### TABLE OF CONTENT

| 1. PROJECT OVERVIEW                  | 2  |
|--------------------------------------|----|
| 2. INTODUCTION                       | 2  |
| 3. DATA INGESTION AND TRANSFORMATION | 3  |
| 4. INGESTING JSON DATA               | 7  |
| 5. COMBINING DATA                    | 9  |
| 6. DATA CLEANING                     | 10 |
| 7. DATA ANALYSING                    | 16 |
| 8. BUSINESS QUESTION                 | 20 |
| 9. CONCLUSION                        | 22 |

#### 1. PROJECT OVERVIEW:

For this project, I used Snowflake's Data Lakehouse architecture to analyse YouTube trending and category data across multiple countries, including the United States, Great Britain, Germany, Canada, France, Brazil, Mexico, South Korea, Japan, and India. The dataset included trending video data in CSV format and category information in JSON format. The trending data captured key details such as video titles, channel names, view counts, likes, dislikes, comments, and trending dates, offering insights into popular content and viewing habits across different regions. The category data provided a comprehensive classification of content, with titles and unique IDs for various video formats.

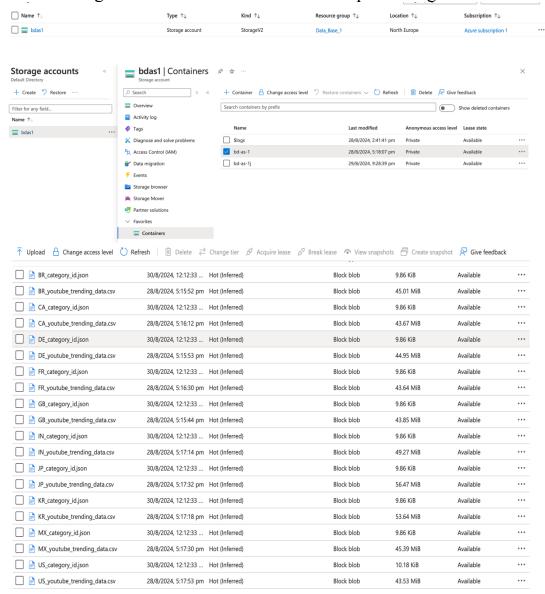
The dataset was initially uploaded to Azure cloud storage and then ingested as external tables into Snowflake. I transformed this data into structured internal tables, ensuring correct data types and integrity. Trending and category data were merged using unique identifiers generated with the UUID\_STRING() method to create the final table. Data cleaning involved handling duplicates, filling missing values, and ensuring consistency. The analysis focused on identifying top-performing content by region, examining trends in video popularity, and analyzing category preferences over time. This approach offered valuable insights into global digital content consumption trends and highlighted opportunities for optimizing YouTube content strategy.

#### 2. INTRODUCTION

The primary goal of this project is to leverage Snowflake's Data Lakehouse architecture to analyse trending and category data on YouTube from various countries. Utilising databases with trending video data and category specifics from the United States, India, Germany, and Japan, among other nations, our goal was to identify regional variations in video popularity and viewer preferences. The procedure comprised transferring the data to cloud storage on Azure, ingesting it as external tables into Snowflake, transforming the data, and carrying out comprehensive analysis. The results shed light on worldwide content trends, explain what makes videos popular on YouTube, and give content producers tactical advice.

#### 3. DATA INGESTION AND TRANSFORMATION:

The dataset, consisting of YouTube trending and category data, was initially downloaded from provided links and subsequently uploaded to Microsoft Azure Blob Storage. Azure Blob Storage was chosen for its scalable and secure data management capabilities, which facilitated smooth integration with Snowflake for further processing.



The datasets were accurately loaded from Azure Blob Storage into Snowflake, transformed into the required formats, and combined into a comprehensive final table. This setup provides a robust foundation for subsequent data analysis tasks. The use of Azure Blob Storage facilitated efficient data management, while Snowflake's powerful features enabled effective data integration and transformation

#### 3.1 DATABASE AND STAGE CREATION:

An essential first step in effectively organising and managing datasets in Snowflake is to create a new database. Tables, schemas, and other database objects are logically contained within databases in Snowflake. It aids in data structure that promotes security, efficiency in querying, and data integrity. For this project, a new snowflake database named 'assignment' was created

```
/* Creating a DataBase called Assignment1 */
```

```
CREATE OR REPLACE DATABASE assignment_1;
USE DATABASE assignment_1;
```

The 'stage\_assignment' stage was created to link snowflake with azure blob storage. Allowing snowflake to access the data files directly from the cloud storage.

```
CREATE OR REPLACE STAGE stage_assignment
URL='azure://bdas1.blob.core.windows.net/bd-as-1'
CREDENTIALS=(AZURE_SAS_TOKEN='?sv=2022-11-
02&ss=b&srt=co&sp=rwdlaciytfx&se=2024-12-30T14:19:42Z&st=2024-
08-
28T07:19:42Z&spr=https&sig=Fq9YilGA6fq18X2w9dnJJoAiTn732LYAvqGO
mPrpmP4%3D')
```

#### 3.2 INGESTING CSV DATA:

FILE FORMAT: Defining a file format in Snowflake is a critical step when ingesting data from external sources. This process ensures that Snowflake correctly interprets and processes the data according to its structure and characteristics.

```
--WE going to create a file format skipping the header, adding delimeters and type as csv to create a external table CREATE OR REPLACE FILE FORMAT file_format_csv

TYPE = 'CSV'
FIELD_DELIMITER = ','
SKIP_HEADER = 1
NULL_IF = ('\\N', 'NULL', 'NUL', '')
FIELD_OPTIONALLY_ENCLOSED_BY = '"'
:
```

A file format for CSV was defined to correctly interpret the structure of the data files.

#### **3.3 EXTERNAL TABLE CREATION:**

An external table 'ex\_table\_youtube\_trending' was created to read the csv files from azure blob storage based on the file format we created.

```
--Created a external table using the file format we created

CREATE OR REPLACE EXTERNAL TABLE ex_table_youtube_trending

WITH LOCATION = @stage_assignment

FILE_FORMAT = file_format_csv

PATTERN = '[A-Z]{2}_youtube_trending_data.csv'; -- patten will help us to bring all the csv file from the storage.
```

