

# Songs Recommendation System

The "Song Recommendation System " project aims to revolutionize music streaming experiences by harnessing the power of Spotify's API and AWS services. With the exponential growth of digital music consumption, personalized recommendations have become essential for engaging users and enhancing their enjoyment. By leveraging real-time data from Spotify, the system dynamically generates tailored song suggestions based on user preferences and trends. Through a combination of data collection, feature extraction, ETL processing, and recommendation generation, the system provides an intuitive and seamless user experience. With a focus on scalability, data integrity, and automation, this project sets out to redefine how users discover and enjoy music in the digital age.

## Data Sources:

### Live streaming Data:

Spotify API Data: to create a dynamic recommendation engine by continuously gathering real-time data from Spotify's API. Leveraging live streaming data enables the system to provide up-to-date and personalized song recommendations to users, enhancing their music streaming experience.

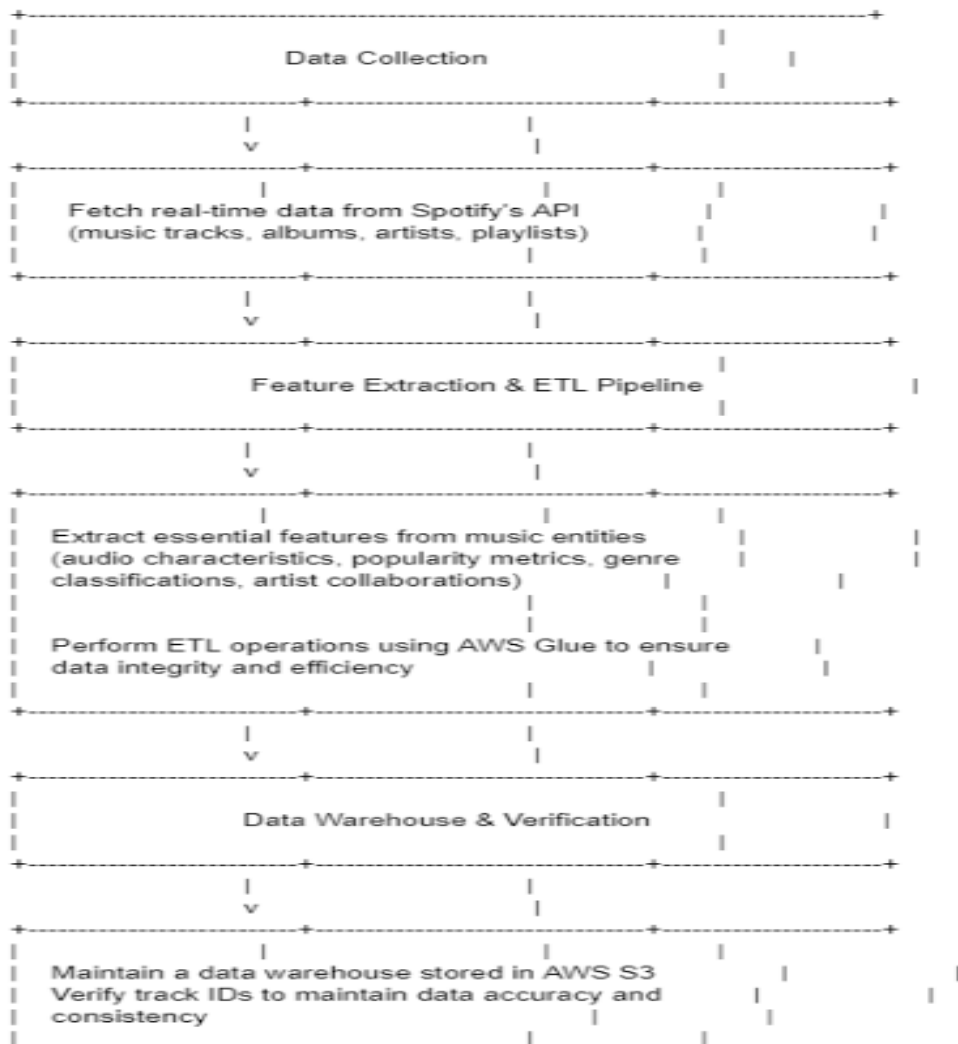
(<https://developer.spotify.com/documentation/web-api> ),

