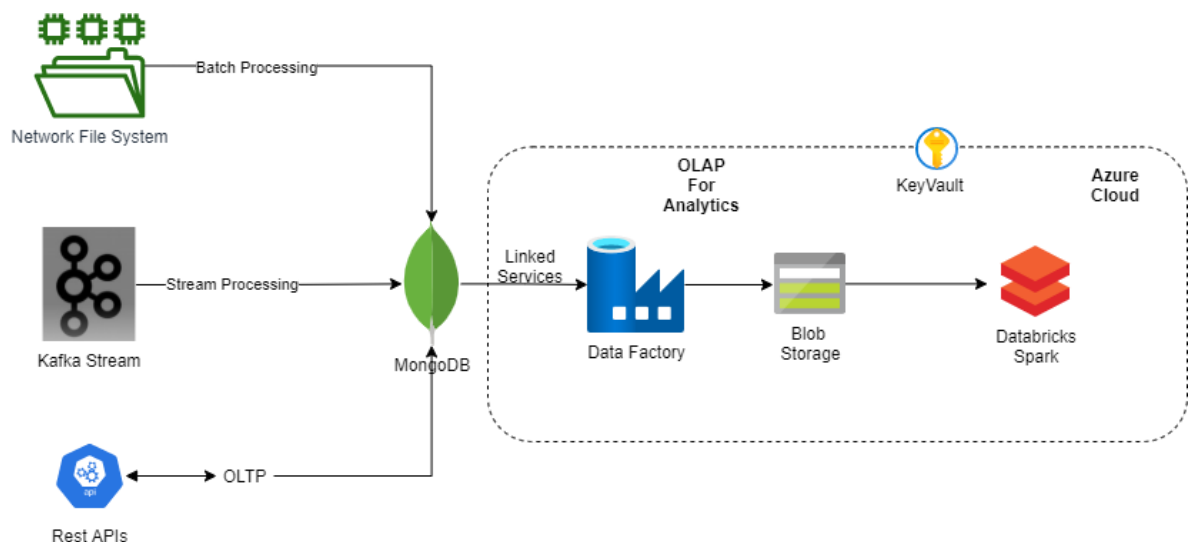# Problem: Spotify Music Recommendation Platform

With multiple audio features such as danceability, energy, and valence, this dataset can be used to build a music recommendation system. By analyzing the preferences of users for certain genres and characteristics of tracks, the system can suggest similar tracks or even recommend new genres that users might enjoy.

# Prepared By : Bhavani Shankar Sunkara(2022OG04022) , Pooja Sharma (2022OG04007)

## Assignment Part1 : Design



# Architecture and Technology Tech Stack :

**Design Components :**

For the backend services Design, APIs, data models and the databases for Spotify Music Recommendation Platform we consider the following Design

## Backend Services :

**Spark Rest API :**
    A Spark REST API facilitates interaction with Spark's functionalities through HTTP requests. It involves defining endpoints to trigger Spark operations like data processing, analytics, or machine learning tasks. Proper serialization, error handling, security measures, and performance optimization are crucial for a reliable and efficient API. This API acts as an interface, allowing users to interact with Spark's capabilities via HTTP calls.
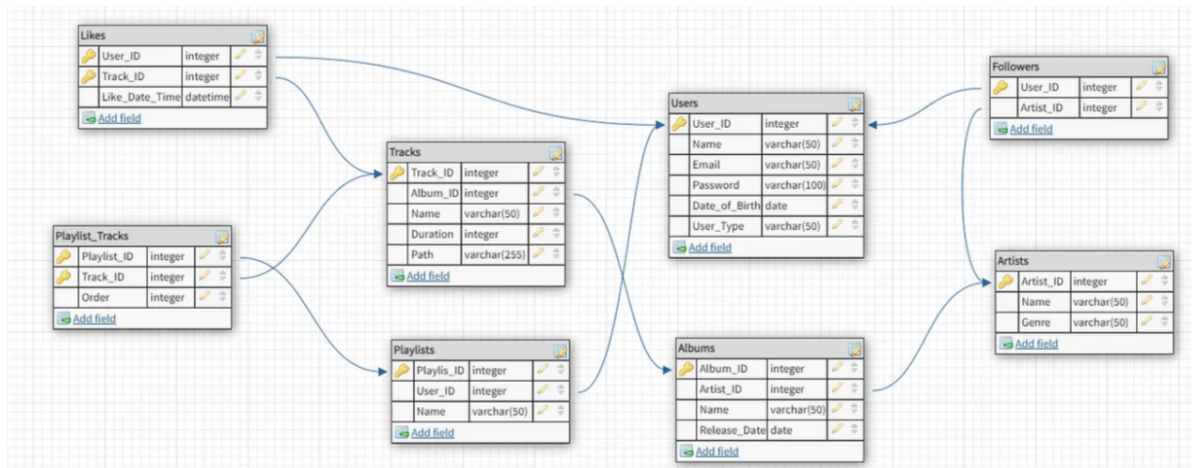
## Rest API :

**Recommendation API :**
    Accept User Input (e.g. track_id,album_name, Artists, track_genre) and returns recommendted track or genres based on User preference and track characteristics.

