist name
- Dim_album: Contains the information regarding the album id and album name
- Dim_date: Contains the fulldate, month and year for further analysis
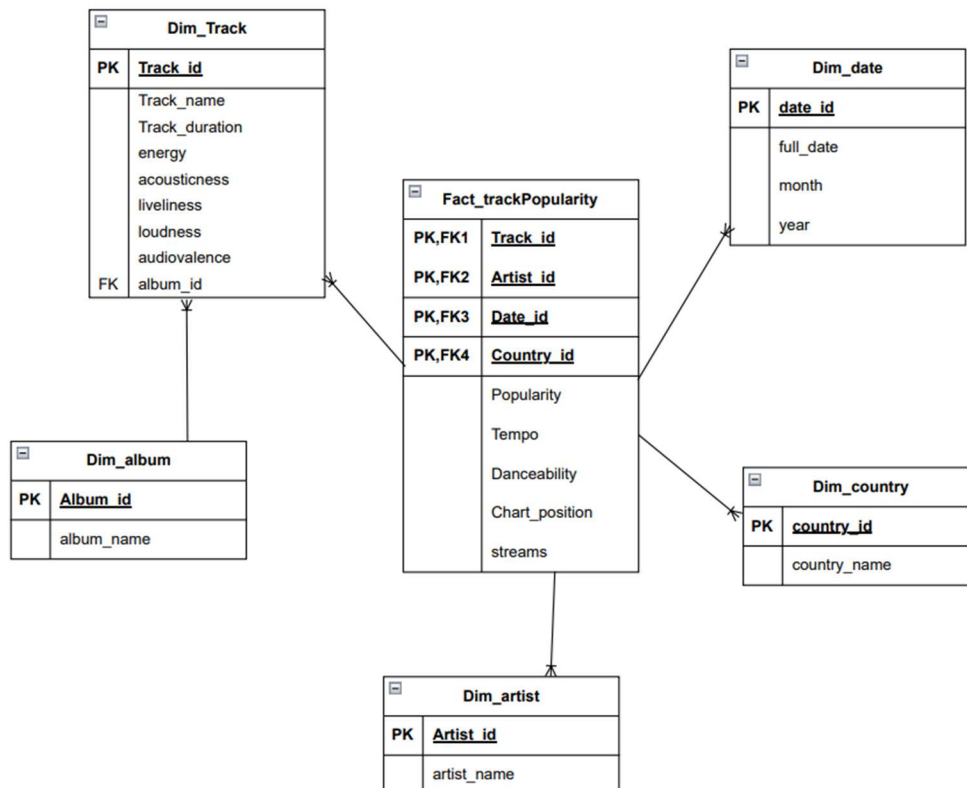- Dim_country: contains the country_id and country name


Fact Table:

- Fact_TrackPopularity: Contains the keys of the dimension table along with measures such as chart position, popularity and number of streams

# Conceptual model



**Dim_track**
- Track_id
- Track_name
- Track_duration
- acousticness
- energy
- liveliness
- loudness
- audio_valence

**Fact_TrackPopularity**
- Track_id
- Artist_id
- Date_id
- Country_id
- measures
- popularity
- tempo
- danceability
- chart_position
- streams

**Dim_date**
- date_id
- full_date
- month
- year

**Dim_album**
- album_id
- album_name

**Dim_artist**
- Artist_id
- Artist_name

**Dim_country**
- country_id
- country_name

# Logical Model



**Dim_Track**

| PK | Track_id |
|----|----------|
|  | Track_name |
|  | Track_duration |
|  | energy |
|  | acousticness |
|  | liveliness |
|  | loudness |
|  | audiovalence |
| FK | album_id |

**Dim_date**

| PK | date_id |
|----|---------|
|  | full_date |
|  | month |
|  | year |

**Fact_trackPopularity**

| PK,FK1 | Track_id |
|--------|----------|
| PK,FK2 | Artist_id |
| PK,FK3 | Date_id |
| PK,FK4 | Country_id |
|  | Popularity |
|  | Tempo |
|  | Danceability |
|  | Chart_position |
|  | streams |

