

# **SPOTIFY TRACKS**

Full-stack Business Intelligence solution

---

# **SPOTIFY TRACKS BUSINESS INTELLIGENCE SOLUTION**

---

Kyrellos El Naggar  
23/7/2025

# ***Table of Contents***

## ***Contents***

Table of Contents .....	2
1. Data Overview .....	3
◆ Identification & Classification Attributes .....	3
◆ Popularity & Engagement Metrics .....	3
◆ Audio Feature Metrics .....	3
◆ Musical Theory Attributes .....	4
◆ Duration Attribute .....	4
2. Conceptual & Logical & Physical Model .....	5
2.1 Conceptual Model .....	5
2.2 Logical Model .....	6
2.3 Physical Model .....	7
3. SSIS Overview .....	8
4. SSAS Overview .....	11
5. Power BI Overview .....	12

# 1. Data Overview

The project is based on a comprehensive Spotify Tracks Dataset comprising 232,725 tracks distributed evenly across 26 music genres (~10,000 tracks per genre). The dataset is used to build a full-stack Business Intelligence solution using SSIS for ETL, SSAS Tabular for analytics, and Power BI for interactive reporting.

## Data Source

- **Input Format:** Excel file
- **Attributes:**
  - track\_id, track\_name, artist\_name, genre
  - **Audio features:** popularity, danceability, energy, key, mode, speechiness, acousticness, instrumentalness, liveness, valence, tempo, duration\_ms, time\_signature

### ◆ Identification & Classification Attributes

Attribute	Description
track_id	A unique identifier assigned to each track by Spotify (e.g., 4uLU6hMCjMI75M1A2tKUQC). Used as the primary key.
track_name	The title of the track. May be multilingual or contain special characters.
artist_name	Name of the performing artist or band associated with the track.
genre	High-level musical genre classification (e.g., Jazz, Rock, Ska). Used for grouping, filtering, and summarization.

### ◆ Popularity & Engagement Metrics

Attribute	Description
popularity	A score from <b>0 to 100</b> provided by Spotify based on play counts, recency, and user interactions. Higher values indicate higher popularity.

### ◆ Audio Feature Metrics (**All values between 0.0 and 1.0 unless specified otherwise**)

Attribute	Description
danceability	Describes how suitable a track is for dancing based on tempo, rhythm stability, and beat strength. Higher means more danceable.
energy	Reflects the intensity and activity level of a track (e.g., fast, loud, and noisy tracks have higher energy).
speechiness	Measures the presence of spoken words in a track. Higher values indicate more spoken content (e.g., podcasts, rap).
acousticness	Confidence score that a track is acoustic (i.e., not synthesized). A value of <b>1.0</b> means it is likely fully acoustic.
instrumentalness	Predicts whether a track contains no vocals. Higher scores indicate less vocal content (approaching 1.0 = purely instrumental).
liveness	Detects the likelihood of an audience presence. High values (closer to 1.0) suggest the track is live-recorded.

Attribute	Description
valence	Describes the musical "positiveness" or emotional tone of a track. A score close to 1.0 indicates a happy, cheerful track; closer to 0.0 indicates a sad or melancholic track.

---

## ◆ **Musical Theory Attributes**

Attribute	Description
key	The musical key of the track (e.g., C, D#, F). A discrete value representing pitch class (12 total).
mode	Indicates the modality: <b>Major</b> (1) or <b>Minor</b> (0), reflecting the general harmonic characteristics of the track.
tempo	Estimated tempo in <b>beats per minute (BPM)</b> . Used for rhythm analysis and beat matching.
time_signature	Notation that specifies how many beats are in each bar (e.g., 4/4, 3/4). Represented as an integer from <b>0 to 5</b> , corresponding to standard musical signatures (e.g., 0 = unknown, 4 = 4/4).

---

## ◆ **Duration Attribute**

Attribute	Description
duration_ms	Length of the track in <b>milliseconds</b> . Used to calculate playback time. Converted to duration_min during ETL for easier reporting and filtering in dashboards.

