# EXTRACTION, TRANSFORMATION, LOADING AND VISUALIZATION OF COMBINED TWITTER AND SPOTIFY DATA IN A SCALABLE ARCHITECTURE

Thesis of the Master's degree in Business Intelligence and Big Data
in Cyber-Secure Environments

July 2022

*Adrián Riesco Valbuena*

*Tutor: Álvar Arnaiz González*

# Table of contents

# INTRODUCTION

# Introduction

**Objective:** Extraction, Transformation and Loading process to capture and visualize data from Twitter and Spotify APIs.

**Using:** Software tools from the Big Data domain and with a life cycle driven by agile methodologies.

**Highlights:**

- Service containerization.

- Data flow orchestration.

- ETL processing (data search, cleaning, merging and loading).

- Data storage.

- End user visualizations.

# OBJECTIVES

# Main objectives

👉🚩 Ability to obtain data in real time.

👉🚩 Combination of at least two different data sources.

👉🚩 Potential to scale in both technology and data volume.

👉🚩 Involvement of various technologies in the Big Data field.

👉🚩 Use of open source tools

👉🚩 Design and implementation of two visualizations for the end user.

# Functional requirements

**FR1** Data must be obtained from the Twitter hashtag #NowPlaying every 30 minutes.

**FR2** There must be at least two different visualizations and one of them must provide the ability to view all of the stored data.

**FR3** At least one of the visualizations must show last songs name, artist and audio features.

**FR4** At least one of the visualizations must have a link to the source tweet.

**FR5** At least one of the visualizations must have the ability to compare different metrics.

**FR6** At least one of the visualizations must combine two different types of visualizations.

**FR7** Both visualizations must provide sorting capabilities.

**FR8** Both visualizations must be responsive to different screen sizes.

# Technical requirements

**TR1** The development must have the ability to be deployed in different environments with minimum effort.

**TR2** The data flow must be automated, with the entire process orchestrated by a single tool.

**TR3** The execution of the ETL process must be done with a tool that can scale and run in distributed environments.

**TR4** The data warehouse must have the ability to escalate in terms of a Big Data problem.

**TR5** The web application must be designed with widely recognized tools.

**TR6** All the tools used must be open source.

# TECHNIQUES AND TOOLS

# Techniques and tools

# RELEVANT ASPECTS

# Project plan

➢ Agile methodology with sprints of two weeks and POV and backlog refinement meetings.

➢ Stored in GitHub.

GitHub

*Sprint == Milestone.*

*Tasks == Issues.*

*January*
**31**
Project start

*February*
**28**
Beginning of development

*April*
April break

*June*
**26**
Fully functional product

*July*
**13**
Project delivery

# Data Twitter

**Endpoint:** Recent search

- 🐦 **id.** Tweet id (integer), useful uniquely identify the tweet.

- 🐦 **text.** Tweet text, useful to identify the song played.

- 🐦 **entities.** Useful to clean the text and remove hashtasg, cashtags, mentions and urls.

- 🐦 **created_at.** Tweet creation date.

# Data Spotify

**Endpoint:** Search for Item

- id.
- name.
- popularity.
- artists' id.
- artists' name.

**Endpoint:** Get Tracks' Audio Features

- id.
- danceability.
- energy.
- key.
- loudness.
- mode.
- speechiness.
- acousticness.

- instrumentalness.
- liveness.
- valence.
- tempo.
- duration_ms.
- time_signature.

# Cleaned data

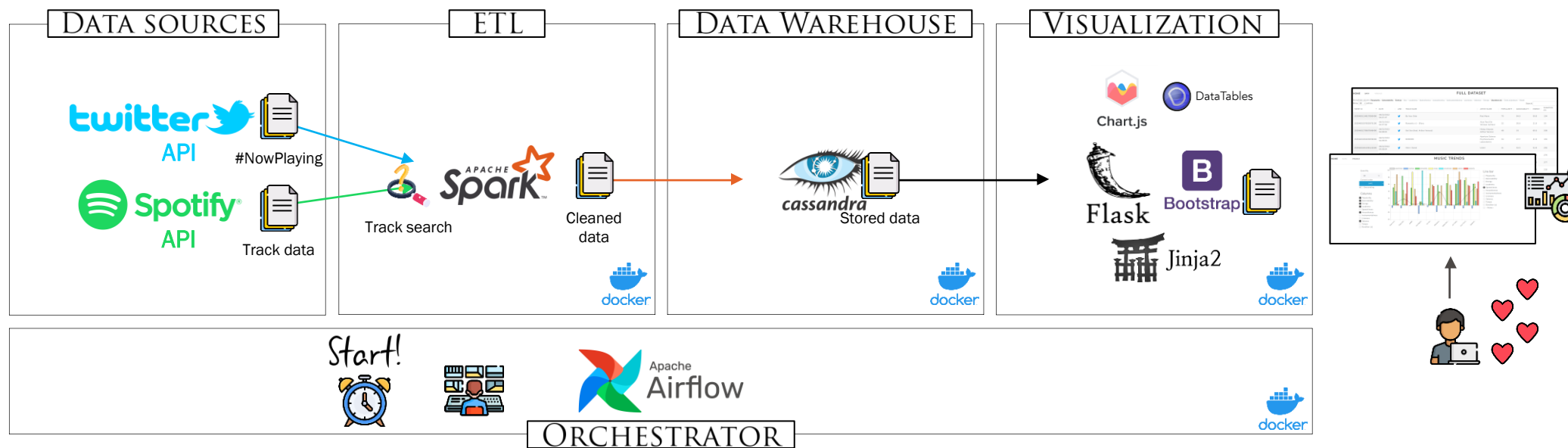| | | |
|---|---|---|
| 🐦 id_tweet | bigint | |
| 🐦 text | text | |
| 🐦 created_at | timestamp | |
| 🐦 url_tweet | text | |
| 🟢 id_track | text | |
| 🟢 name | text | |
| 🟢 popularity | int | |
| 🟢 artists_id | text | |