|    | VALUE  |
|----|--|
| 1  | { "c1": "uq5LCIQN3cE", "c10": "105756", "c11": "139946", "c2": "연녕하세요 보접입니다", "c3": "2020-08-09T09:32:48Z", "c4": "UCu9BCtGlEr73LXZsKmoujKw", "c5": "보접 BK", "c6": "24", "c7": "2020-08-12T00:00   |
| 2  | { "c1": "I-ZbZCHsHD0", "c10": "494", "c11": "3339", "c2": "早早左의 계획 [德契명 프리閩]", "c3": "2020-08-12T09:00:082", "c4": "UCRuSxVu4iqTK5kCh90ntAgA", "c5": "德契명", "c6": "1", "c7": "2020-08-12T00:00:  |
| 3  | { "c1": "9d7jNUjBoss", "c10": "68898", "c11": "50688", "c2": "행정 반성하면서 설렜습니다.", "c3": "2020-08-10T09:54:13Z", "c4": "UCMVC92EOs9yDJ05JS-CMesQ", "c5": "양평 YangPang", "c6": "22", "c7": "2020-08-10T09:54:13Z", "c4": "UCMVC92EOs9yDJ05JS-CMesQ", "c5": "양평 YangPang", "c6": "22", "c7": "2020-08-10T09:54:13Z", "c4": "UCMVC92EOs9yDJ05JS-CMesQ", "c5": "양평 YangPang", "c6": "22", "c7": "2020-08-10T09:54:13Z", "c4": "UCMVC92EOs9yDJ05JS-CMesQ", "c5": "양평 YangPang", "c6": "22", "c7": "2020-08-10T09:54:13Z", "c4": "UCMVC92EOs9yDJ05JS-CMesQ", "c5": "양평 YangPang", "c6": "22", "c7": "2020-08-10T09:54:13Z", "c4": "UCMVC92EOs9yDJ05JS-CMesQ", "c5": "양평 YangPang", "c6": "22", "c7": "2020-08-10T09:54:13Z", "c4": "UCMVC92EOs9yDJ05JS-CMesQ", "c5": "양평 YangPang", "c6": "22", "c7": "2020-08-10T09:54:13Z", "c4": "UCMVC92EOs9yDJ05JS-CMesQ", "c5": "양평 YangPang", "c6": "22", "c7": "2020-08-10T09:54:13Z", "c4": "UCMVC92EOs9yDJ05JS-CMesQ", "c5": "양평 YangPang", "c6": "22", "c7": "2020-08-10T09:54:13Z", "c4": "UCMVC92EOs9yDJ05JS-CMesQ", "c5": "양평 YangPang", "c6": "22", "c7": "2020-08-10T09:54:13Z", "c4": "UCMVC92EOs9yDJ05JS-CMesQ", "c5": "양평 YangPang", "c6": "22", "c7": "2020-08-10T09:54:13Z", "c4": "UCMVC92EOs9yDJ05JS-CMesQ", "c5": "양평 YangPang", "c6": "22", "c7": "2020-08-10T09:54:13Z", "c4": "UCMVC92EOs9yDJ05JS-CMesQ", "c5": "양평 YangPang", "c6": "22", "c7": "2020-08-10T09: "22", "c4": "2020-08-10T09: "22", "c5": "2020-08-10T09: "22", "c5": "2  |
| 4  | { "c1": "3pLL3-sM/g", "c10": "1798", "c11": "8751", "c2": "인녕하세요 퍼푸릅입니다.", "c3": "2020-08-11T15:00:58Z", "c4": "UCkQCwnkQfgSuPTTnw_Y7v7w", "c5": "퍼푸릅 Quaddurup", "c6": "24", "c7": "2020-08-12  |
| 5  | { "c1": "zrsB]YukE8s", "c10": "9763", "c1": "23405", "c2": "박전영 (J.Y. Park) When We Disco (Duet with 선미) M/V", "c3": "2020-08-11T09:00:13Z", "c4": "UCa06TYtiC8U5ttz62hTrZgg", "c5": "YP Entertainment"  |
| 6  | { "c1": "jbGRowa5tlk", "c10": "15176", "c11": "31040", "c2": "ITZY "Not Shy" M/V TEASER", "c3": "2020-08-11T15:00:13Z", "c4": "UCaO6TYtlC8U5ttz62hTrZgg", "c5": "JYP Entertainment", "c6": "10", "c7": "2020-08-11T15:00:13Z", "c4": "UCaO6TYtlC8U5ttz62hTrZgg", "c5": "JYP Entertainment", "c6": "10", "c7": "2020-08-11T15:00:13Z", "c4": "UCaO6TYtlC8U5ttz62hTrZgg", "c5": "JYP Entertainment", "c6": "10", "c7": "2020-08-11T15:00:13Z", "c4": "UCaO6TYtlC8U5ttz62hTrZgg", "c5": "JYP Entertainment", "c6": "10", "c7": "2020-08-11T15:00:13Z", "c4": "UCaO6TYtlC8U5ttz62hTrZgg", "c5": "JYP Entertainment", "c6": "10", "c7": "2020-08-11T15:00:13Z", "c4": "UCaO6TYtlC8U5ttz62hTrZgg", "c5": "JYP Entertainment", "c6": "10", "c7": "2020-08-11T15:00:13Z", "c4": "UCaO6TYtlC8U5ttz62hTrZgg", "c5": "JYP Entertainment", "c6": "10", "c7": "2020-08-11T15:00:13Z", "c4": "UCaO6TYtlC8U5ttz62hTrZgg", "c5": "JYP Entertainment", "c6": "10", "c7": "2020-08-11T15:00:13Z", "c4": "UCaO6TYtlC8U5ttz62hTrZgg", "c5": "JYP Entertainment", "c6": "10", "c7": "2020-08-11T15:00:13Z", "c4": "UCaO6TYtlC8U5ttz62hTrZgg", "c5": "JYP Entertainment", "c6": "10", "c7": "2020-08-11T15:00:13Z", "c4": "UCaO6TYtlC8U5ttz62hTrZgg", "c5": "JYP Entertainment", "c6": "10", "c7": "2020-08-11T15:00:13Z", "c4": "UCaO6TYtlC8U5ttz62hTrZgg", "c5": "JYP Entertainment", "c6": "10", "c7": "2020-08-11T15:00:13Z", "c4": "UCaO6TYtlC8U5ttz62hTrZgg", "c5": "JYP Entertainment", "c6": "10", "c7": "2020-08-11T15:00:13Z", "c4": "UCaO6TYtlC8U5ttz62hTrZgg", "c5": "JYP Entertainment", "c6": "10", "c7": "2020-08-11T15:00:13Z", "c5": "10", "c7": "2020-08-11T15:00:13Z", "c5": "10", "c5": "10" |
| 7  | { "c1": "X-TPQPEyRGo", "c10": "393", "c1": "834", "c2": "혹인 마동석이 출근하는 마트를 털러온 강도들 에게 벌어진 일 ππ", "c3": "2020-08-10T09:37:33Z", "c4": "UCpCilDf9UrfRqte55FHWIYQ", "c5": "드림텔러(DreamTelling ) 대표 기계 보이지 않는데 하는데 보다 되었다면 되었다면 하는데 보다 되었다면 되었다면 하는데 보다 되었다면 하는데 보다 되었다면 되었다면 하는데 보다 되었다면 되었다면 되었다면 되었다면 되었다면 되었다면 되었다면 되었다  |
| 8  | { "c1": "REUA4roJndU", "c10": "13337", "c11": "18286", "c2": "진심으로 사과드립니다. 좌舍합니다.", "c3": "2020-08-11T14:00:01Z", "c4": "UCwZTeeKyb1hT6sg8tf0aZiA", "c5": "임다TV", "c6": "24", "c7": "2020-08-12T   |
| 9  | { "c1": "7F0/40uehC4", "c10": "338", "c11": "5697", "c2": "집에 혼자 있을 때 하는 짓ㅋㅋㅋㅋㅋㅋㅋ [명꼬발활]", "c3": "2020-08-12T09:00:02Z", "c4": "UCCJkwrmhlqWkSFV-sQol4Qw", "c5": "명꼬발활Mingggo", "c6": "23"   |
| 10 | { "c1": "Odsnm3n6ZdM", "c10": "103", "c11": "2324", "c2": "정윤호가 무려 6시간 공들여 만든 세기의 발명품 [발명왕] Ep.1", "c3": "2020-08-11T09:30:00Z", "c4": "UC0SoPwEH3idvemSDvKaYgGA", "c5": "달라스튜디오", "c6": "달라스튜디오", "c6": "103", "c10": "103", "c10", "c1  |
| 11 | { "c1": "rFwZqtPc-Ss", "c10": "344", "c11": "1622", "c2": [책이벤트] 진짜 인플루먼서로 사는 법   황태환 @비교부부 Bgeul Bubu 하준파파 에이처유지 대표   처유 회복 사랑 컴패선   서바시 1217회", "c3": "2020-08-10T09:00:122",  |
| 12 | { "c1": "7Y8Vv_KHp7I", "c10": "1309", "c11": "1166", "c2": "ண ☑ SUB] [몰카] 누가범도 강도인데 전혀 의심을 안하는 직원을 본다면?!!! 옆 커플들 제대로 터졌음 ㅋㅋㅋㅋㅋ - [동네놈들   HOODBOYZ]", "c3": "2020-08-09T11:30:01Z   |
| 13 | { "c1": "rRaBKB9gDSQ", "c10": "2863", "c11": "16283", "c2": "뒷장고 논란에 대한 해명 및 전할 말씀이 있습니다", "c3": "2020-08-09T16:12:59Z", "c4": "UCBloXzDidCnpbM_7uyG0_Tg", "c5": "HONG SOUND", "c6": "24",   |
| 14 | { "c1": "TadJcNU1dTA", "c10": "378", "c11": "236", "c2": 세계 1위 재벌의 손자가 난처당하자 밝어진 충격적인 일, ㄷㄷ(살화", "c3": "2020-08-11107:00:042", "c4": "UCpr2S3SBmvivrx904pLUZHw", "c5": "movie trip 무비트   |