## Data Quality Handling

- key, mode, time\_signature corrected for type issues.
- Null values handled across key fields (genre, artist\_name, track\_name, etc.).
- duration\_ms converted to minutes for better interpretability.

## **2. Conceptual & Logical & Physical Model**

This solution implements a **star schema** design in the Data Warehouse layer for high performance and ease of analysis.

### **2.1 Conceptual Model**

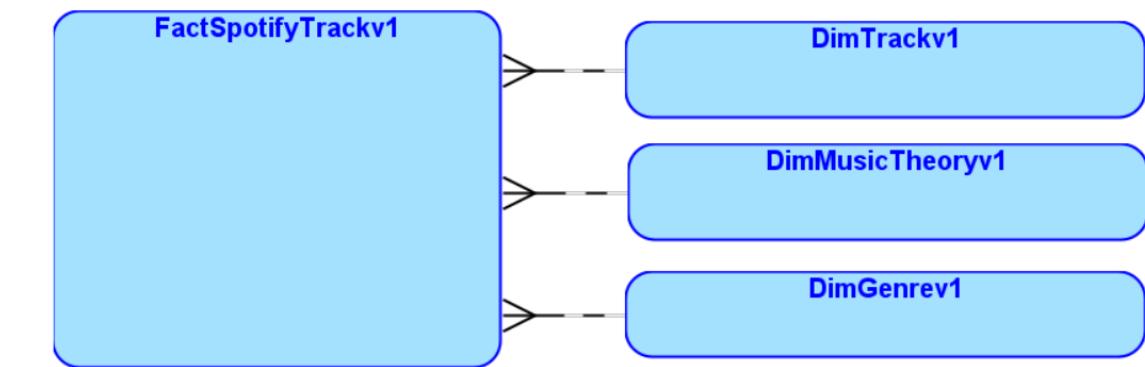
The conceptual model defines the **core entities**, their **roles**, and **relationships** within the business context—abstracted away from technical implementation. It describes **what data exists (Tables)** and **how it's related**, without getting into physical storage or data types.

The conceptual model consists of **one central fact entity** and **three surrounding dimension entities**, forming a **classic star schema** structure.

#### **Entities (Blocks):**

- FactSpotifyTrack (*Fact Table*)
- DimTrack (*Dimension*)
- DimGenre (*Dimension*)
- DimMusicalTheory (*Dimension*)

All dimensions related to fact in 1:M relationship.



## 2.2 Logical Model

The **logical model** is defined using the following entities:

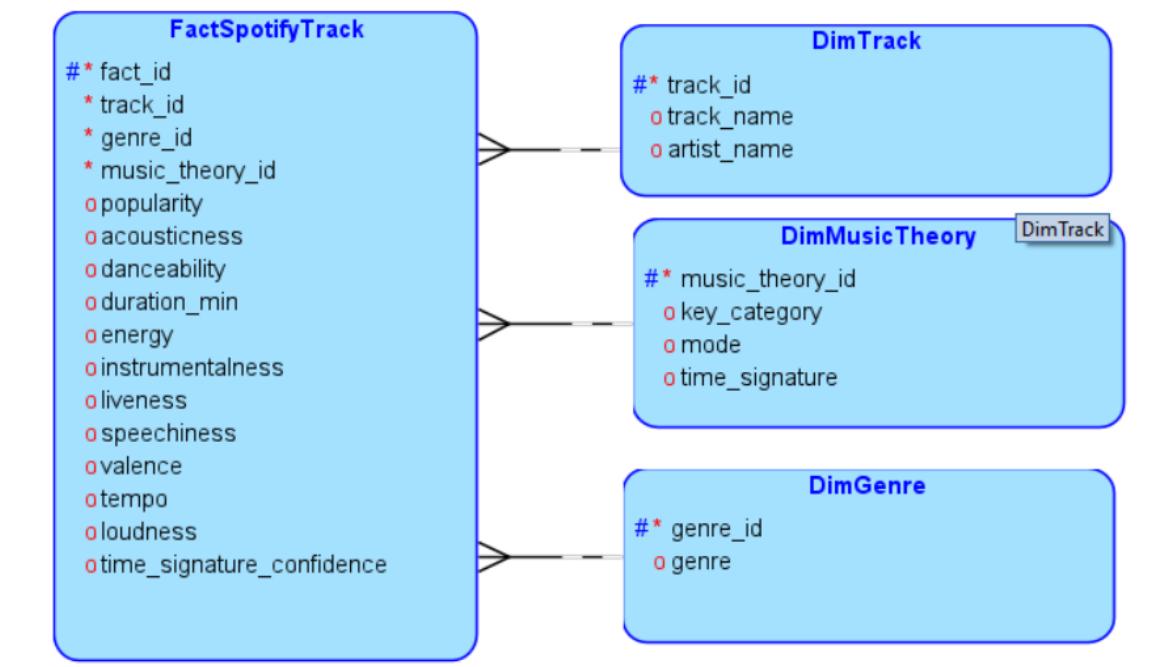
### Dimensions:

- DimTrack: Track metadata including track\_id, track\_name, artist\_name
- DimGenre: Genre classifications
- DimMusicalTheory: Musical attributes like key, mode, time\_signature

### Fact Table:

- FactSpotifyTrack: Fact table referencing all dimensions and housing measurable audio features and calculated durations.

Each dimension is designed to support one-to-many relationships with the fact table.



## 2.3 Physical Model

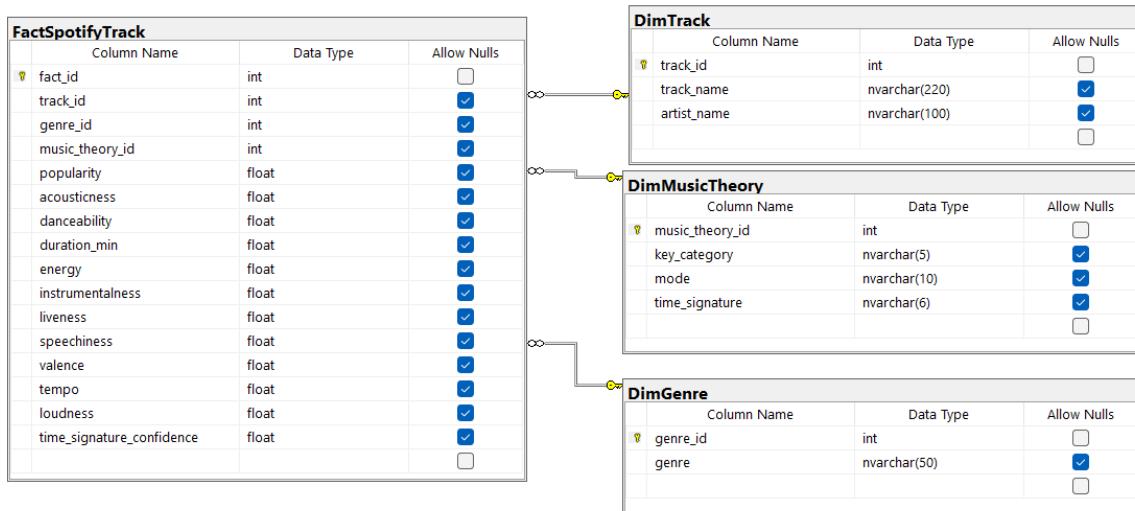
The physical model implements the **logical star schema** within a SQL Server-based Data Warehouse. It focuses on **how data is stored** in the database, **table structures**, and **physical organization**—not on how data gets there.

