



# Harmony Hub: Data Warehouse for Music Streaming

**Team: Inside Data**

Abhishek Shinde, Sejal Sardal, Rushikesh Shinde, Shashank Guda

## Project Overview

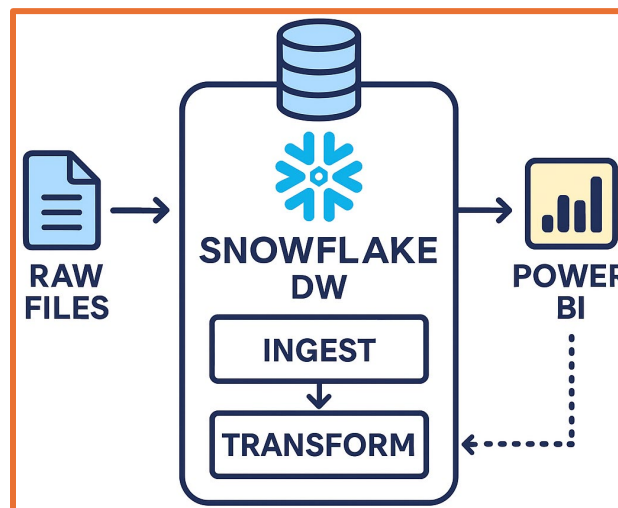
Harmony Hub is a data warehouse solution built on Snowflake for a music streaming platform. The goal of this project is to design, implement, and optimize a centralized data infrastructure to manage and analyze vast volumes of music-related data, such as:

- Tracks and genres
- Artists and albums
- User subscriptions
- Streaming session logs

The platform supports advanced analytical queries for business intelligence, user behavior tracking, subscription modeling, and genre-based insights.

## Data Warehouse Pipeline Architecture

The Harmony Hub project follows a streamlined ELT (Extract-Load-Transform) data pipeline architecture built on Snowflake Data Warehouse, enabling efficient ingestion, transformation, and analytics of large-scale music streaming data. The pipeline supports end-to-end data flow from raw file ingestion to visualization in Power BI, ensuring high performance, scalability, and modularity.





## Pipeline Components

### 1. Raw Data Ingestion

The pipeline begins with the ingestion of raw files in formats such as CSV, JSON, or Parquet. These files include critical datasets such as:

- User demographics and subscription details
- Streaming session logs
- Track and genre metadata
- Artist and album information

These files are manually uploaded into Snowflake's staging area.

### 2. Snowflake Data Warehouse (DW)

The core processing takes place in Snowflake, a cloud-native data warehouse that serves both as the ingestion layer and transformation engine.

- **Ingestion Layer:** Raw files are loaded into raw tables using Snowflake's native COPY INTO operations from internal or external stages.
- **Transformation Layer:** Using SQL-based transformation logic, data is cleaned, joined, enriched, and structured into:
  - **Fact Tables:** Including *fact\_sessions*, *fact\_user\_subscriptions*, and *bridge\_track\_genres*
  - **Dimension Tables:** Such as *dim\_users*, *dim\_tracks*, *dim\_genres*, *dim\_artists*, and *dim\_albums*

These transformations implement business rules, generate surrogate keys, and ensure referential integrity across the schema.