The PATTERN function specifies which files to include based on their names, using regular expressions to filter data files.

#### **3.4 SCHEMA ADJUSTMENT:**

In Snowflake, correcting the external table schema involves adjusting data types and extracting relevant information to ensure that the data is accurately represented and accessible.

```
CREATE OR REPLACE EXTERNAL TABLE ex_table_youtube_trending
Videoid STRING as (value:c1::STRING),
TITLE STRING as (value:c2::STRING),
PUBLISHEDAT DATE as (value:c3::DATE),
CHANNELID STRING as (value:c4::STRING),
CHANNELTITLE STRING as (value:c5::STRING),
CATEGORYID INT as (value:c6::INT),
TRENDING_DATE DATE as (value:c7::DATE),
VIEW_COUNT INT as (value:c8::INT),
LIKES INT as (value:c9::INT),
DISLIKES INT as (value:c10::INT),
COMMENT_COUNT INT as (value:c11::INT),
COUNTRY STRING AS SUBSTRING (METADATA $FILENAME, 1, 2)
WITH LOCATION = @stage_assignment
FILE_FORMAT = file_format_csv
PATTERN = '[A-Z]{2}_youtube_trending_data.csv';;
```

|    | VALUE                  | VIDEOID     | TITLE                     | PUBLISHEDAT | CHANNELID            | CHANNELTITLE        | CATEGORYID | TRENDING_DATE | VIEW_CC |
|----|------------------------|-------------|---------------------------|-------------|----------------------|---------------------|------------|---------------|---------|
| 1  | { "c1": "J78aPJ3VyNs", | J78aPJ3VyNs | I left youtube for a mor  | 2020-08-11  | UCYzPXprvl5Y-Sf0g4v  | jacksepticeye       | 24         | 2020-08-12    | 203     |
| 2  | { "c1": "9nidKH8cM38   | 9nidKH8cM38 | TAXI CAB SLAYER KILL      | 2020-08-11  | UCFMbX7frWZfuWdjAl   | Eleanor Neale       | 27         | 2020-08-12    | 23      |
| 3  | { "c1": "M9Pmf9AB4M    | M9Pmf9AB4Mo | Apex Legends   Stories    | 2020-08-11  | UC0ZV6M2THA81QT9     | Apex Legends        | 20         | 2020-08-12    | 238     |
| 4  | { "c1": "kgUV1MaD_M8   | kgUV1MaD_M8 | Nines - Clout (Official ' | 2020-08-10  | UCvDkzrj8ZPIBqRd6fl) | Nines               | 24         | 2020-08-12    | 61      |
| 5  | { "c1": "49Z6Mv4_WC/   | 49Z6Mv4_WCA | i don't know what im d    | 2020-08-11  | UCtinbF-Q-fVthA0qrF  | CaseyNeistat        | 22         | 2020-08-12    | 94      |
| 6  | { "c1": "ua4QMFQATcc   | ua4QMFQATco | CGP Grey was WRON(        | 2020-08-11  | UC2C_jShtL725hvbm1   | CGP Grey            | 27         | 2020-08-12    | 105     |
| 7  | { "c1": "x-KbnJ9fvJc", | x-KbnJ9fvJc | Kya Baat Aa : Karan Au    | 2020-08-11  | UCm9SZAI03Rev9sFw    | Rehaan Records      | 10         | 2020-08-12    | 1130    |
| 8  | { "c1": "3C66w5Z0ixs"  | 3C66w5Z0ixs | I ASKED HER TO BE M       | 2020-08-11  | UCvtRTOMP2TqYqu51    | Brawadis            | 22         | 2020-08-12    | 151     |
| 9  | { "c1": "ZNfeMbO_AHc   | ZNfeMbO_AHo | Popek ft. Dr Alban - It': | 2020-08-12  | UC8Mh9UmrlaQPEcyb    | KrólAlbaniiTV       | 24         | 2020-08-12    | 27      |
| 10 | { "c1": "VIUo6yapDbc"  | VIUo6yapDbc | Ultimate DIY Home Mc      | 2020-08-11  | UCDVPcEbVLQgLZX0F    | Mr. Kate            | 26         | 2020-08-12    | 112     |
| 11 | { "c1": "pAEa5BnAUgQ   | pAEa5BnAUgQ | Champions League pre      | 2020-08-11  | UC2-0sE0YbQFuaURc    | OneFootball English | 17         | 2020-08-12    | 21      |
| 12 | { "c1": "cSlrKfmgR1s", | cSIrKfmgR1s | Sensible Transfers: Ars   | 2020-08-12  | UCGYYNGmyhZ_kwBF     | Tifo Football       | 17         | 2020-08-12    | 12      |
| 13 | { "c1": "w-hVKXmib1c"  | w-hVKXmib1c | I've Seen Something K     | 2020-08-10  | UCHhfSXoDG6gSgpOv    | Behzinga            | 20         | 2020-08-12    | 87      |
| 14 | { "c1": "xTpr7piQu2M". | xTpr7piQu2M | FUT 21   Official Trailer | 2020-08-10  | UCovaxd5LQSuP4Chk    | EA SPORTS FIFA      | 20         | 2020-08-12    | 121     |