**Dim_album**

| PK | Album_id |
|----|----------|
|  | album_name |

**Dim_country**

| PK | country_id |
|----|-----------|
|  | country_name |

**Dim_artist**

| PK | Artist_id |
|----|-----------|
|  | artist_name |

**ETL PIPELINE**

- **Storage layer**
    - Amazon S3 bucket:
    - Raw data path:
    - Processed data path:
- **Processing layer**
    - AWS Lambda Trigger: Activates Glue ETL jobs when all required files are present. Sequential job execution (dimension tables first, fact tables second)
    - AWS Glue ETL Jobs: Creates dimension tables from raw data and creates fact tables referencing dimension tables
- **Data Crawling**
    - Processed-data crawler: Scans the processed data files to create catalog
    - Raw-data crawler: Scans the raw data files to create catalog tables
- **Serving layer**
    - AWS Athena is used as a SQL query interface for the analytics. It enables direct querying of the processed data

**Data Extraction:**

- An S3 bucket is created for storing data.
- Manually uploaded the csv and the json files to the raw data folder created in the bucket
- A lambda function will be enabled on the bucket, once the files are uploaded it will trigger the ETL jobs

☰ Amazon S3 › Buckets ① ⊡ ⊘

## General purpose buckets  All AWS Regions        Directory buckets

### General purpose buckets (1) Info
Buckets are containers for data stored in S3.

⟳  📋 Copy ARN   Empty   Delete   **Create bucket**

🔍 Find buckets by name                                          ‹ 1 › ⚙

| | Name ▲ | AWS Region ▽ | Creation date ▽ |
|---|---|---|---|
| ○ | spotify-popularity-analytics | US East (N. Virginia) us-east-1 | November 21, 2025, 15:38:11 (UTC-05:00) |

▶ **Account snapshot** Info                              View dashboard
Updated daily
Storage Lens provides visibility into storage usage and activity trends.

▶ **External access summary - *new*** Info
Updated daily
External access findings help you identify bucket permissions that allow public access or access from other AWS accounts.

---

☰ Amazon S3 › Buckets › spotify-popularity-analytics ① ⊡ ⊘

## spotify-popularity-analytics Info

**Objects**  Metadata  Properties  Permissions  Metrics  Management  Access Points

### Objects (2)
Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory ↗ to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. Learn more ↗

⟳  📋 Copy S3 URI   📋 Copy URL   ⬇ Download   Open ↗   Delete   Actions ▾   Create folder   ⬆ **Upload**

🔍 Find objects by prefix                                          ‹ 1 › ⚙

| | Name ▲ | Type ▽ | Last modified ▽ | Size ▽ | Storage class ▽ |
|---|---|---|---|---|---|
| ☐ | 📁 processed-data/ | Folder | - | - | - |
| ☐ | 📁 raw-data/ | Folder | - | - | - |

## Data Crawling

- o The datasets in the raw data are made available in the data catalogue using the AWS Glue Crawler. Which crawls all the sub folders present in the raw data for data.

## Data Transformation

- o The AWS Glue python script is used to transform the raw data into processed data, that is the dimension tables and fact table
  - o Dimsong
  - o Dimartist
  - o Dimalbum
  - o Dimdate
  - o Dimcountry
  - o FactSongPopularity
- o Created two scripts for the data transformation. Create_dim_tables and create_fact_table for the transformation of raw data to processed data.

