

# Cấu trúc thư mục bài nộp

Bài nộp của nhóm gồm 3 folder **Slides**, **Code**, và **Dashboard** như sau:

- **Slides**: folder chứa slide thuyết trình
  - Slides.pdf: slides ở định dạng pdf
- **Code**: folder chứa code dùng để phân tích

Quá trình chạy các file

## 1. Youtube-CollectingData.ipynb

- a. **Mục tiêu**: dùng để thu thập dữ liệu về video và channel trên Youtube.
- b. **Output**: các file csv chứa dữ liệu về video (*videos\_day\_1.csv*, *videos\_day\_2.csv*, *videos\_day\_3.csv*...) và các file csv chứa dữ liệu về channel (*channels\_day\_1.csv*, *channels\_day\_2.csv*, *channels\_day\_3.csv*...) được lưu ở folder **Youtube-CollectingData-Data**

## 2. ProccessingConflictData.ipynb

- a. **Mục tiêu**: dùng để đồng nhất dữ liệu giữa các file csv video do có những file bị thiếu thông tin về 'duration', 'description' và 'tag'.
- b. **Input**: các file csv video trong folder **Youtube-CollectingData-Data**
- c. **Output**: file *dataNew.csv* trong folder **ProcessingConflictData-Data**

## 3. Youtube-Preprocessing\_1.ipynb

- a. **Mục tiêu**: dùng để gán thể loại (Genre) cho các video
- b. **Input**: file *dataNew.csv* trong folder **ProcessingConflictData-Data**
- c. **Output**: file *finalData.csv* trong folder **Youtube-Preprocessing-Data**

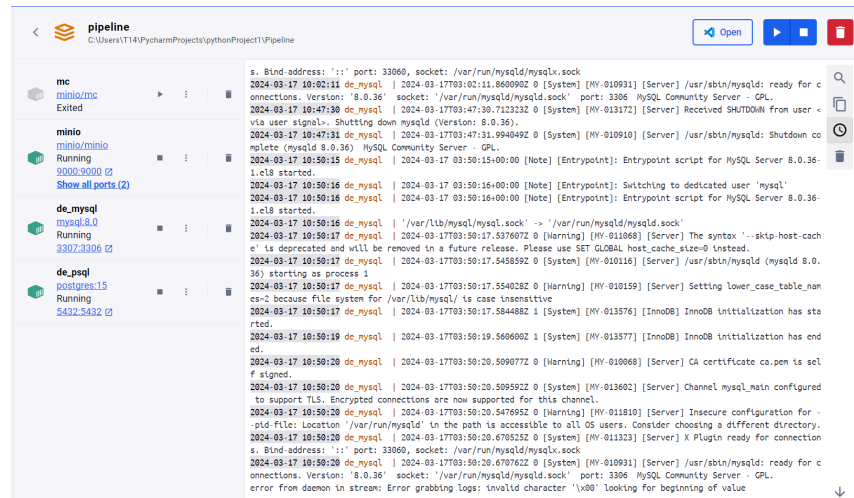
## 4. Youtube-Preprocessing\_2.ipynb

- a. **Mục tiêu**: dùng để gán thể loại (Genre) cho các video nhưng còn thiếu ở lần Preprocessing\_1
- b. **Input**: file *finalData.csv* trong folder **Youtube-Preprocessing-Data**
- c. **Output**: file *finalData\_1.csv* trong folder **Youtube-Preprocessing-Data**

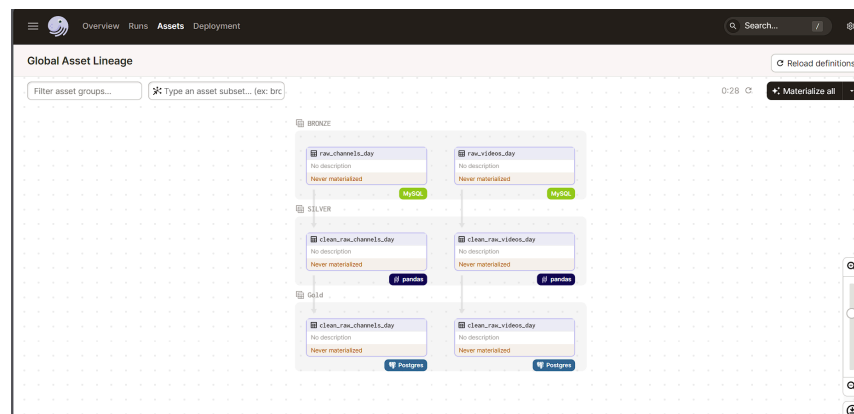
## 5. Thư mục Pipeline

### a. Mục tiêu:

- Run docker: docker compose , docker build



- Import Data vào MySQL: file import\_data.py
- Run Pipeline: Run dagster dev -f assets/one\_table.py



Click Materialize all để thực hiện các assets đưa dữ liệu vào PostgreSQL.