The external table schema was corrected to match the data types and using substring function we extract the country information from the file name And included in the table.

#### 3.5 CREATE INTERNAL TABLE:

An internal table 'table\_youtube\_trending' was created and all the data was moved from external table to internal. Since we are allowed to modify the

internal table, The VALUE column was dropped since it was no longer needed.

```
CREATE OR REPLACE TABLE table_youtube_trending AS
SELECT *
FROM ex_table_youtube_trending
;
ALTER TABLE table_youtube_trending DROP COLUMN VALUE ;
SELECT*
FROM table_youtube_trending;
```

|    | VIDEOID     | TITLE  | PUBLISHEDAT | CHANNELID                | CHANNELTITLE        | CATEGORYID | TRENDING_DATE |
|----|-------------|--|-------------|--------------------------|---------------------|------------|---------------|
| 1  | J78aPJ3VyNs | I left youtube for a month and THIS is what happened.              | 2020-08-11  | UCYzPXprvl5Y-Sf0g4vX-m6g | jacksepticeye       | 24         | 2020-08-12    |
| 2  | 9nidKH8cM38 | TAXI CAB SLAYER KILLS 'TO KNOW HOW IT FEELS'                       | 2020-08-11  | UCFMbX7frWZfuWdjAML0babA | Eleanor Neale       | 27         | 2020-08-12    |
| 3  | M9Pmf9AB4Mo | Apex Legends   Stories from the Outlands – "The Endorsement"       | 2020-08-11  | UC0ZV6M2THA81QT9hrVWJG3A | Apex Legends        | 20         | 2020-08-12    |
| 4  | kgUV1MaD_M8 | Nines - Clout (Official Video)                                     | 2020-08-10  | UCvDkzrj8ZPIBqRd6flxdhTw | Nines               | 24         | 2020-08-12    |
| 5  | 49Z6Mv4_WCA | i don't know what im doing anymore                                 | 2020-08-11  | UCtinbF-Q-fVthA0qrFQTgXQ | CaseyNeistat        | 22         | 2020-08-12    |
| 6  | ua4QMFQATco | CGP Grey was WRONG   | 2020-08-11  | UC2C_jShtL725hvbm1arSV9w | CGP Grey            | 27         | 2020-08-12    |
| 7  | x-KbnJ9fvJc | Kya Baat Aa : Karan Aujla (Official Video) Tania   Sukh Sanghera D | 2020-08-11  | UCm9SZAl03Rev9sFwloCdz1g | Rehaan Records      | 10         | 2020-08-12    |
| 8  | 3C66w5Z0ixs | I ASKED HER TO BE MY GIRLFRIEND                                    | 2020-08-11  | UCvtRTOMP2TqYqu51xNrqAzg | Brawadis            | 22         | 2020-08-12    |
| 9  | ZNfeMbO_AHo | Popek ft. Dr Alban - It's My Life (prod. Claysteer)                | 2020-08-12  | UC8Mh9UmrlaQPEcybdWvQsOg | KrólAlbaniiTV       | 24         | 2020-08-12    |
| 10 | VIUo6yapDbc | Ultimate DIY Home Movie Theater for The LaBrant Family!            | 2020-08-11  | UCDVPcEbVLQgLZX0Rt6jo34A | Mr. Kate            | 26         | 2020-08-12    |
| 11 | pAEa5BnAUgQ | Champions League preview: Messi v Boateng part 2! ▶ OneFootb       | 2020-08-11  | UC2-0sEOYbQFuaURd_AU6Krg | OneFootball English | 17         | 2020-08-12    |
| 12 | cSIrKfmgR1s | Sensible Transfers: Arsenal  | 2020-08-12  | UCGYYNGmyhZ_kwBF_lqqXdAQ | Tifo Football       | 17         | 2020-08-12    |
| 13 | w-hVKXmib1c | I've Seen Something KSI Hasn't Shown You                           | 2020-08-10  | UCHhfSXoDG6gSgpOvLH4wrRw | Behzinga            | 20         | 2020-08-12    |
| 14 | xTpr7piQu2M | FUT 21   Official Trailer  | 2020-08-10  | UCovaxd5LQSuP4ChkxK0pnZQ | EA SPORTS FIFA      | 20         | 2020-08-12    |

#### 4. INGESTING JSON DATA:

### 4.1 CREATINHG AN EXTERNAL TABLE:

An 'ex\_table\_youtube\_category' external table is created to store and query data directly from Json files located on the azure blob storage specified by the 'stage\_assignment' stage. We mentioned file format as Json. It specifies the data files are in Json format. The PATTERN clause uses a regular expression to filter and match only specific files in the storage location. The pattern'. \_category\_id[.]Json' targets all JSON files whose filenames end with \_category\_id. json.

```
CREATE OR REPLACE EXTERNAL TABLE ex_table_youtube_category
WITH LOCATION = @stage_assignment
FILE_FORMAT = (TYPE=JSON)|
PATTERN = '.*_category_id[.]json';
```

```
VALUE

1 { "etag": "kBCr3l9kLHHU79W4lp5196LDptl", "items": [ { "etag": "lfNa37JGcqZs-jZeAyFGkbeh6bc", "id": "1", "kind": "youtube#videoCategory", "snippet": { "assignable": true, "channelld": "UCBR8-f
```

### **4.2 CREATING THE INTERAL TABLE BY USING LATERAL FLATTEN:**

By extracting and converting JSON data from an external table called ex\_table\_youtube\_category, the SQL statement produces a table called table\_youtube\_category.

Split\_part extracts the country code from the file name from the file and allow us to include in the table.