**AWS Glue**

- Getting started
- ETL jobs
  - Visual ETL
  - **Notebooks**
  - Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)
- Zero-ETL integrations **New**
- ▼ **Data Catalog**
  - Databases
    - Tables
  - Stream schema registries
    - Schemas
  - Connections
  - Crawlers
    - Classifiers
  - Catalog settings
- ▶ **Data Integration and ETL**
- ▶ **Legacy pages**

What's New ↗
Documentation ↗
AWS Marketplace

🔵 Enable compact mode
🔵 Enable new navigation

# AWS Glue Studio Info

## Create job Info

| | Author in a visual interface focused on data flow. | | Author using an interactive code notebook. | | Author code with a script editor. |
|---|---|---|---|---|---|
| ▦ | **Visual ETL** | ⬚ | **Notebook** | { } | **Script editor** |

▶ **Example jobs** Info      [ Create example job ]

## Your jobs (2) Info

🔄   [ Actions ▾ ]   [ Run job ]

🔍 Filter jobs by property      < 1 > ⚙

| ☐ | Job name ▽ | Type | Created by | Last modified ▽ | AWS Glue version ▽ | Action ▽ |
|---|---|---|---|---|---|---|
| ☐ | create_fact_table | Glue ETL | Script | 11/21/2025, 5:08:02 PM | 5.0 | - |
| ☐ | create_dim_tables | Glue ETL | Script | 11/21/2025, 5:04:37 PM | 5.0 | - |

---

**AWS Glue**

- Getting started
- ETL jobs
  - Visual ETL
  - **Notebooks**
  - Job run monitoring
- Data Catalog tables
- Data connections
- Workflows (orchestration)
- Zero-ETL integrations **New**
- ▼ **Data Catalog**
  - Databases
    - Tables
  - Stream schema registries
    - Schemas
  - Connections
  - Crawlers
    - Classifiers
  - Catalog settings
- ▶ **Data Integration and ETL**
- ▶ **Legacy pages**

What's New ↗
Documentation ↗
AWS Marketplace

🔵 Enable compact mode
🔵 Enable new navigation

# create_dim_tables

Last modified on 11/21/2025, 5:04:37 PM   [ Actions ▾ ]   [ Save ]   [ **Run** ]

Script    Job details    **Runs**    Data quality    Schedules    Version Control

## Job runs (1/15) Info

Last updated (UTC)
November 21, 2025 at 22:06:38   🔄   [ View details ]   [ Stop job run ]   [ ✦ Troubleshoot with AI ]    [ **Table View** | Card View ]

🔍 Filter job runs by property      < 1 > ⚙

| | Run status ▽ | Retries ▽ | Start time (Local) ▽ | End time (Local) ▽ | Duration ▽ | Capacity (DPUs) ▽ | Worker type ▽ | Glue version ▽ |
|---|---|---|---|---|---|---|---|---|
| 🔘 | ✅ Succeeded | 0 | 11/21/2025 17:04:40 | 11/21/2025 17:06:09 | 1 m 13 s | 10 DPUs | G.1X | 5.0 |
| ⭕ | ❌ Failed | 0 | 11/21/2025 17:02:27 | 11/21/2025 17:03:54 | 1 m 3 s | 10 DPUs | G.1X | 5.0 |
| ⭕ | ❌ Failed | 0 | 11/21/2025 16:59:55 | 11/21/2025 17:01:09 | 1 m 7 s | 10 DPUs | G.1X | 5.0 |
| ⭕ | ✅ Succeeded | 0 | 11/21/2025 16:54:27 | 11/21/2025 16:55:57 | 1 m 22 s | 10 DPUs | G.1X | 5.0 |

**Run details**    Input arguments (9)    Logs    Run insights    Metrics    Troubleshooting analysis - *preview*    Spark UI

| | | | |
|---|---|---|---|
| Job name | Start time (Local) | Glue version | Last modified on (Local) |
| create_dim_tables | 11/21/2025 17:04:40 | 5.0 | 11/21/2025 17:06:09 |
| Id | End time (Local) | Worker type | Log group name |
| jr_c2ae95f290eac899eb16949d2069f35488c6956e91099b 0413e7c79385b77c42 📋 | 11/21/2025 17:06:09 | G.1X | /aws-glue/jobs |
| Run status | Start-up time | Max capacity | Number of workers |
| ✅ Succeeded | 15 seconds | 10 DPUs | 10 |
| Retry attempt number | Execution time | Execution class | Timeout |
| Initial run | 1 minute 13 seconds | Standard | 480 minutes |
| Trigger name | Security configuration | Cloudwatch logs | Usage profile |
| - | - | • Output logs ↗ | - |
| | | • Error logs ↗ | |

## Data Catalog

- A glue crawler is used to make the data in the processed S3 bucket available in the Data Catalog by scanning all the subfolders within the specified bucket.