### 3. Analytics & Visualization (Power BI)

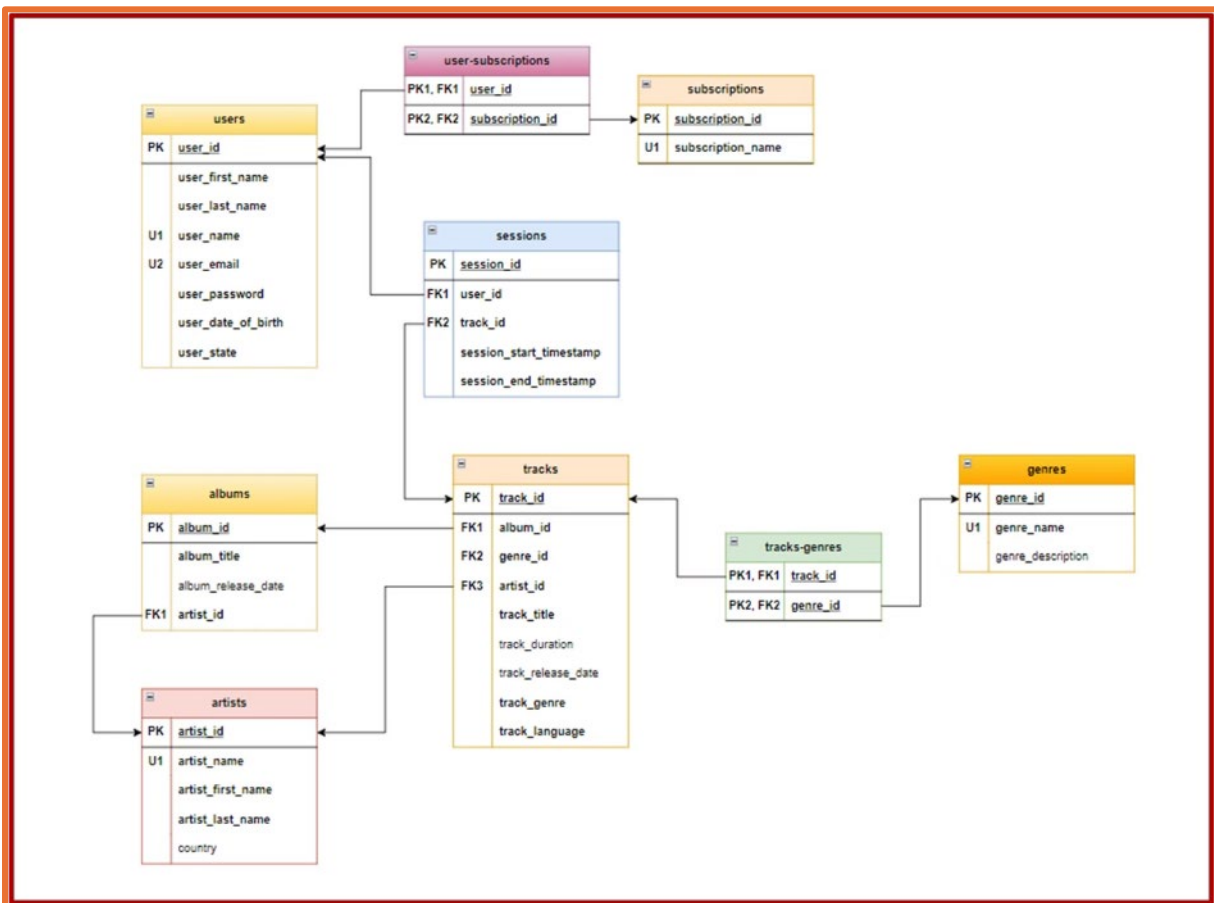
Once data is structured in Snowflake, it is connected to Power BI for advanced reporting and dashboarding. Power BI utilizes Import Mode to fetch data from Snowflake, enabling real-time and historical insights.

## Entity Relationship Design

The logical design includes a normalized relational schema to ensure data integrity and reduce redundancy.

### Key Entities:

- **Users:** Contains personal details and subscription linkage
- **Tracks:** Captures metadata of each song
- **Sessions:** Logs user activity with timestamps
- **Genres, Artists, Albums:** Support metadata organization
- **User-Subscriptions:** Many-to-many bridge for user plans
- **Tracks-Genres:** Handles the many-to-many mapping between songs and their genres