The lateral flatten function is used to expand the nested json arrays. It enables each item within the array to become a separate row in thwe output. This function is crucial for handling json data structures, allowing you to extract fields like 'id' and 'snippet:title' for each category in normalized format.

```
CREATE OR REPLACE TABLE table_youtube_category AS
SELECT
   split_part(metadata$filename,'_',1)as country,
   l.value:id::string as CATEGORYID,
   l.value:snippet:title::string as CATEGORY_TITLE
FROM ex_table_youtube_category, LATERAL FLATTEN(value:items) l;
```

|    | COUNTRY | CATEGORYID | CATEGORY_TITLE   |
|----|---------|------------|------------------|
| 1  | KR      | 1          | Film & Animation |
| 2  | KR      | 2          | Autos & Vehicles |
| 3  | KR      | 10         | Music            |
| 4  | KR      | 15         | Pets & Animals   |
| 5  | KR      | 17         | Sports           |
| 6  | KR      | 18         | Short Movies     |
| 7  | KR      | 19         | Travel & Events  |
| 8  | KR      | 20         | Gaming           |
| 9  | KR      | 21         | Videoblogging    |
| 10 | KR      | 22         | People & Blogs   |
| 11 | KR      | 23         | Comedy           |
| 12 | KR      | 24         | Entertainment    |
| 13 | KR      | 25         | News & Politics  |
| 14 | KR      | 26         | Howto & Stvle    |

### **5.COMBINING DATA:**

The final internal table table\_youtube\_final was created by merging table\_youtube\_trending with table\_youtube\_category using a LEFT JOIN. A unique identifier was added to each record using UUID\_STRING(). It generates generating a distinct UUID for each row

```
CREATE OR REPLACE TABLE table_youtube_final AS
SELECT

UUID_STRING() AS id, --UUID_STRING() function is executed for each row in the table, generating a distinct UUID for each row
t.Videoid AS VIDEO_ID,
t.TITLE AS TITLE,
t.PUBLISHEDAT AS PUBLISHEDATE,
t.CHANNELID AS CHANNEL_ID,
t.CHANNELID AS CHANNEL_TITLE,
t.CATEGORY_TITLE AS CHANNEL_TITLE,
t.CATEGORY_TITLE AS CATEGORY_ID,
c.CATEGORY_TITLE AS CATEGORY_TITLE,
t.TRENDING_DATE AS TRENDING_DATE,
t.VIEW_COUNT AS VIEW_COUNT,
t.LIKES AS LIKES,
t.DISLIKES AS DISLIKES,
t.COMMENT_COUNT AS COMMENT_COUNT,
t.COUNTRY AS COUNTRY
FROM table_youtube_trending t

LEFT JOIN table_youtube_trending t

LEFT JOIN table_youtube_final;

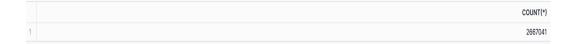
--ROW COUNT

SELECT *FROM table_youtube_final;

FROM table_youtube_final;
```

|    | ID                                   | VIDEO_ID    | TITLE   | PUBLISHEDATE | CHANNEL_ID               | CHAI  |
|----|--------------------------------------|-------------|---|--------------|--------------------------|-------|
| 1  | 0b87e4dd-c8d6-4f3d-b1fe-17a2fb5bc47d | 7rlwxSPUcQk | ON EST POSITIF AU COVID-19 (coronavirus)                          | 2020-08-11   | UCpWaR3gNAQGsX48cllQC0qw | Tibo  |
| 2  | 60a1e41f-dfc8-46c8-a064-3bd2ee894b1b | AcBd_RH9JSw | PASSER UNE NUIT DANS LA PISCINE DE LA VILLA AVEC @Michou!         | 2020-08-11   | UCUI7mwOyySfZzUkq4H29nug | LeBo  |
| 3  | 0f7c3e9a-2766-452c-bbc8-694e2537a268 | JVm8P6kKgD0 | FRANGLISH - My Salsa feat. Tory Lanez (clip officiel)             | 2020-08-12   | UCnwfc0O-LGEg8y9bEQTaSgQ | Fran  |
| 4  | c024c526-a472-40a0-a26e-8a4619798fd2 | JfgeD7xWy-w | L'ÉTÉ LE PLUS ECLATÉ DE MA LIFE                                   | 2020-08-11   | UCMj2VQ3-8zpyeAl7dU0T-Bg | Fahd  |
| 5  | 1b807a99-f410-4d78-92c4-ab1455ed67ad | P3lkBxra3a8 | JE L'ESSAYE ENFIN ! DUCATI HYPERMOTARD                            | 2020-08-12   | UC-uy7_1hColMpQ_2rs-kY6w | KIKA  |
| 6  | ea3acc67-595b-4ae9-83bf-67943498d4f8 | AwdfloSJrtA | MÉLI-MÉLO le 3.   | 2020-08-09   | UCDPK_MTu3uTUFJXRVcTJcEw | Mcfly |
| 7  | 77ea8e4a-ce8d-4b23-a9b8-6c784b8c51fe | J9iT2z0-ITc | Inoxtag attend Léa devant la Villa Il n'a toujours pas compris #7 | 2020-08-10   | UCNGq4mP3Ds50UGjPo8IJOcw | Valou |
| 8  | b589d8d2-bf05-4bce-8875-d7c3ee9d6efa | zatMgEXo4Ko | LDC - MES PRONOS DU FINAL 8 !                                     | 2020-08-12   | UCopwMeeGavRDZbo3EFBTNUg | Bruce |
| 9  | af2a93a9-d05a-46df-b10f-4960b09a3256 | 56HGTFZabxA | Inoxtag - Funkinox ft. Kazzey (Clip officiel)                     | 2020-08-09   | UCL9aTJb0ur4sovxcppAopEw | Inoxt |
| 10 | 408b9608-a07d-457c-90df-3fe5caba190a | n0rGxQw5lcs | MICHOU - DANS LE CLUB (PARODIE FORTNITE) SAISON 3 CHAPITRE        | 2020-08-12   | UChGhSbkm9jS_u1thVCeGpiw | Kebo  |
| 11 | 134c625d-297d-4f75-b274-3be86cea9b60 | DTs12U36×54 | VOS SOIRÉES QUI ONT FAILLI (TRÈS) MAL FINIR ! (Feat Amixem)       | 2020-08-09   | UCow2lGnug1l3Xazkrc5jM_Q | JOYO  |
| 12 | d723cc9b-95d3-4b46-b613-1bafd913499d | jbGRowa5tlk | ITZY "Not Shy" M/V TEASER   | 2020-08-11   | UCaO6TYtlC8U5ttz62hTrZgg | JYP E |
| 13 | 58934716-0baf-430b-b590-a86a050e2f23 | emuvKXZo-f4 | MON HOUSE TOUR! ( Je vous présente ma maison )                    | 2020-08-11   | UCTHkQBKxQsytdNccmKl8pGw | Lytha |
| 14 | 6c23fa3c-7dc7-4842-ac7b-52f9c766b95b | Fi4CK7VniUo | le meilleur anniversaire de toute sa vie                          | 2020-08-11   | UCxDsG8fSVisYiTE_ga6YJ7Q | Esile |

The row count of table\_youtube\_final was verified to ensure completeness, resulting in 2,667,041 rows.