| | | |
|---|---|---|
| 🟢 artists_name | text | |
| 🟢 danceability | float | |
| 🟢 energy | float | |
| 🟢 key | int | |
| 🟢 loudness | float | |
| 🟢 mode | float | |
| 🟢 speechiness | float | |
| 🟢 Acousticness | float | |

| | | |
|---|---|---|
| 🟢 Instrumentalness | float | |
| 🟢 liveness | float | |
| 🟢 valence | float | |
| 🟢 tempo | float | |
| 🟢 duration_ms | int | |
| 🟢 time_signatura | float | |

# Data flow

# Data flow

# Implementation highlights and challenges

**Developer keys.** Twitter and Spotify APIs required developer keys.

**APIs' rate limits.** Maximum of 500.000 tweets per month.

**Airflow operators.** There is no specific operator for Apache Cassandra.

**Airflow and Spark connection.** Environmental variable to set connection.

**Database schema configuration at launch.** Additional container to create schema.
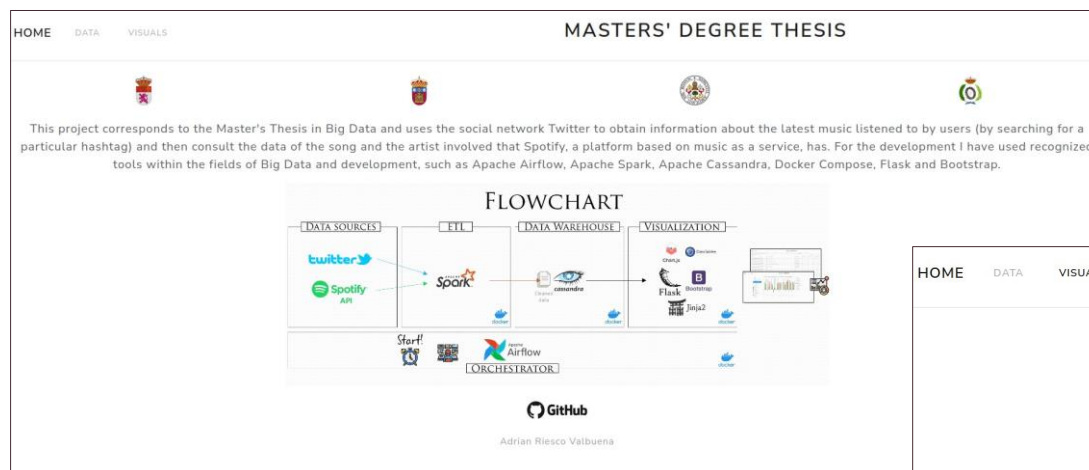
**Data representation.** Use of Datatables and Chart.js.

**Mismatched tracks.** Several languages and only the first result collected from Spotify.

# User interface

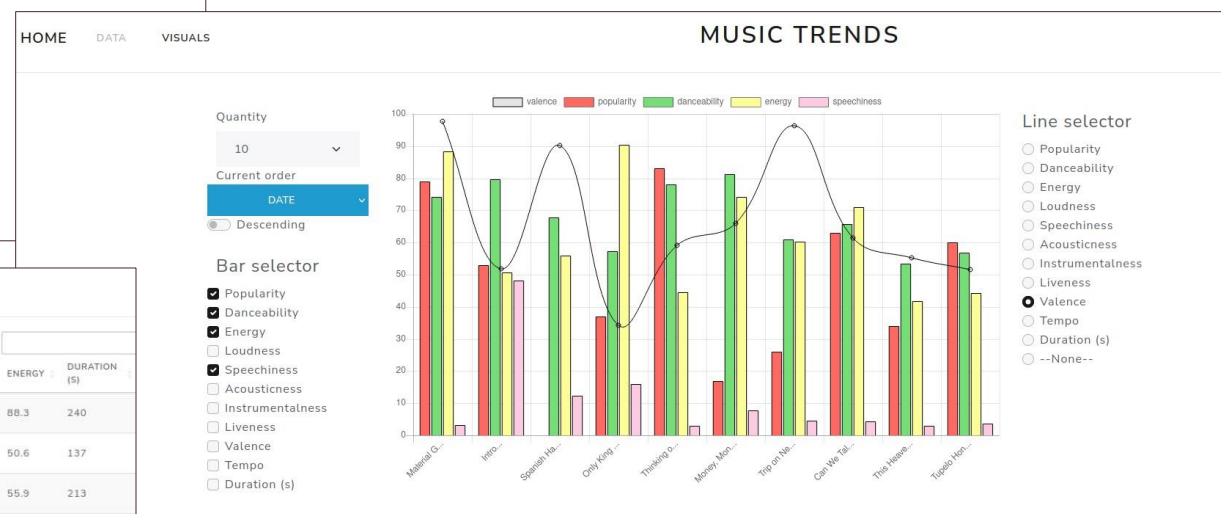# CONCLUSION AND FUTURE WORK LINES

# Conclusion and future work lines

Final result is considered a success 🏆

Future work lines:

Improve the percentage of correctly identified tracks.

Ensure that we capture as much data as possible.

Add visualizations in the front-end layer.

Upload from history functionality.

Replace Docker Compose with a more suitable tool.

Refactor the code.

Partial or total migration to the cloud.

# Questions?

# Thank you for your time!

🧑🏽‍🏫

Adrián Riesco Valbuena