## Dimensional Modeling & Star Schema

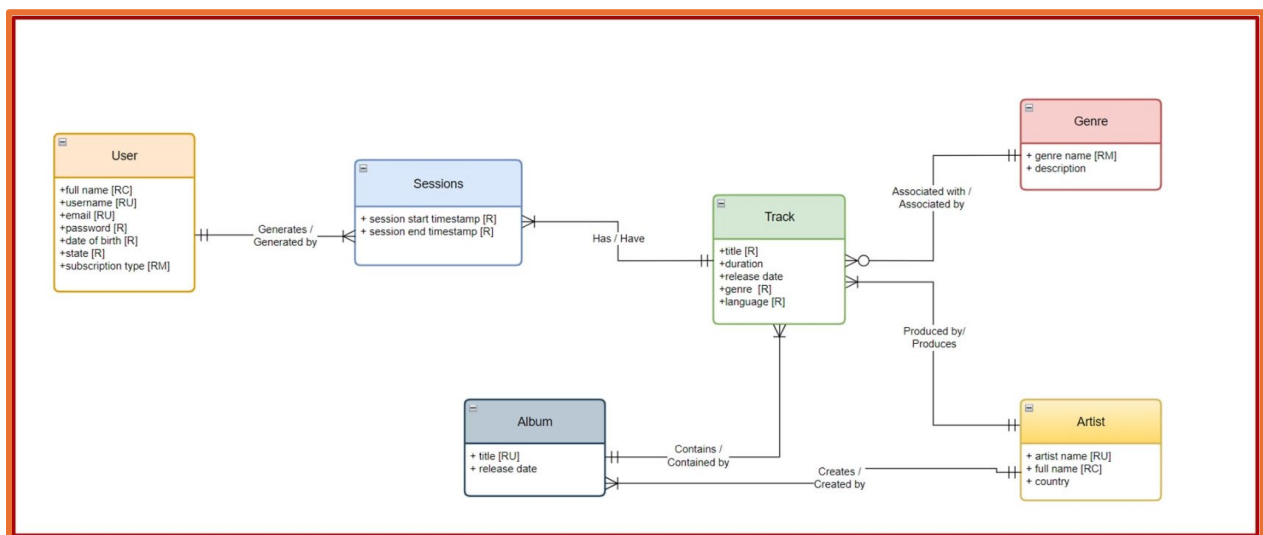
To support analytical workloads, a **star schema** was created with clearly defined fact and dimension tables. This structure enhances query performance for aggregations and time-series analysis.

### Dimensions:

- **DIM\_USERS** – User information
- **DIM\_ARTISTS** – Artist information
- **DIM\_ALBUMS** – Album metadata
- **DIM\_TRACKS** – Track details
- **DIM\_GENRES** – Genre descriptions
- **DIM\_SUBSCRIPTIONS** – Subscription type metadata
- **DIM\_DATE** – Date dimension for temporal analysis

### Fact Tables:

- **FACT\_SESSIONS** – Records of user streaming activity
- **FACT\_USER\_SUBSCRIPTIONS** – User-subscription history
- **BRIDGE\_TRACK\_GENRES** – Connects tracks and genres (many-to-many)





## Bus Matrix Design

The Bus Matrix helps in identifying the analytical capabilities across business processes and aligns them with respective fact and dimension tables.

Business Process	Fact Table	Grain	Type
User Subscription Tracking	user_subscriptions	One row per user per subscription	Transaction
Session Activity	sessions	One row per session	Transaction
Track Info	tracks	One row per track	Periodic Snapshot
Album Release Analysis	albums	One row per album	Periodic Snapshot
Genre Coverage	tracks-genres	One row per track per genre	Factless

The matrix guides which dimensions relate to which fact tables, facilitating reusable dimension modeling.

Business Process Name	Fact Table	Grainularity	Fact Grain Type	Dimension Tables						
				dim_users	dim_tracks	dim_albums	dim_artists	dim_genres	dim_date	dim_subscriptions
User Subscription Tracking	user_subscriptions	One row per user per subscription	Transaction	X						X
Session Activity	sessions	One row per session record	Transaction	X	X				X	
Track Information	tracks	One row per track	Periodic Snapshot		X	X	X	X		
Album Release Analysis	albums	One row per album	Periodic Snapshot			X	X		X	
Genre Coverage Analysis	tracks-genres	One row per track per genre	Factless Fact Table		X			X		



## Detailed Dimensional Modeling

We conducted attribute-level design for each table. Each dimension was enriched with surrogate keys, relevant metadata, and sourced directly from raw input tables. All date/time fields were normalized for compatibility with the DIM\_DATE table.