### Databases Involved:

- **DWH\_Spotify** → Final Data Warehouse containing structured data for analytics
- **Tables Stored Physically:**
  - **FactSpotifyTrack** (*Fact table storing measurable metrics*)
  - **DimTrack** (*Descriptive details of tracks and artists*)
  - **DimGenre** (*Genre classifications*)
  - **DimMusicalTheory** (*Musical key, mode, time signature*)

### Structure:

- Surrogate keys used for all dimension tables
- Foreign key relationships enforced between fact and dimension tables
- Star schema optimized for analytical querying and SSAS consumption
- All text fields normalized and typecast to avoid BLOB issues in reporting tools



### 3. SSIS Overview

The SSIS (SQL Server Integration Services) layer orchestrates the **ETL (Extract, Transform, Load)** process in three clearly defined stages, aligned with the architecture's layered design:

#### ETL Database Layers:

1. **ODS\_Spotify** – Raw operational data store
2. **STG\_Spotify** – Cleaned and conformed staging layer
3. **DWH\_Spotify** – Final dimensional warehouse

#### A. ODS Package – Raw Data Load:

**Objective:** Ingest raw Spotify Excel data and apply basic cleaning and conversions before loading it into the ODS layer.

##### ETL Flow:

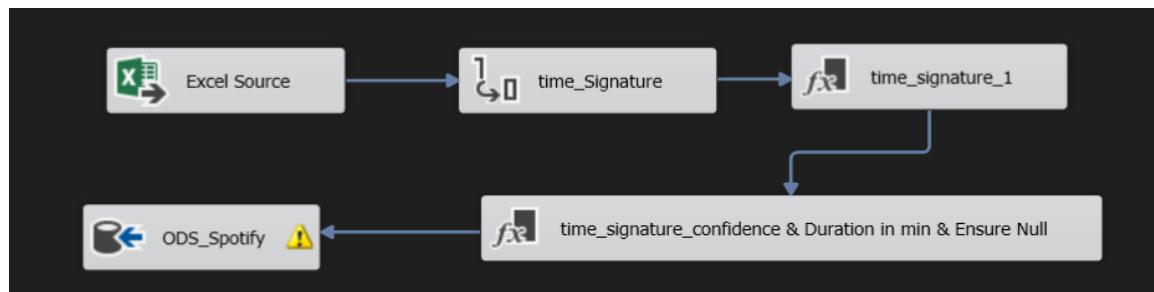
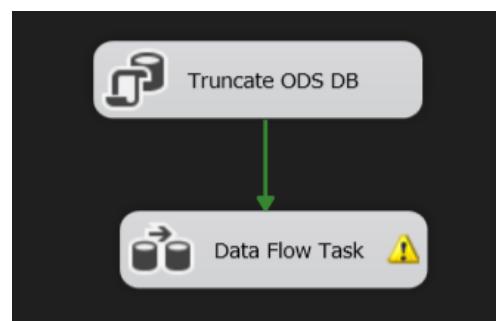
- **Source:** Excel file with 232,725 tracks
- **Destination:** ODS\_Spotify database

##### Key Tasks:

- Truncate existing data (Execute SQL Task)
- Convert time\_signature from date to string (DT\_WSTR)
- Replace NULL values in key descriptive fields:
  - track\_name, artist\_name, genre, key, mode, time\_signature
- Derive duration\_min from duration\_ms for usability in analytics

##### SSIS Components Used:

- Excel Source
- Data Conversion
- Derived Column
- Data Flow Task
- OLE DB Destination