## Data Models:

Tracks Genere : Stores audio features and meta data for each track , including genre information.



DataBase:
    NoSQL DataBase (MongoDB) : Spotify DB Stores track data , including audio features and meta data. Offers structured data storageand efficient querying capabilities

# Design Support both Realtime stream processing and batch analytics

Batch Processing : Apache Spark can be used for batch processing of data.
Apache Spark (pyspark and NOSQL for OLAP) : Provides a big data batch processing environment from Kafka for performing analytics queries on the dataset.

**The CAP theorem refers to the trade-offs between Consistency, Availability, and Partition Tolerance in distributed systems. In designing a system, balancing these factors is essential:**

Consistency: Ensures all nodes in a distributed system have the most recent data updates. Strong consistency may lead to increased latency or reduced availability.
- **Availability:** Guarantees every request receives a response, even if some nodes fail. Prioritizing availability might lead to sacrificing consistency.
Partition Tolerance: Ensures the system continues to operate despite network partitions or communication failures.
In the context of a Spotify-like music recommendation system:

- **Consistency vs. Availability:** For real-time recommendations, availability might be prioritized over strong consistency to ensure responsiveness. Sacrificing immediate consistency could ensure that users receive recommendations promptly, even if data across nodes isn't fully synchronized.

- **Partition Tolerance:** Essential in a distributed system like Spotify, ensuring uninterrupted service despite potential network disruptions. Prioritizing partition tolerance allows the system to handle failures or network partitions without compromising its operation.
Design choices regarding data consistency, availability, and partition tolerance must consider the specific requirements of real-time streaming and batch analytics use cases. Striking a balance between these aspects is crucial to ensure system performance, data integrity, and user experience in a distributed environment.

Why Apache Spark? over Pig or Hive because of the following reasons -

- Apache Spark stands out compared to tools like Pig or Hive due to several reasons:
- **Speed and Performance:** Spark's in-memory computation and DAG (Directed Acyclic Graph) execution model significantly boost performance compared to MapReduce-based systems like Pig or Hive. Spark's ability to cache data in memory and perform iterative computations efficiently results in faster processing times.
- **Versatility and Rich APIs:** Spark offers a wide range of libraries (Spark SQL, MLlib, GraphX) that support various data processing tasks like SQL queries, machine learning, graph processing, and streaming analytics. This versatility allows Spark to handle diverse workloads within a single framework, which Pig or Hive might handle separately.
- **Ease of Use and Development:** Spark provides user-friendly APIs in multiple languages (Scala, Python, Java), making it accessible to developers with different skill sets. Its high-level APIs simplify complex tasks, reducing development effort.
- **Real-time and Batch Processing:** Spark supports both real-time streaming (Spark Streaming) and batch processing (Spark batch processing), offering a unified platform for handling both types of workloads. This capability provides more comprehensive and flexible data processing options compared to Pig or Hive, which are primarily batch-oriented.
- **Interactive Analysis:** Spark's interactive mode (Spark Shell) allows ad-hoc queries and exploratory data analysis, providing a more interactive and responsive environment for data exploration compared to Hive or Pig.
- **Community and Ecosystem:** Spark has a large and active community, extensive documentation, and a broad ecosystem of tools and integrations, offering robust support and resources for users.


Stream Processing (Kafka):

In Spotify, Kafka powers real-time stream processing by efficiently handling continuous streams of user interactions, music data, and updates. It ensures scalability, fault tolerance, and seamless integration with data processing frameworks, enabling immediate responses like personalized recommendations and real-time updates based on user actions. Kafka serves as the backbone for processing ongoing streams of data, optimizing the music streaming experience for users.


Rationale and Trade Offs:

Why Mongo DB ?  Over Neo4j, DynamoDb, HBase
Each database has unique strengths, so the choice should align with specific use case needs, performance, and architecture requirements
There are some Feature for using MongoDb for sopity MusicRecommendation system.
- High Performance for Read-heavy Workload.
- Scalability and Sharding.
- Rich Query Language and Aggregation Framework and Geospatial Capabilities.
- MongoDB offers these advantages, the choice of database ultimately depends on specific use cases, scalability needs, team expertise, and the overall architecture of the application.

    Other NoSQL databases like Cassandra, DynamoDB, or graph databases like Neo4j could also be suitable alternatives based on the specific requirements of the Spotify-like music recommendation system.