Youtube Analytics:

Data:

## Trending YouTube Video Statistics

Data Card    Code (2676)    Discussion (43)    ▲ 5116    New Notebook    ⬇ Download (211 MB)    ⬤    ⋮

The Dataset contains information of trending videos on the most popular video streaming services, YouTube. The data was collected through YouTube API and contains regional information about trending videos fōrom around 2017-18. There are 10 regions in total where data was fetched from. For each region, we have two files- one JSON and CSV.

{i} CA_category_id.json
▥ CAvideos.csv
{i} DE_category_id.json
▥ DEvideos.csv
{i} FR_category_id.json
▥ FRvideos.csv
{i} GB_category_id.json
▥ GBvideos.csv
{i} IN_category_id.json
▥ INvideos.csv
{i} JP_category_id.json
▥ JPvideos.csv
{i} KR_category_id.json
▥ KRvideos.csv
{i} MX_category_id.json
▥ MXvideos.csv
{i} RU_category_id.json
▥ RUvideos.csv
{i} US_category_id.json

Because of the difference in file format/structure, I have to combine the csv and json files for each region to make a final analysis model.

Approach:
For this I utilized AWS Cloud Services. AWS offers several advantages over on premises infrastructure including Scalability, Cost Effectiveness, Serverless architecture, Services made only for Data Analytics, BDA.

For my data storage, I needed tools to save structured and semi-structured data which is why I thought of S3 Data Lake formation. S3-based data lakes offer cost advantages compared to traditional data warehouses.
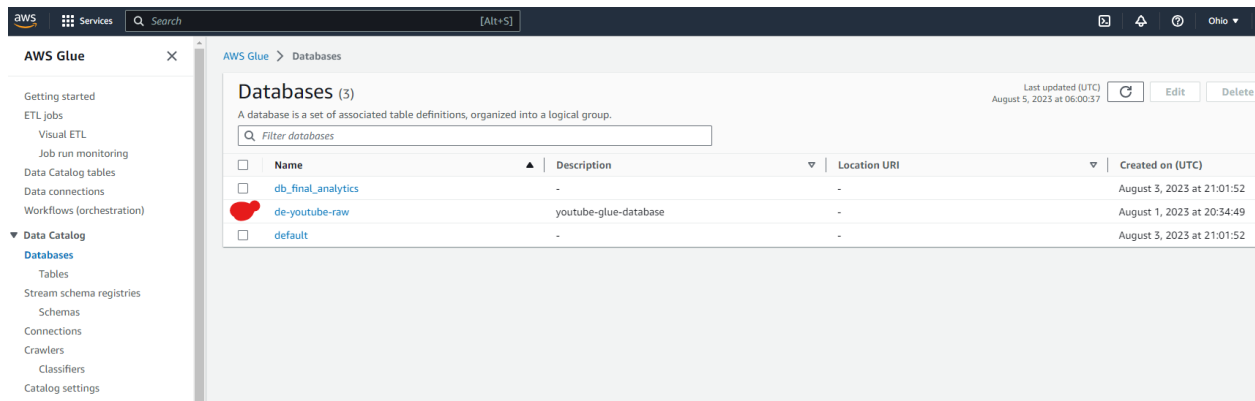
Data lakes are well-suited for storing and managing diverse data types, including structured, semi-structured (like JSON), and unstructured data (like CSV). Since the YouTube dataset contains both JSON and CSV files, a data lake provides the flexibility to store and process this mixed data without the need to convert it to a rigid schema upfront.

# Pipeline 1:
Processing JSON files:

After downloading files from Kaggle, I uploaded all json files to my S3 raw data bucket through AWS CLI. Before doing this, I created an IAM user with my root account and created an S3 raw Bucket in it which will contain my raw data (JSON in this case).

My data is now uploaded to S3, to understand the schema and do further analysis in Athena, I need to create a Glue crawler that will crawl through my data and store metadata in the Glue catalog.