### Highlights:

- Consistent data types and naming conventions
- Use of surrogate keys for dimensional stability
- Cardinality checks for fact-dimension linkage
- Pre-calculated durations and states for faster analysis

Column Name	Description	DataType	P	uniq	not_n	Source
genrekey	dimension key		x	x	x	surrogate key
genre_id	primary key of the source systems (busir	INT	x	x	x	genres.genre_id
genres_genre_name	Genre name	VARCHAR(255)		x	x	genres.genres_genre_name
genres_description	Genre description	VARCHAR(255)				genres.genres_description
subscriptionkey	dimension key		x	x	x	surrogate key
subscription_id	primary key of the source systems (busir	INT	x	x	x	subscriptions.subscription_id
subscriptions_subscription_	Subscription name	STRING		x	x	subscriptions.subscriptions_subscription
userkey	dimension key		x	x	x	surrogate key
user_id	primary key of the source systems (busir	INT	x	x	x	users.user_id
users_first_name	User first name	STRING		x	x	users.users_first_name
users_last_name	User last name	STRING		x	x	users.users_last_name
users_user_name	User name	STRING(50)		x	x	users.users_user_name
users_user_email	User email	STRING(255)		x	x	users.users_user_email
users_date_of_birth	User date of birth	DATE		x	x	users.users_date_of_birth
users_state	User state	STRING(50)		x	x	users.users_state
users_subscription_type	User subscription type	STRING(50)		x	x	users.users_subscription_type
artistkey	dimension key		x	x	x	surrogate key
artist_id	primary key of the source systems (busir	INT	x	x	x	artists.artist_id
artists_artist_name	Artist name	STRING		x	x	artists.artists_artist_name
artists_first_name	Artist first name	STRING				artists.artists_first_name
artists_last_name	Artist last name	STRING				artists.artists_last_name
artists_country	Artist country	STRING(50)				artists.artists_country
albumkey	dimension key		x	x	x	surrogate key
album_id	primary key of the source systems (busir	INT	x	x	x	albums.album_id
albums_title	Album title	STRING		x	x	albums.albums_title
albums_release_date	Album release date	DATE				albums.albums_release_date
albums_artist_id	Artist ID (FK)	INT			x	albums.albums_artist_id
trackkey	dimension key		x	x	x	surrogate key
track_id	primary key of the source systems (busir	INT	x	x	x	tracks.track_id
tracks_album_id	Album ID (FK)	INT			x	tracks.tracks_album_id
tracks_genre_id	Genre ID (FK)	INT			x	tracks.tracks_genre_id
tracks_artist_id	Artist ID (FK)	INT			x	tracks.tracks_artist_id
tracks_title	Track title	STRING		x	x	tracks.tracks_title
tracks_duration	Track duration	INT				tracks.tracks_duration
tracks_release_date	Track release date	DATE				tracks.tracks_release_date
tracks_genre	Track genre	STRING(50)		x	x	tracks.tracks_genre
tracks_language	Track language	STRING(50)		x	x	tracks.tracks_language
sessionkey	dimension key		x	x	x	surrogate key
session_id	primary key of the source systems (busir	INT	x	x	x	sessions.session_id
sessions_user_id	User ID (FK)	INT			x	sessions.sessions_user_id
sessions_track_id	Track ID (FK)	INT			x	sessions.sessions_track_id
sessions_start_timestamp	Session start timestamp	TIMESTAMP_NTZ	x		x	sessions.sessions_start_timestamp
sessions_end_timestamp	Session end timestamp	TIMESTAMP_NTZ				sessions.sessions_end_timestamp



## Business Problem: Regional Genre Popularity Analysis

In a competitive music streaming industry, delivering personalized and geographically relevant content is essential for user engagement, retention, and monetization. Understanding regional listening preferences enables streaming services, record labels, and advertisers to align their strategies with the cultural and behavioral nuances of users across different states. However, deriving such insights from large-scale user activity data requires a robust and scalable data infrastructure. This business problem is addressed through the implementation of a Snowflake-based data warehouse in the Harmony Hub project, which supports structured analytical workflows and efficient querying.

### Role of the Data Warehouse

The Harmony Hub data warehouse is designed to consolidate and model streaming activity data in a way that supports complex, high-value analytical queries. By integrating user demographics, session activity, and track metadata, the warehouse enables organizations to generate actionable insights. Specifically, the warehouse schema includes:

- **User Information** (dim\_users) with location and subscription data
- **Streaming Sessions** (fact\_sessions) capturing session-level engagement
- **Track and Genre Metadata** (dim\_tracks, dim\_genres, bridge\_track\_genres)

Through this architecture, SQL queries can be executed efficiently to join session data with genre classifications and user state information, enabling a detailed breakdown of musical preferences across regions.