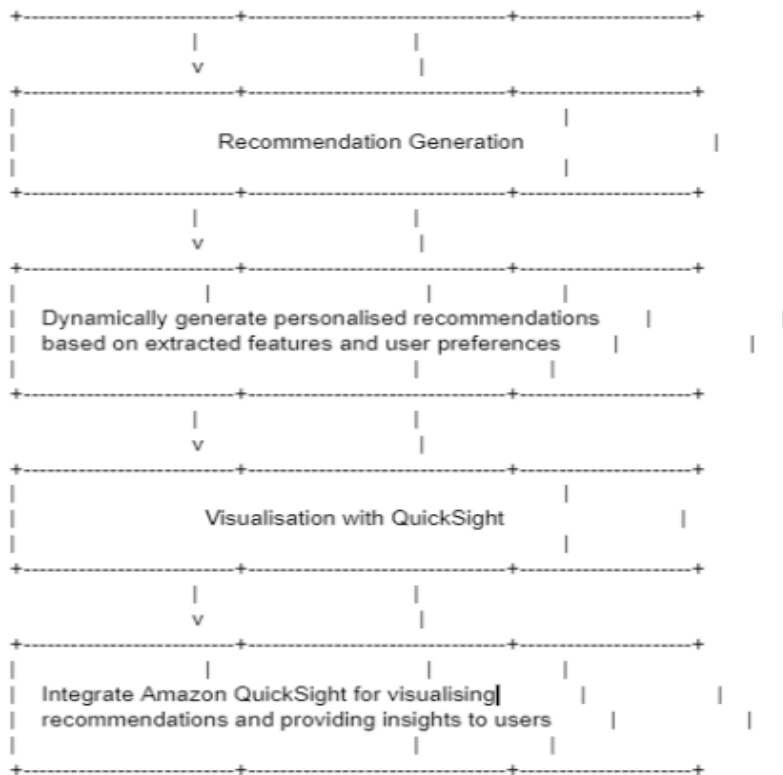
Dataset: Top Spotify Songs  
2023](<https://www.kaggle.com/datasets/nelgiriyeewithana/top-spotify-songs-2023>)

## The tools used:

- Spotify API: Real-time data on music tracks, albums, artists, and playlists. Primary source for gathering song and user preference information.
- AWS S3 Buckets: Store collected data, ensuring scalability, accessibility, and efficient management.
- AWS Glue: Orchestrate ETL operations, ensuring seamless data processing and transformation.
- Amazon Athena: Query and analyse structured data directly in S3, enabling efficient data retrieval and manipulation.
- Amazon QuickSight: Visualise recommendations and insights, providing an intuitive interface for exploration and discovery.
- Python Libraries:
  - Spotipy: Interact with Spotify's API to fetch real-time data.

Pandas: Manipulate and analyse data retrieved from Spotify's API.