AWS Glue > Crawlers > processeddata

**AWS Glue** ‹

Getting started
ETL jobs
    Visual ETL
    Notebooks
    Job run monitoring
Data Catalog tables
Data connections
Workflows (orchestration)
Zero-ETL integrations  New
▼ Data Catalog
Databases
    Tables
Stream schema registries
    Schemas
Connections
Crawlers
    Classifiers
Catalog settings
▶ Data Integration and ETL
▶ Legacy pages

What's New ↗
Documentation ↗
AWS Marketplace

🔵 Enable compact mode
🔵 Enable new navigation

✅ **Crawler successfully starting**
The following crawler is now starting: "processeddata"                                      ✕

# processeddata
Last updated (UTC)  November 21, 2025 at 23:08:17  ↻  | Run crawler | Edit | Delete |

## Crawler properties

| | | | |
|---|---|---|---|
| **Name** processeddata | **IAM role** LabRole ↗ | **Database** spotify-popularity | **State** READY |
| **Description** - | **Security configuration** - | **Lake Formation configuration** - | **Table prefix** - |
| **Maximum table threshold** - | | | |

▶ Advanced settings

| **Crawler runs** | Schedule | Data sources | Classifiers | Tags |

### Crawler runs (1)
The list of crawler runs for this crawler.

| 🔍 Filter data | | 📅 Filter by a date and time range |

↻ | Stop run | View CloudWatch logs ↗ | View run details |

‹ 1 › ⚙

| | Start time (UTC) ▲ | End time (UTC) ▽ | Current/last duration ▽ | Status ▽ | DPU hours ▽ | Table changes ▽ |
|---|---|---|---|---|---|---|
| ⚪ | November 21, 2025 at 23:08:43 | November 21, 2025 at 23:09:45 | 01 min 02 s | ✅ Completed | - | - |

CloudShell  Feedback                                                © 2025, Amazon Web Services, Inc. or its affiliates.   Privacy   Terms   Cookie preferences

🔵 43°F
Mostly cloudy
6:09 PM
11/21/2025

---

AWS Glue > Databases

**AWS Glue** ‹

Getting started
ETL jobs
    Visual ETL
    Notebooks
    Job run monitoring
Data Catalog tables
Data connections
Workflows (orchestration)
Zero-ETL integrations  New
▼ Data Catalog
Databases
    Tables
Stream schema registries
    Schemas
Connections
Crawlers
    Classifiers
Catalog settings
▶ Data Integration and ETL
▶ Legacy pages

What's New ↗
Documentation ↗
AWS Marketplace

🔵 Enable compact mode
🔵 Enable new navigation

# Databases (2)
Last updated (UTC)  November 21, 2025 at 23:11:07  ↻  | Edit | Delete | **Add database** |

A database is a set of associated table definitions, organized into a logical group.

| 🔍 Filter databases |

‹ 1 › ⚙

| | Name ▲ | Description ▽ | Location URI ▽ | Created on (UTC) ▽ |
|---|---|---|---|---|
| ☐ | spotify | - | - | November 21, 2025 at 20:48:00 |
| ☐ | spotify-popularity | - | - | November 21, 2025 at 23:07:25 |

CloudShell  Feedback                                                © 2025, Amazon Web Services, Inc. or its affiliates.   Privacy   Terms   Cookie preferences

🔵 43°F
Mostly cloudy
6:11 PM
11/21/2025

## Data Loading:

o The target tables are then loaded into Athena for querying and for popularity analytics of the Spotify music