### Query Insight: Most Popular Genres per State

This specific analytical query provides a state-wise breakdown of the most popular music genres, evaluated through both the number of listening sessions and the total listening duration. The insight produced is highly valuable across multiple business functions:

- **Artists and Record Labels:** Gain visibility into where their genre of music is performing best, aiding in targeted marketing and tour planning.
- **Advertising and Marketing Teams:** Can optimize campaign targeting by aligning messaging with the musical preferences of specific geographic segments.

**Top 3 popular genres per State.**

```
-- Query 1: Top 3 popular genres per State.
WITH ranked_genres AS (
  SELECT
    users.users_state,
    genres.genres_genre_name AS genre,
    COUNT(*) AS playback_count,
    SUM(DATEDIFF(minute, sessions.sessions_start_timestamp, sessions.sessions_end_timestamp)) AS total_listening_minutes,
    ROW_NUMBER() OVER (PARTITION BY users.users_state ORDER BY SUM(DATEDIFF(minute, sessions.sessions_start_timestamp, sessions.sessions_end_timestamp)) DESC)
    AS genre_rank
  FROM sessions
  JOIN users ON sessions.sessions_user_id = users.user_id
  JOIN tracks ON sessions.sessions_track_id = tracks.track_id
  JOIN tracks_genres ON tracks.track_id = tracks_genres.tracks_genres_track_id
  JOIN genres ON tracks_genres.tracks_genres_genre_id = genres.genre_id
  GROUP BY users.users_state, genres.genres_genre_name
)
SELECT users_state, genre, playback_count, total_listening_minutes
FROM ranked_genres
WHERE genre_rank <= 3
ORDER BY users_state, genre_rank;
```

**Query**

**Output**

	users_state	genre	playback_count	total_listening_minutes
1	Arizona	Ambient House	3	11
2	Arizona	Hip Hop	3	11
3	Arizona	World	4	10
4	California	Blues	3	12
5	California	Grime	3	12
6	California	Classical	3	9
7	Colorado	Electropop	2	7
8	Colorado	Alternative	2	7
9	Colorado	Indie	1	5
10	Florida	World	5	13
11	Florida	Funk	3	12
12	Florida	House	5	9

For example, if the analysis shows that 'Ambient House' has the highest listening duration in Arizona, while 'Blues' is dominant in California, it signals strong regional differentiation in musical taste. Such findings inform product, content, and marketing decisions at both strategic and operational levels.

## Business Process Context

This analysis aligns with the business process of **“Genre Coverage Analysis”**, as outlined in the project's bus matrix. The process leverages the tracks-genres factless fact table and joins it with relevant dimension tables to derive insights. The use of dimensional modeling ensures scalability and reusability of components across other analytical processes within the platform.

## Business Problem: Artist Popularity by Age Demographics

In the highly personalized world of music streaming, understanding how different demographic segments interact with artists is key to effective audience engagement. One of the most valuable dimensions in audience analysis is age group, which often correlates strongly with musical taste, content preferences, and consumption patterns. For artists, record labels, and streaming platforms, insights into generational listening behavior are essential for targeting and outreach.





This business problem is addressed within the Harmony Hub data warehouse through age-based segmentation and user activity analysis, enabled by a robust dimensional model and efficient query design.

## Role of the Data Warehouse

The Harmony Hub data warehouse enables this analysis through the integration of user demographic data with artist-level streaming metrics. The relevant data points include:

- **User Demographics** from the dim\_users table (age, gender, location)
- **Streaming Sessions** from the fact\_sessions table (track played, duration, timestamp)
- **Artist Metadata** via the dim\_artists and dim\_tracks tables (track-to-artist mapping)

The data is structured to allow for analytical queries that group session data by artist and user age group, making it possible to identify which artists resonate with which generational cohorts.