### **6.DATA CLEANING:**

This report describes the procedures and queries that were used in Snowflake's 'table\_youtube\_final' and 'table\_youtube\_category' tables to organise, clean, and deduplicate data. Ensuring data consistency and integrity was the goal because these are essential for precise analysis of trends in YouTube videos.

# 1. In table\_youtube\_category, Which category\_title Has Duplicates if We Don't Take Into Account the categoryid(Return Only a Single Row)?

This query identifies category\_title entries that have multiple categoryid values, ignoring the category id's uniqueness. By grouping by CATEGORY\_TITLE and using HAVING COUNT, this query returns any title associated with more than one category\_id, highlighting potential inconsistencies in how categories are labeled across different datasets or countries.

|   | CATEGORY_TITLE |  |
|---|----------------|--|
| 1 | Comedy         |  |

### 2. In table\_youtube\_category, Which category\_title Only Appears in One Country?

```
-- 2: In "table_youtube_category" which category_title only
appears in one country?

SELECT
        CATEGORY_TITLE
FROM
        table_youtube_category
GROUP BY
        CATEGORY_TITLE
HAVING
        COUNT(DISTINCT COUNTRY) = 1;
```

This query aims to identify category\_title entries that are unique to a single country. By grouping by CATEGORY\_TITLE and filtering with HAVING COUNT(DISTINCT COUNTRY) = 1, it pinpoints titles that appear exclusively within one country's dataset. This information is useful for understanding regional specificity in content categorization.

```
CATEGORY_TITLE

1 Nonprofits & Activism
```

## 3. In table\_youtube\_final, what is the categoryid of the Missing category\_titles?

```
-- 3: In "table_youtube_final", what is the categoryid of the missing category_titles?

SELECT CATEGORY_ID

FROM table_youtube_final
```

This query selects the CATEGORY\_ID from records where CATEGORY\_TITLE is NULL in table\_youtube\_final. Identifying these entries is crucial for data completeness, as missing category\_titles may lead to gaps in data analysis or misinterpretation of trends.

|    | CATEGORY_ID |
|----|-------------|
| 1  | 29          |
| 2  | 29          |
| 3  | 29          |
| 4  | 29          |
| 5  | 29          |
| 6  | 29          |
| 7  | 29          |
| 8  | 29          |
| 9  | 29          |
| 10 | 29          |
| 11 | 29          |
| 12 | 29          |
| 13 | 29          |
|    |             |

### 4.Update the table\_youtube\_final to Replace the NULL Values in category\_title with the Answer from the Previous Question

```
--4.Update the table_youtube_final to replace the NULL values
in category_title with the answer from the previous question.

UPDATE table_youtube_final
SET category_title = (
    SELECT CATEGORY_TITLE FROM table_youtube_category
    GROUP BY CATEGORY_TITLE
    HAVING COUNT(DISTINCT COUNTRY)=1
)
WHERE category_title IS NULL;
```

This update operation replaces NULL values in category\_title in the table\_youtube\_final table. The subquery identifies unique CATEGORY\_TITLE entries that appear in only one country from table\_youtube\_category. Filling these gaps ensures that all records are fully populated, maintaining the dataset's integrity and completeness.

|   | number of rows updated | numb |
|---|------------------------|------|
| 1 | 1563                   |      |
|   |                        |      |

### 5. In table\_youtube\_final, Which Video Doesn't Have a channeltitle (Return Only the Title)?

```
--5: In "table_youtube_final", which video doesn't have a channeltitle (return only the title)?

SELECT TITLE
FROM table_youtube_final
WHERE CHANNEL_TITLE IS NULL;
```

This query identifies videos that lack a CHANNEL\_TITLE, returning only their titles. Detecting these gaps is essential for data completeness and ensuring every video is properly attributed to a channel, which is important for analysing channel-specific trends.

```
TITLE

1 Kala Official Teaser | Tovino Thomas | Rohith V S | Juvis Productions | Adventure Company
```

### 6. Delete from table\_youtube\_final, Any Record with video\_id = "#NAME?"

```
--6:Delete from "table_youtube_final", any record with video_id = "#NAME?"

DELETE FROM table_youtube_final
WHERE video_id = '#NAME?';
```

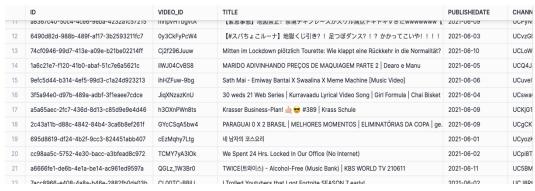
This delete statement removes records with video\_id values set to #NAME?. These entries likely result from data import errors or placeholders and can skew analysis if left unaddressed. Cleaning these records ensures the dataset's accuracy.

```
number of rows deleted

32081
```

### 7. Create a New Table Called table\_youtube\_duplicates Containing Only the "Bad" Duplicates Using the ROW\_NUMBER() Function

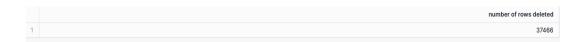
from table\_youtube\_final. The ROW\_NUMBER() function is used to rank each record within partitions defined by video\_id, country, and trending\_date. The ORDER BY view\_count DESC clause ensures that the highest view count gets a row\_num of 1, marking others as duplicates. Rows with row\_num > 1 are considered duplicates and stored in the new table.



8. Delete the Duplicates in table\_youtube\_final Using table youtube duplicates

```
--8 Delete records from table_youtube_final that are found in table_youtube_duplicates
DELETE FROM table_youtube_final
WHERE id IN (
    SELECT id
    FROM table_youtube_duplicates);
```

This delete statement removes duplicate entries from table\_youtube\_final by matching records against table\_youtube\_duplicates. Using the id field ensures that only unwanted duplicates are deleted while retaining the most accurate records for each unique combination of video\_id, country, and trending\_date.

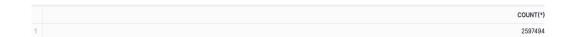


9. Count the Number of Rows in table\_youtube\_final and Check That It Is Equal to 2,597,494 Rows

```
--9 .Count the number of rows in "table_youtube_final" and check that it is equal to 2,597,494 rows. rows.

SELECT COUNT(*)
FROM table_youtube_final;
```

This final check ensures the table\_youtube\_final table contains the expected 2,597,494 rows after the cleaning process. Verifying the row count confirms that all duplicates have been correctly managed and that no unintended deletions occurred, ensuring the dataset's integrity and readiness for further analysis.



