

Spotify's End-to-End Insight ETL including replication with Debezium

Rahmadiyan Muhammad





Rahmadiyan M

Education

Bina Sarana Informatika - Sistem Informasi (ongoing)

Working

PT. Bank BTPN - Big Data Infrastructure & Operation Support





Project Background



An end-to-end solution for ingesting, processing, analyzing, and visualizing Spotify user listening data. It's goal is to provide comprehensive insight into user listening patterns and preferences by leveraging various data engineering and analytics techniques.

Project repo:

https://github.com/rahmadiyann/db_final_project





Problem Statement



We employ API integration, database management, data quality checks, data processing and transformation, as well as email notification and efficient data storage and management.

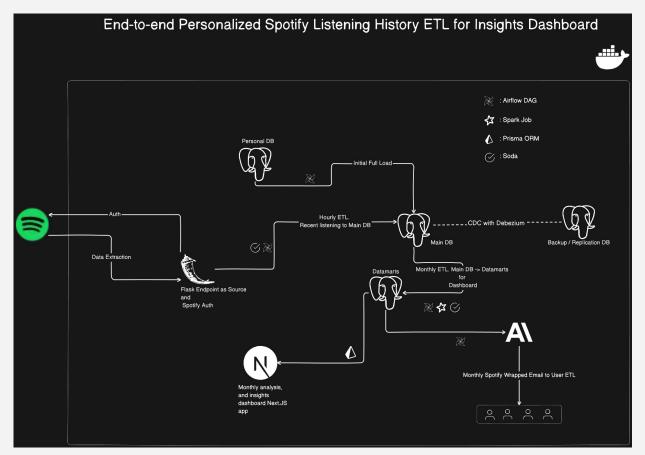
The success of the project can be measured such as data accuracy and completeness, data processing performance and the quality and relevance of generated insights and recommendation





Data Platform Understanding

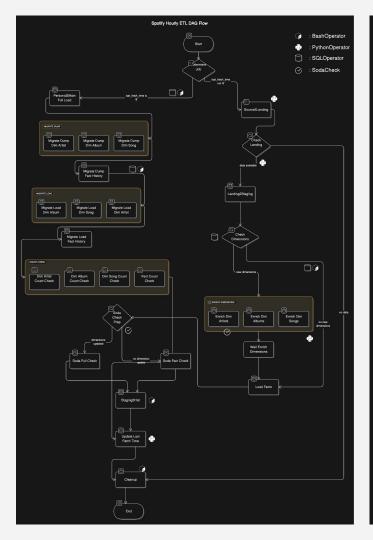


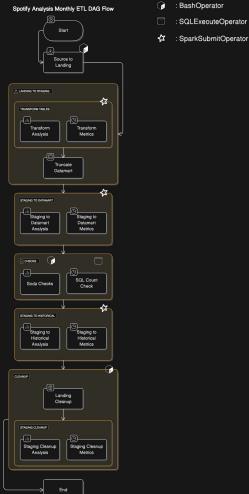




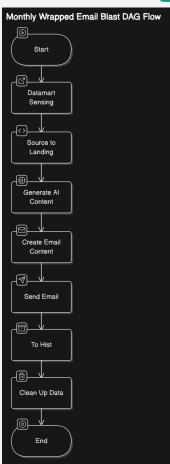


Transformation & Consideration

















Explicit_Preference	~
explicit	
play_count	
percentage	



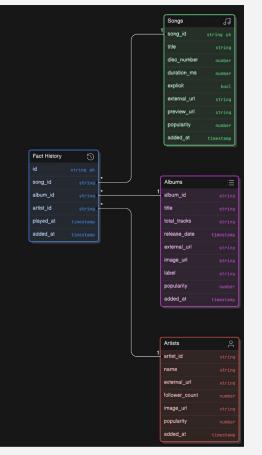
Song_Popularity_Distribution	u
popularity_bracket	
play_count	
popularity_range	
percentage	

Hour_of_Day_Listening_Distribution	0
hour_of_day	
play_count	
percentage	

Album_Release_Year_Play_Count	∷≣
release_year	
play_count	

Session_Between_Songs	เา
session_type	
percentage	
count	

Day_of_Week_Listening_Distribution	0
day_of_week	
unique_songs	
unique_artists	
song_variety_percentage	
artist_variety_percentage	



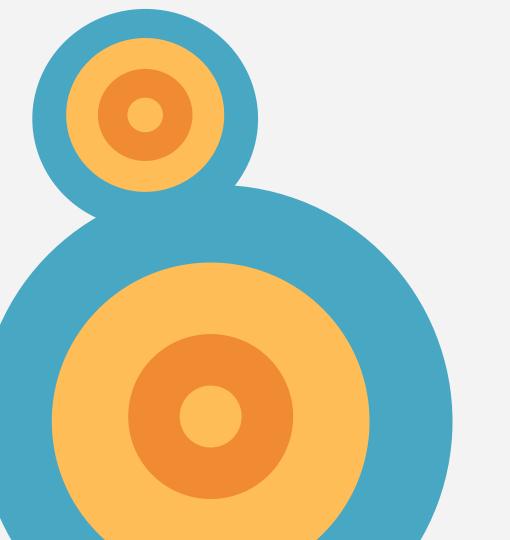






This platform successfully ingest, process, and visualize user's spotify listening data using relevant Data Engineering tools such as Spark, and Airflow as well as maintaining backup using Debezium and data quality check using Soda Core. The output is as expected in a beautiful manner that user can see through the Next.JS dashboard.





Thank you.