## Query Insight: Artist Popularity Across Age Groups

### Artist Popularity Across Different Age Groups

```
-- Query 3: Artist Popularity Across Different Age Groups
SELECT
  artists.artists_artist_name,
  CASE
    WHEN DATEDIFF(year, users.users_date_of_birth, GETDATE()) < 18 THEN 'Under 18'
    WHEN DATEDIFF(year, users.users_date_of_birth, GETDATE()) BETWEEN 18 AND 25 THEN '18-25'
    WHEN DATEDIFF(year, users.users_date_of_birth, GETDATE()) BETWEEN 26 AND 35 THEN '26-35'
    WHEN DATEDIFF(year, users.users_date_of_birth, GETDATE()) BETWEEN 36 AND 50 THEN '36-50'
    ELSE 'Above 50'
  END AS age_group,
  COUNT(*) AS session_count
FROM sessions
JOIN users ON sessions.sessions_user_id = users.user_id
JOIN tracks ON sessions.sessions_track_id = tracks.track_id
JOIN artists ON tracks.tracks_artist_id = artists.artist_id
GROUP BY artists.artists_artist_name,
  CASE
    WHEN DATEDIFF(year, users.users_date_of_birth, GETDATE()) < 18 THEN 'Under 18'
    WHEN DATEDIFF(year, users.users_date_of_birth, GETDATE()) BETWEEN 18 AND 25 THEN '18-25'
    WHEN DATEDIFF(year, users.users_date_of_birth, GETDATE()) BETWEEN 26 AND 35 THEN '26-35'
    WHEN DATEDIFF(year, users.users_date_of_birth, GETDATE()) BETWEEN 36 AND 50 THEN '36-50'
    ELSE 'Above 50'
  END
ORDER BY artists.artists_artist_name, session_count DESC;
```

Query

	artists_artist_name	age_group	session_count
1	Adele	26-35	7
2	Adele	36-50	3
3	Aerosmith	26-35	11
4	Aerosmith	36-50	3
5	Alicia Keys	36-50	10
6	Alicia Keys	26-35	5
7	Alicia Keys	18-25	1
8	Ariana Grande	26-35	11
9	Ariana Grande	36-50	6
10	Beyoncé	26-35	9
11	Beyoncé	36-50	4
12	Billie Eilish	36-50	5
13	Billie Eilish	26-35	3

Output

This analysis identifies the most popular artists within distinct age brackets, such as “Under 18,” “18–30,” “31–50,” and “Above 50.” It aggregates listening metrics—such as session count and total play duration—by user age and artist.

These insights enable the following strategic outcomes:



- **Artists and Labels:** Can segment their marketing campaigns and tailor content releases to the age groups most aligned with their current fanbase.
- **Streaming Services:** Can improve personalization algorithms and playlist generation by factoring in age-specific preferences.
- **Event Planners and Promoters:** Can use this data to forecast demand for concerts or promotions among specific age demographics.

For example, if an artist is shown to be particularly popular with listeners aged 18–30, their promotional efforts such as social media campaigns or concert tours can be optimized to target that audience. Conversely, if another artist shows a strong following among users over 50, outreach could focus on more traditional channels or curated experiences for mature listeners.

### **Business Process Context**

This analysis supports the business process of **“User Segmentation by Demographics”**, enabling data-driven decision-making for audience targeting. It leverages facts from the fact\_sessions table and dimensions from dim\_users and dim\_artists, demonstrating the power of the warehouse’s star schema design in enabling multidimensional insights.

---

## **Business Problem: Track Popularity and Listening Behavior by Time of Day**

In a data-driven music streaming environment, understanding temporal listening behavior is vital for curating relevant content, improving user experience, and optimizing system performance. Listeners’ preferences and engagement levels often vary significantly depending on the time of day, reflecting lifestyle patterns such as work, study, exercise, or relaxation routines.

This business problem is addressed through time-based session analytics within the Harmony Hub data warehouse. The ability to measure and analyze track popularity and session duration across different hours of the day offers key insights for strategic content planning and user engagement.

### **Role of the Data Warehouse**

The Harmony Hub data warehouse integrates timestamped session data with track metadata, enabling time-of-day analysis through the following components:



- **Session Logs** from fact\_sessions, capturing session start times, end times, and associated track IDs
- **Track Metadata** from dim\_tracks for context on each track
- **Date and Time Dimensions** via dim\_date and an optional dim\_time or derived hour field

The data model supports efficient filtering and aggregation to analyze listening activity across 24-hour intervals, enabling precise insights into user behavior throughout the day.