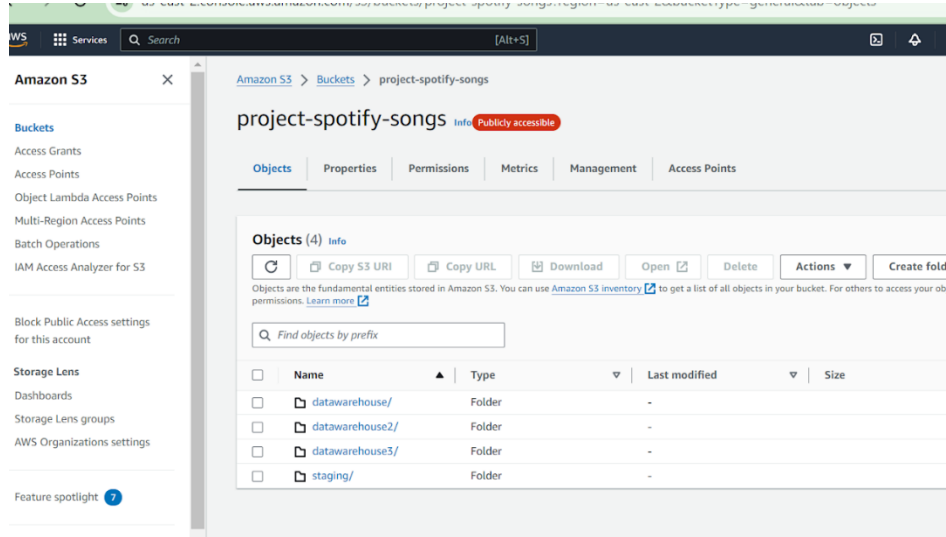
## IMPLEMENTATION

### Spotify Api:

- We head over to the Spotify Developer Dashboard and sign in or create an account.
- Once logged in, we create a new application by giving it a name (Song Recommender) and brief description.
- After creating our app, we'll get a Client ID and Client Secret, which are like special keys which we used to connect our app to Spotify.

### AWS S3 Buckets:

1. Bucket Information:
  - Bucket Name: project-spotify-songs
  - Access: Not publicly accessible
2. Objects in the Bucket:
  - There are the objects listed within the bucket:
    - data warehouse/ - for storing the output recommendations.
    - staging/ - for storing the input datasets for giving the output.



## AWS Glue:

### 1. Data Sources:

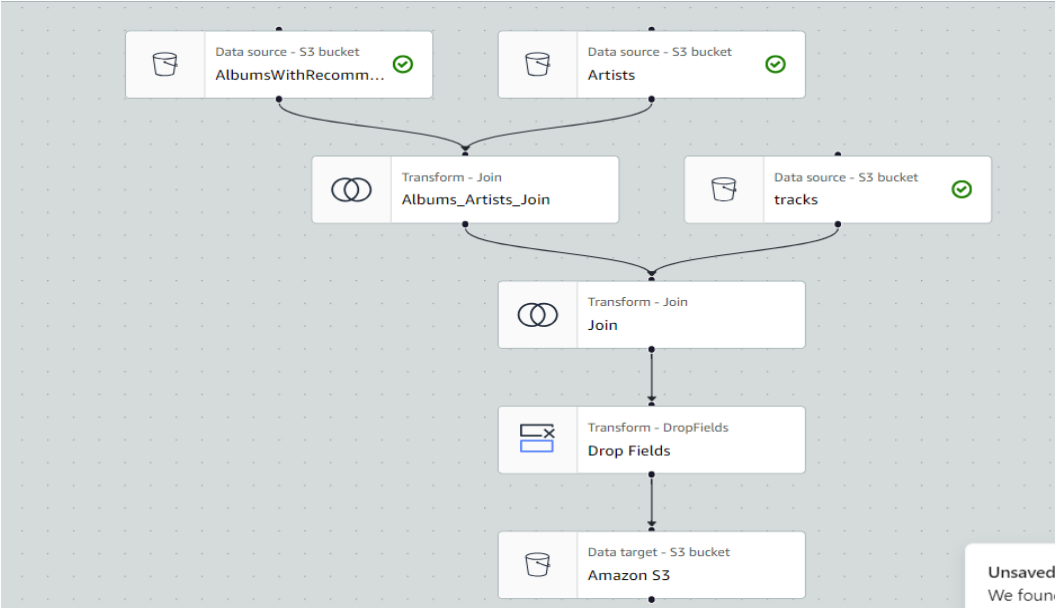
- There are three data sources represented as S3 buckets:
  - Albums With Recommendations
  - Artists
  - tracks

### 2. Transformation Steps:

- Two join transformations are performed:
  - Albums\_Artists\_Join: Combines data from the “Albums With Recommendations” and “Artists” sources.
  - The output of the first join is then joined with the “tracks” data source.
- Following the joins, there is a transformation step labelled “Drop Fields,” which likely removes unnecessary fields or columns.

### 3. Data Target:

- The final output is stored in another S3 bucket named “Amazon S3.”



**Crawlers:**

we employed crawlers to automatically discover and catalogue the schema of the data stored in our S3 buckets.