#### 7. DATA ANALYSIS:

1. What are the 3 Most Viewed Videos for Each Country in the "Gaming" Category for trending\_date = '2024-04-01'?

```
--1. What are the 3 most viewed videos for each country in the
"Gaming" category for the trending_date = ''2024-04-01''.
Order the result by country and the rank

SELECT country, title, channel_title, view_count, RANK() OVER
(PARTITION BY country ORDER BY view_count DESC) AS RK
FROM
    table_youtube_final
WHERE
    category_title = 'Gaming'
    AND trending_date = '2024-04-01'

QUALIFY
    RK <= 3
ORDER BY
    country,
    RK;</pre>
```

The query identifies the top three most viewed videos in the "Gaming" category for each country on April 1, 2024. It uses the RANK() window function partitioned by country and ordered by view\_count in descending order to assign a rank to each video. The QUALIFY clause ensures only the top three videos (ranked 1 to 3) for each country are included in the results. This helps highlight the most popular gaming content across different regions on a specific date.

### 2. For Each Country, Count the Number of Distinct Videos with a Title Containing the Word "BTS"

```
--2.For each country, count the number of distinct video with a title containing the word "BTS" and order the result by count in a descending order

SELECT country, COUNT(DISTINCT VIDEO_ID) AS Number_Count FROM table_youtube_final
WHERE title LIKE '%BTS%'|
GROUP BY country
ORDER BY Number_Count DESC;
```

This query counts the number of unique videos containing the word "BTS" (case insensitive) in their title for each country. The COUNT(DISTINCT VIDEO\_ID) function ensures that each video is counted only once, even if it appears multiple times. The results are then ordered in descending order based on the count, highlighting countries with the most "BTS"-related content.

# 3. For Each Country, Year, and Month (Combined into a Single Column), Identify the Most Viewed Video and Its likes\_ratio

```
--3. For each country, year and month (in a single column), which video is the most viewed and what is its likes_ratio (defined as the percentage of likes against view_count) truncated to 2 decimals. Order the result by year_month and country.
WITH VideoStats AS (
          country
          EXTRACT(YEAR FROM trending_date) AS year,
          EXTRACT(MONTH FROM trending_date) AS month,
          video_id,
          title,
channel_title,
category_title,
          view_count.
          likes,
dislikes,
               WHEN view_count = 0 THEN 0
               ELSE ROUND((likes::NUMERIC / view_count) * 100, 2)
          END AS likes_ratio,
          RANK() OVER [PARTITION BY country, EXTRACT(YEAR FROM trending_date), EXTRACT(MONTH FROM trending_date) ORDER BY view_count DESC) AS
          table_youtube_final
SELECT
     TO_DATE(year || '-' || TO_CHAR(month, 'FM00') || '-01', 'YYYY-MM-DD') AS year_month,
     title, channel_title AS CHANNELTITLE,
     category_title,
view_count,
     likes_ratio
FROM
     VideoStats
WHERE
  video_rank = 1
  AND year = 2024
ORDER BY
  year_month,
  country;
```

The query calculates the most viewed video and its likes\_ratio (likes as a percentage of view count) for each country, year, and month for 2024. Using a WITH clause, it extracts the year and month from trending\_date, computes the likes\_ratio, and ranks videos by view count. The final selection filters for the top-ranked (most viewed) video for each month and country, providing insight into top-performing videos and their viewer engagement.

|    | COUNTRY | YEAR_MONTH | TITLE  | CHANNELTITLE      | CATEGORY_TITLE | VIEW_COUNT | LIKES_RATIO |
|----|---------|------------|--|-------------------|----------------|------------|-------------|
| 1  | BR      | 2024-01-01 | Survive 100 Days Trapped, Win \$500,000  | MrBeast           | Entertainment  | 139504939  | 3.20        |
| 2  | CA      | 2024-01-01 | Still Here   Season 2024 Cinematic - League of Legends (ft. Forts, Tiffany Aris, and 2WE | League of Legends | Gaming         | 104159411  | 1.69        |
| 3  | DE      | 2024-01-01 | Still Here   Season 2024 Cinematic - League of Legends (ft. Forts, Tiffany Aris, and 2WE | League of Legends | Gaming         | 104159411  | 1.69        |
| 4  | FR      | 2024-01-01 | Still Here   Season 2024 Cinematic - League of Legends (ft. Forts, Tiffany Aris, and 2WE | League of Legends | Gaming         | 104159411  | 1.69        |
| 5  | GB      | 2024-01-01 | Still Here   Season 2024 Cinematic - League of Legends (ft. Forts, Tiffany Aris, and 2WE | League of Legends | Gaming         | 104159411  | 1.69        |
| 6  | IN      | 2024-01-01 | Protect \$500,000 Keep It!   | MrBeast           | Entertainment  | 85458562   | 4.21        |
| 7  | JP      | 2024-01-01 | Survive 100 Days Trapped, Win \$500,000  | MrBeast           | Entertainment  | 137639799  | 3.22        |
| 8  | KR      | 2024-01-01 | Survive 100 Days Trapped, Win \$500,000  | MrBeast           | Entertainment  | 143955997  | 3.16        |
| 9  | MX      | 2024-01-01 | Survive 100 Days Trapped, Win \$500,000  | MrBeast           | Entertainment  | 137639799  | 3.22        |
| 10 | US      | 2024-01-01 | Grand Theft Auto VI Trailer 1  | Rockstar Games    | Gaming         | 166323421  | 6.73        |
| 11 | BR      | 2024-02-01 | Face Your Biggest Fear To Win \$800,000  | MrBeast           | Entertainment  | 126846652  | 3.54        |
| 12 | CA      | 2024-02-01 | Face Your Biggest Fear To Win \$800,000  | MrBeast           | Entertainment  | 119170728  | 3.66        |
| 13 | DE      | 2024-02-01 | Face Your Biggest Fear To Win \$800,000  | MrBeast           | Entertainment  | 114978689  | 3.72        |
|    | 50      |            | 5 W 5 V 5 T W 400000   |                   |                | 44 4070000 |             |

### 4. For Each Country, Which category\_title Has the Most Distinct Videos Before 2022?

```
-- 4.For each country, which category_title has the most distinct videos and what is its percentage (12 decimals) out of the total distinct number of videos of that country? Only look at the data from before 2022. Order the result by category_title and country
 WITH VideoCategories AS (
      SELECT
            country,
           category_title,
COUNT(DISTINCT video_id) AS total_category_video
            table_youtube_final
      EXTRACT(YEAR FROM trending_date) >= 2022
GROUP BY
            country.
 TotalCountryVideos AS (
            country,
COUNT(DISTINCT video_id) AS total_country_video
      FROM table_youtube_final
      EXTRACT(YEAR FROM trending_date) >= 2022
GROUP BY
           country
    vc.country AS COUNTRY,
vc.category_title AS CATEGORY_TITLE,
    vc.total_category_video AS TOTAL_CATEGORY_VIDEO,
tc.total_country_video AS TOTAL_COUNTRY_VIDEO,
ROUND((vc.total_category_video::NUMERIC / tc.total_country_video) * 100, 2) AS PERCENTAGE
FROM
    VideoCategories vc
    TotalCountrvVideos to
     vc.country = tc.country
     (vc.country, vc.total_category_video) IN (
              country
              MAX(total_category_video)
         FROM VideoCategories
         GROUP BY
ORDER BY
    vc.country;
. Which channeltitle has produced the most distinct videos and what is this number?
WITH ChannelVideoCounts AS (
SELECT
         COUNT(DISTINCT video_id) AS distinct_video_count
    FROM table_youtube_final
    GROUP BY
         channel_title
     channel_title,
     distinct_video_count
      distinct_video_count = (SELECT MAX(distinct_video_count) FROM ChannelVideoCounts);
```