### Query Insight: Track Popularity and Session Duration by Time of Day

#### Track Popularity and Session Duration by Time of Day

```
-- Query 4: Track Popularity and Session Duration by Time of Day
WITH ranked_tracks AS (
    SELECT
        DATEPART(hour, sessions.sessions_start_timestamp) AS hour_of_day,
        tracks.tracks_title,
        COUNT(*) AS session_count,
        SUM(DATEDIFF(minute, sessions.sessions_start_timestamp, sessions.sessions_end_timestamp)) AS total_listening_minutes,
        ROW_NUMBER() OVER (PARTITION BY DATEPART(hour, sessions.sessions_start_timestamp) ORDER BY SUM(DATEDIFF(minute, sessions.sessions_start_timestamp, sessions.sessions_end_timestamp)) DESC) AS track_rank
    FROM sessions
    JOIN tracks ON sessions.sessions_track_id = tracks.track_id
    GROUP BY tracks.tracks_title, DATEPART(hour, sessions.sessions_start_timestamp)
)
SELECT hour_of_day, tracks_title, session_count, total_listening_minutes
FROM ranked_tracks
WHERE track_rank <= 3
ORDER BY hour_of_day, track_rank;
```

Query

Results

	hour_of_day	tracks_title	session_count	total_listening_minutes
1	0	Into You	2	9
2	0	Lucy in the Sky with Diamonds	2	8
3	0	In My Feelings	2	7
4	1	Blackbird	2	7
5	1	Circles	1	5
6	1	Bad Romance	1	5
7	2	The Pretender	2	7
8	2	Needed Me	1	5
9	2	Brown Sugar	2	5
10	3	Back in Black	3	9
11	3	Just Dance	2	8
12	3	Sorry	2	8
13	4	Go Your Own Way	3	7
14	4	Shake It Off	2	6

Output

This query provides a breakdown of the most popular tracks and average session durations across different hours (e.g., early morning, midday, evening, late night). It enables the identification of **hour-specific trends**, such as:

- Tracks with the highest engagement during **morning commutes** (e.g., energetic or motivational music)
- Relaxing tracks gaining popularity during **evening hours**
- Increased session durations during **late-night streaming**, possibly indicating deep engagement or passive listening



These insights enable several business outcomes:

- **Streaming Platforms:** Can dynamically update homepage recommendations and playlists based on real-time listening trends
- **Content Teams:** Can schedule feature placements or track promotions during peak engagement windows
- **Advertisers:** Can target campaigns with contextually appropriate messaging (e.g., fitness ads during morning sessions, wellness ads in the evening)

For instance, discovering that lo-fi chillhop tracks peak in the 9 PM to 11 PM window can inform content promotion strategies and targeted playlist updates.

### **Business Process Context**

This analysis is aligned with the business process of “**Temporal Listening Analysis**”, which is central to enhancing user engagement through contextual personalization. It utilizes fact\_sessions with derived time-of-day groupings and joins with track-level metadata, demonstrating how time-based behavioral analytics are supported by the dimensional warehouse design.

## **Power BI Dashboards: Business Value Through Data Warehousing**

The Power BI dashboards built for the Harmony Hub project serve as a powerful visualization layer atop the Snowflake Data Warehouse. They translate raw streaming, subscription, and user data into actionable insights, addressing key business challenges across user engagement, content strategy, and operational planning.

### **1. Executive Overview Dashboard: Monitoring Platform Health**

This dashboard functions as the centralized reporting layer for key performance indicators. It supports stakeholders in evaluating platform activity, subscription adoption, content performance, and user engagement trends.