## B. STG Package – Staging and Transformation

**Objective:** Prepare and organize cleaned data into dimensions and fact tables for further warehousing.

### ETL Flow:

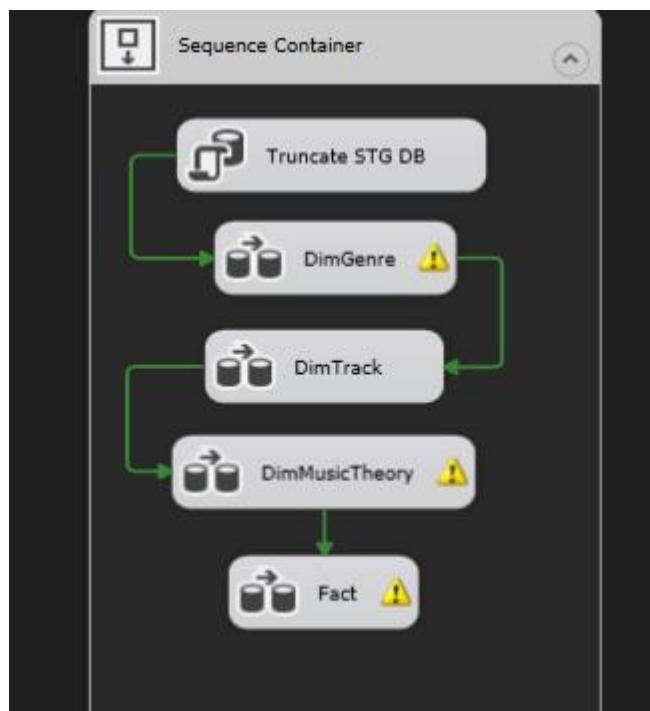
- **Source:** ODS\_Spotify database
- **Destination:** STG\_Spotify database

### Key Tasks:

- Truncate STG dimension and fact tables
- Store and structure data into:
  - DimGenre
  - DimTrack → with validation of track\_name, artist\_name (null & length checks)
  - DimMusicalTheory → validated time\_signature, key, and mode
  - Fact\_Spotify → includes references and prep for lookup in DWH

### SSIS Components Used:

- Sequence Container
- Data Flow Tasks
- Derived Column
- OLE DB Destination



## C. DWH Package – Warehouse Load & Lookup

**Objective:** Populate the dimensional model in DWH\_Spotify with clean, fully integrated data.

### ETL Flow:

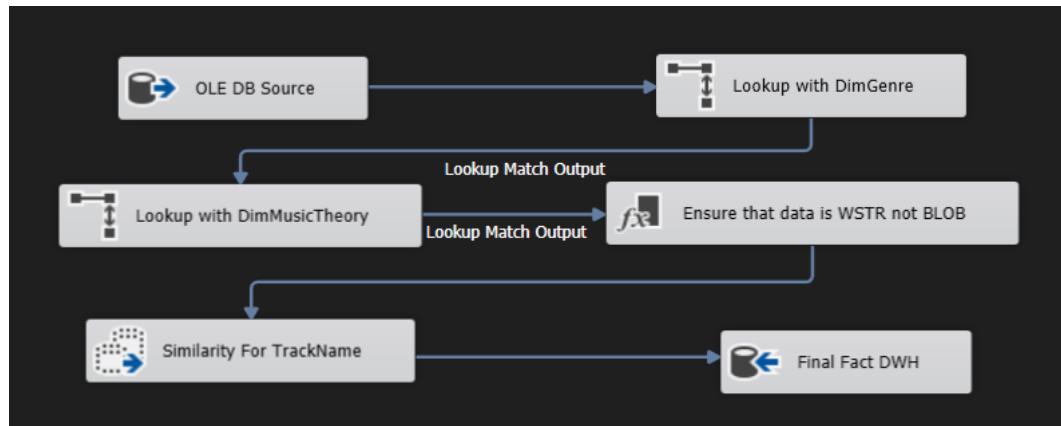
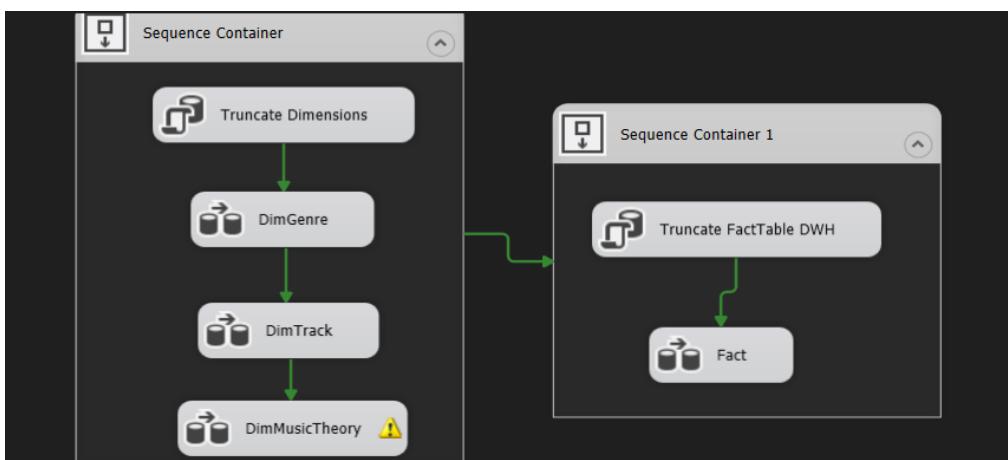
- **Source:** STG\_Spotify database
- **Destination:** DWH\_Spotify database