Table details		Advanced properties	
Name datawarehouse	Description -	Database spotify	Classification Parquet
Location s3://project-spotify-songs/datawarehouse/	Connection -	Deprecated -	Last updated April 21, 2024 at 06:58:07
Input format org.apache.hadoop.hive.q1.io.parquet.MapredParquetInputFormat	Output format org.apache.hadoop.hive.q1.io.parquet.MapredParquetOutputFormat	Serde serialization lib org.apache.hadoop.hive.q1.io.parquet.serde.ParquetHiveSerDe	

Schema		Partitions	Indexes	Column statistics - new
Schema (13) View and manage the table schema.				
<input type="text" value="Filter schemas"/>				
#	Column name	Data type	Partition key	Comment
1	followers	string	-	-
2	track_id	string	-	-
3	artist_popularity	string	-	-
4	.track_name	string	-	-
5	artist_id	string	-	-

**Athena:**

We leveraged Amazon Athena for querying and analyzing this structured data directly in S3.

**Data**

Data source: AwsDataCatalog  
Database: spotify

Tables and views: **Tables (4)**  
 - datawarehouse  
 - datawarehouse2  
 - datawarehouse3  
 - songs

**Query 1**

```
1 SELECT * FROM datawarehouse
```

Run again | Explain | Cancel | Clear | Create

Query results | Query stats

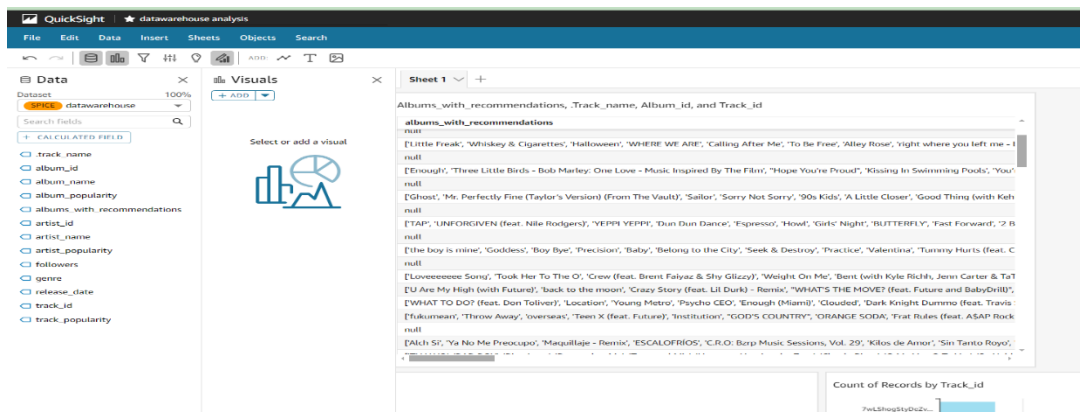
Completed | Time in queue: 57 ms | Run time: 548 ms | Data scanned: 92.31 KB

Results (238)

#	followers	track_id	artist_popularity	track_name	artist_id
1	3773086	7jrQIBanUJOM56RyC46jZ	77	Another Love	2txHhyCwHjUEpJjWfEyyqX
2	82063955	7sVbKoBdhXtYCE06qC15N	95	We Still Don't Trust You	1Xyo4u8uXC1ZmMpatF05PJ
3	82063955	7sVbKoBdhXtYCE06qC15N	95	We Still Don't Trust You	1Xyo4u8uXC1ZmMpatF05PJ
4	82063955	7sVbKoBdhXtYCE06qC15N	95	We Still Don't Trust You	1Xyo4u8uXC1ZmMpatF05PJ
5	7745826	7sVbKoBdhXtYCE06qC15N	94	We Still Don't Trust You	0lEtisbK0Kxa5IF7G42Z0p
6	7745826	7sVbKoBdhXtYCE06qC15N	94	We Still Don't Trust You	0lEtisbK0Kxa5IF7G42Z0p

## QuickSight:

This enabled us to transform the insights gained from Athena into interactive visualisations and dashboards.



In conclusion, our Song Recommendation System with Spotify API project has successfully demonstrated the power of leveraging real-time music data to enhance the user experience in music streaming platforms. Through meticulous planning and execution, we were able to design and implement a robust infrastructure utilising AWS services such as S3 buckets, Glue, Athena, and QuickSight, along with Python libraries like Spotipy and Pandas.