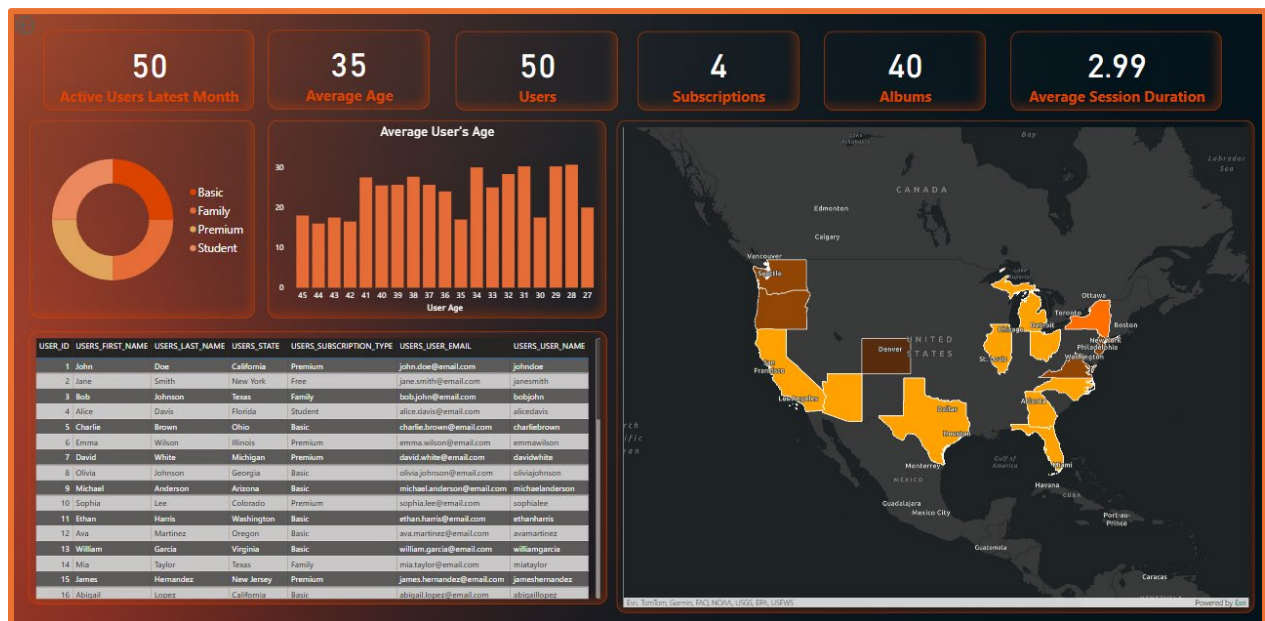
#### **Business Problems Addressed:**

- How is the overall user engagement trending over time?
- Are subscription plans being adopted consistently across the platform?
- Which music genres and album releases are gaining traction month-over-month?
- How effective are user retention strategies based on session duration?



The underlying warehouse consolidates user activity logs, genre classifications, album metadata, and subscription types into well-structured dimension and fact tables. This structure enables seamless aggregation, time-series comparison, and cross-category filtering within Power BI.

## 2. Users & Subscriptions Dashboard: Understanding Audience Segments





This dashboard provides a deep dive into user demographics, subscription preferences, and geographic distribution. It facilitates segmentation-based strategy formulation for product, marketing, and customer experience teams.

### Business Problems Addressed:

- What are the dominant age groups engaging with the platform?
- How do subscription preferences vary across different demographic or geographic segments?
- Which regions show higher engagement, and where are growth opportunities?

Data from *dim\_users*, *dim\_subscriptions*, and *fact\_sessions* is joined within the warehouse to offer a complete profile for each user. This allows Power BI to filter and visualize user behavior by age, region, and subscription type supporting personalized marketing and regional campaign planning.

### 3. Artists & Tracks Dashboard: Informing Content and Licensing Strategy



This dashboard focuses on artist popularity, genre preferences, and the lifecycle of track engagement. It is especially useful for content managers and record labels aiming to optimize catalog curation and promotional efforts.



### Business Problems Addressed:

- Which artists are most influential across different listener segments?
- How does genre popularity shift over time?
- Are track releases aligned with seasonal engagement patterns?

By modeling *dim\_artists*, *dim\_tracks*, *dim\_genres*, and their relationships with *fact\_sessions*, the data warehouse enables detailed artist- and genre-level analytics. Power BI consumes this structured data to generate comparative views and performance rankings, supporting smarter content acquisition and promotion.

Across all dashboards, Snowflake's data warehouse architecture ensures:

- **Centralized, accurate, and cleansed data** for trustworthy reporting
- **Optimized schema designs** (star schema and bridge tables) for performance
- **Time-based snapshotting and trend analysis** using dimensional modeling

These dashboards not only answer key business questions but also serve as tools for **continuous optimization** enabling Harmony Hub to evolve its content and service offerings based on rich, data-driven insights.

---