### Key Tasks:

- Truncate final dimension and fact tables
- Load dimensions: DimTrack, DimGenre, DimMusicalTheory
- Populate FactSpotifyTrack with:
  - **Lookups:** Genre and musical theory IDs
  - **Fuzzy Lookup:** DimTrack based on track\_name and artist\_name (similarity threshold = 0.98) to handle multilingual/inconsistent entries
  - **Data Type Enforcement:** Ensures string fields are compatible with SSAS and Power BI (avoids BLOB output)

### SSIS Components Used:

- Two Sequence Containers (Dimensions / Fact)
- Lookup
- Fuzzy Lookup
- Derived Column
- Data Flow Tasks
- OLE DB Destination



## 4. SSAS Overview

The **SSAS Tabular model** is built on top of DWH\_Spotify using the following:

### Model Structure:

- Tables: DimGenre, DimMusicalTheory, DimTrack, FactSpotifyTrack
- Relationships: Standard star schema relationships between fact and dimensions

### Measures Implemented (Sample):

Measure Type	Example Measures
Aggregates	Total Tracks, Total Artists, Average Popularity
Conditional Counts	DanceableEnergeticTracks, FastTracks, HighlyInstrumentalTracks
Percentages	Major Mode %, Danceability %, Energy %, Popularity %
Custom Metrics	Average Tracks per Artist, Popular Instrumentals

- Measures written in **DAX**
- Supports advanced analysis like artist productivity, track mood, genre energy levels

The screenshot shows the Microsoft Data Modeler interface with the following components:

- Deployment Progress:** A progress bar at the top indicates the deployment operation may take several minutes to complete.
- Deployment Results:** A summary table shows 5 Total items, 0 Cancelled, 5 Success, and 0 Error. The details table lists the successful deployment of metadata and data for DimMusicTheory, DimTrack, FactSpotifyTrack, and DimGenre.
- Data View:** The main view displays a table named "Very Danceable Tracks" with columns including fact\_id, tra..., gen..., music\_theo..., popularity, acousticness, danceability, duration\_min, energy, and instrumentane... The table shows various statistics such as Average Time Signature Confidence, Total Artists, Total Track IDs, Total Genres, and Popular Instrumentals.
- Tabular Model Explorer:** A sidebar on the right lists various measures and dimensions, including Acousticness %, Artists by Genre, and Danceability %.
- Output:** A log window at the bottom shows the build process starting at 03:25, succeeded, took 46.772 seconds, and completed at 03:26.

## 5. Power BI Overview

### Dashboard Navigation Structure:

- **Cover Page:** Contains buttons to:
  1. Overview
  2. Artist
  3. Tracks
  4. Genre/Audio
- **Navigation Tools:**
  - Filter bar
  - Home button (Spotify logo)
  - Forward/Backward navigation
  - Clear filters