This query identifies the category\_title with the most distinct videos for each country using data from before 2022. It calculates the number of distinct videos per category and country, then computes the percentage of each category's videos out of the total videos for that country. This percentage is truncated to 12 decimal places. The results help understand content diversity and dominant categories in different countries.

|    | COUNTRY | CATEGORY_TITLE | TOTAL_CATEGORY_VIDEO | TOTAL_COUNTRY_VIDEO | PERCENTAGE |
|----|---------|----------------|----------------------|---------------------|------------|
| 1  | BR      | Entertainment  | 5417                 | 23760               | 22.80      |
| 2  | DE      | Entertainment  | 7709                 | 30719               | 25.10      |
| 3  | FR      | Entertainment  | 7548                 | 32849               | 22.98      |
| 4  | GB      | Entertainment  | 5643                 | 27855               | 20.26      |
| 5  | IN      | Entertainment  | 21281                | 50250               | 42.35      |
| 6  | JP      | Entertainment  | 5658                 | 17627               | 32.10      |
| 7  | KR      | Entertainment  | 5122                 | 15175               | 33.75      |
| 8  | MX      | Entertainment  | 4195                 | 17532               | 23.93      |
| 9  | CA      | Gaming         | 6594                 | 30869               | 21.36      |
| 10 | US      | Gaming         | 6226                 | 28799               | 21.62      |

### 5. Which channel\_title Has Produced the Most Distinct Videos?

The query finds the channel\_title that has produced the most distinct videos across the entire dataset. It counts the number of unique videos per channel and then selects the channel with the maximum count. This analysis helps identify the most prolific content creators on YouTube, indicating channels with broad content production.



### **8. BUSINESS QUESTION:**

The query first creates a common table expression (CTE) named VideoStats to rank videos by view\_count for each country, year, and month.

It excludes videos in the "Music" and '

```
-- THE GOAL IS TO FIND OUT WHICH CATEGORIES FREQUENTLY HAVE TRENDING VIDEOS. THEN IT WILL EASY FOR ME TO FIGURE OUT WHICH CATEGORY VIDEO I CAN MAKE
THAT WILL APPEAR ON TRENDING LIST OT NOT
WITH VideoStats AS (

SELECT

country,

category_title,

RANK() OVER (PARTITION BY country, EXTRACT(YEAR FROM trending_date), EXTRACT(MONTH FROM trending_date) ORDER BY view_count DESC) AS
TRENDING_CT_RK

FROM

table_youtube_final
WHERE NOT

category_title IN ('Music', 'Entertainment') -- (EXCLUDING MUSIC AND ENTERTAINMENT)
)

SELECT

country,

category_title,

COUNT(*) AS TRENDING_CT_RK

FROM

VideoStats
WHERE

TRENDING_CT_RK = 1
GROUP BY

country,

category_title

ORDER BY

country,

TRENDING_CT_RK DESC;

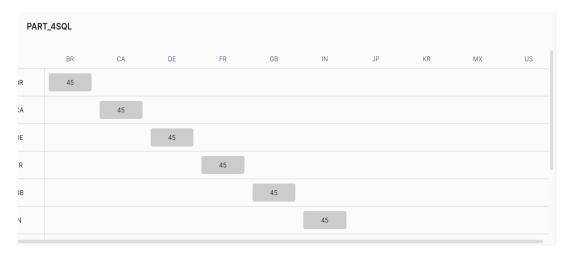
-- ACCORTING TO THE DATA, I WILL CREATE A YOUTUBE CANNNEL AND MAKE A VIDEO IN GAMING CATEGORY. IT HAS A HIGH CHANCE OF APPEAR IN TOP TREND IN YOUTUBE.
```

Entertainment" categories to focus on other content types. The RANK() window function is applied to assign a rank to each video based on its view count within its respective country and time frame.

In the main query, only videos ranked first (TRENDING\_CT\_RK = 1) are considered, representing the top trending videos for each country and period. The results are then grouped by country and category\_title, and ordered by the number of times a category has reached the top trending spot (TRENDING CT RK).

### **Insights and Strategy**

Based on the analysis, the "Gaming" category frequently appears in the top trending videos in the US, making it a strong candidate for a new YouTube channel. This suggests that focusing on creating gaming content could increase the likelihood of trending. However, this strategy may not work universally in every country, as cultural differences and content preferences vary. For example, other countries might have different popular categories, such as "Sports" or "News." Therefore, while "Gaming" is a strong contender in the US, localized content strategies should be developed for success in other regions.



|    | COUNTRY | CATEGORY_TITLE   | TRENDING_CT_RK |
|----|---------|--|----------------|
| 1  | BR      | Gaming   | 12             |
| 2  | BR      | Science & Technology   | 10             |
| 3  | BR      | People & Blogs   | 8              |
| 4  | BR      | Sports   | 8              |
| 5  | BR      | Film & Animation   | 2              |
| 6  | BR      | Comedy   | 1              |
| 7  | BR      | Autos & Vehicles   | 1              |
| 8  | BR      | Education  | 1              |
| 9  | BR      | Travel & Events  | 1              |
| 10 | BR      | Howto & Style  | 1              |
| 11 | CA      | Science & Technology   | 10             |
| 12 | CA      | Gaming   | 7              |
| 13 | CA      | People & Blogs   | 7              |
|    | ^*      | CO. A. A. S. | ^              |

### **8.CONCLUION:**

In conclusion, this project effectively utilized Snowflake's Data Lakehouse architecture to analyze YouTube trending and category data across various countries, providing a comprehensive understanding of global content consumption patterns. By carefully ingesting, transforming, and cleaning the data, we ensured its integrity and accuracy, enabling robust analysis. The insights gained from this project reveal significant regional differences in content preferences and highlight the potential for optimizing content strategies on YouTube. Focusing on high-performing categories, such as Gaming or Sports (excluding Music and Entertainment), can increase the likelihood of creating trending videos. Overall, this analysis underscores the importance of data-driven decision-making in digital content creation and strategy development.