- Bởi vì sử dụng localhost nên sẽ xuất ra file để các thành viên khác thực hiện phân tích. File

### b. Input: finalData\_1.csv trong folder

**Youtube-Preprocessing-Data** và *channels\_day\_1.csv*,  
*channels\_day\_2.csv*, *channels\_day\_3.csv*, *channels\_day\_4.csv*,  
*channels\_day\_5.csv* trong folder **Youtube-CollectingData-Data**

### c. Output: File *clean\_raw\_channels.csv* và *clean\_raw\_videos.csv*

## 6. Youtube-Analysis-Ranked.ipynb

### a. Mục tiêu: Phân cấp rank cho channel và video

- b. Input: file *clean\_raw\_channels.csv* và file *clean\_raw\_videos.csv* trong **Thư mục Pipeline** ở bước trên.

- c. **Output:** file *finalData\_ranked.csv* lưu trong folder **Youtube-Analysis-Ranked-Data**, cũng là file data chuẩn cuối cùng được dùng cho các quá trình phân tích.

## 7. duration\_and\_view.ipynb

- a. **Mục tiêu:** Tìm mối tương quan giữa thời lượng video (Duration) và số lượt xem (View)
- b. **Input:** file *finalData\_ranked.csv* trong folder **Youtube-Analysis-Ranked-Data**
- c. **Output:** Biểu đồ scatter plot thể hiện mối tương quan giữa thời lượng video (Duration) và số lượt xem (View)

## 8. NLP\_Titles.ipynb

- a. **Mục tiêu:** Phân tích các tiêu đề (Title) video bằng tiếng Anh và đếm số lượng từ xuất hiện nhiều nhất (bao gồm từ mang tính tích cực và tiêu cực)
- b. **Input:** file *finalData\_ranked.csv* trong folder **Youtube-Analysis-Ranked-Data**
- c. **Output:**
  - + Biểu đồ Treemap và biểu đồ Word Cloud thể hiện các từ xuất hiện nhiều nhất trong các tiêu đề tiếng Anh.
  - + file *word\_count.csv* trong folder **NLP\_Titles\_Data**, được dùng cho Dashboard

## 9. AnalysisTitlePodcast.ipynb

- a. **Mục tiêu:** Phân tích độ dài title, một số khía cạnh về format của title đối với video podcast như: Title có có chứa tên channel, title được đặt ở dạng câu hỏi, Trích xuất keyword đối với mỗi title podcast và xác định xu hướng dùng từ tích cực và tiêu cực
- b. **Input:** file *finalData\_ranked.csv* trong folder **Youtube-Analysis-Ranked-Data**
- c. **Output:**
  - + Biểu đồ tròn thể hiện các phần trăm độ dài title.
  - + Biểu đồ cột cho thấy tỉ lệ số video có tên channel trong title theo các phân khúc
  - + Biểu đồ tròn tỉ lệ channel tồn tại video có title chứa tên channel theo các phân khúc
  - + Biểu đồ cột cho thấy tỉ lệ số video có title là câu hỏi
  - + Biểu đồ tròn tỉ lệ channel tồn tại video có title là câu hỏi
  - + Biểu đồ tròn tỷ lệ sử dụng từ tiêu cực trong title theo các phân khúc.

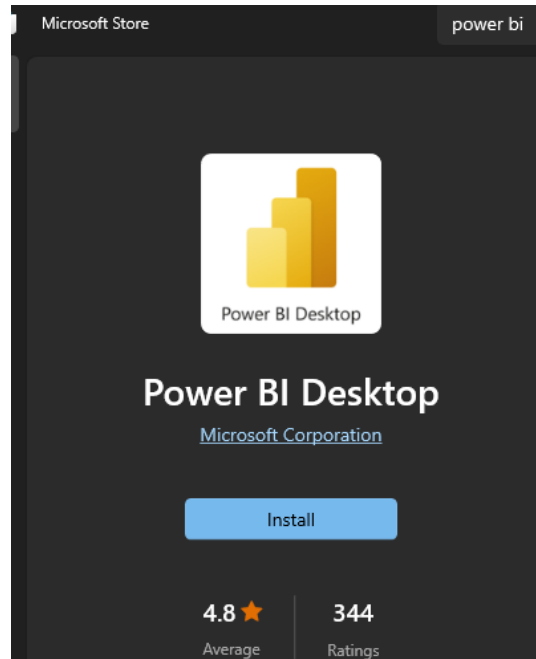
## 10. File chromedriver.exe để chạy trình duyệt Chrome ảo

- **Dashboard**: folder chứa dashboard powerBI

- youtube\_dashboard.pbix

- Cách chạy file(với Windows):

- 1. Cài đặt phần mềm Power BI Desktop từ Microsoft Store



- 2. Click chuột phải lên file youtube\_dashboard.pbix -> Chọn Open with -> Chọn Power BI Desktop