Here is how the schema looked after crawling through my raw s3:



**Roles & Permissions: I gave my glue crawler an IAM role that permits it to access S3 and Glue Service Role.**

Now this catalog can be referred to by Athena for querying data in s3. But when I ran Athena to acces contents of s3, I received an error which required me to convert my multi line json format data to Parquet format, a format which makes Athena process queries efficiently.

For this a ETL pipeline is made with following architecture:

Setting up Lambda Function:



**Roles: I gave my Lambda Function same IAM role that permits it to access S3 and Glue .**

Here Lambda function needs to be triggered whenever an object is put on s3 raw. The code in run time environment will be executed which will convert the multi line JSON into Parquet format with help of environment variables (info about destination bucket, catalog and table).
Additionally, the lambda function will write to the glue data catalog which stores schema of  json converted parquet data in s3. I used Pandas and AWS Wrangler library in python for reading, converting to parquet and loading the file in the target destination. For this I had to upload a zip file of these libraries to set up the lambda layer. The transformed data will be stored in a new s3 named S3 Cleaned.

Lambda Trigger: I created an S3 trigger that triggers my lambda function whenever an object is inserted in the bucket. This will remove manual involvement and will automatically trigger the function, convert the files and store in the new S3 (also create glue catalog).

Using CloudWatch logs to solve trigger problem for lambda:
Problems faced: I couldn't invoke my lambda function whenever objects were inserted. So I created cloud watch event to check status of my function when invoked. It turned out that CloudWatch event were not triggered when new objects were inserted, this was due to incorrect bucket prefix I configured my lambda function with. To solve this I simply removed the prefix preference from trigger.

Now my S3 cleaned bucket contains converted parquet files which can be accessed by Athena. I'm now able to access files in S3 cleaned buckets with help of the catalog that I got from lambda. Here is how it looks after I ran query on S3 cleaned Bucket:



## Pipeline 2:

Processing CSV data:
Now that my JSON file is converted, I also need to do some transformation in my csv files: change data type for columns like likes, comments, views etc. In the process I will also try to convert csv to parquet so that my final bucket contains both json converted parquet files from pipeline 1 and csv to parquet files from pipeline 2.

When uploading csv files to my S3 raw, I want information about the schema of files. So I set up a crawler that will store catalog in Glue Catalog that will crawl through csv files in S3 raw.

This data catalog, combined with the one in pipeline 1 should give Athena enough metadata to combine the csv files and the json turned parquet files from pipeline 1. But after running a join query between two schemas, I got an error stating that the data type for the join was not

compatible. Therefore I need to change the datatype for numeric columns which have string data type instead. I will also try to convert csv files to parquet for easy processing with Athena.

<u>Setting up Glue ETL Job:</u>

Next I needed to set up a ETL job which can be done with AWS Glue. By giving source bucket, IAM role, Transformation via visual job creator (converting data type), choosing destination target bucket where format will be Parquet format.

Here in the Glue ETL script I modified the automatically generated code so that the destination files are categorized with regions and then stored in the destination S3 cleaned bucket.

After doing the above changes, I was receiving an error because I couldn't convert the non english words in title columns from non english native regions.

I utilized Bucket versioning in raw S3 bucket and soft deleted every files except ca(Canada),us (United States) and gb (Great Britain) which are english dominated region.



Now after running this job successfully, my S3 cleaned bucket contains both json converted parquet from pipeline 1 and csv converted parquet with necessary data type transformations. To get a schema for this csv converted parquet files, I need to set up a crawler which will store schema of these files in data catalog.

## Pipeline 3:

Combining both parquet files and performing data transformation to categorize data with region then category_id.

Now my S3 cleaned data contains both parquet files. I will need to combine these both and store in new S3 final bucket which categorizes files based on region and category_id. I will also include "add to catalog option" that will automatically create a data schema for my Final S3 bucket.

I will again use Glue ETL job to perform extract transform loading of data.
The visual from Glue ETL looks like this:

# de-harshits-youtube-final-analytics

Source ▾    Action ▾    Target ▾    Undo    Redo    Remove    ⊕    ⊖    ⊡

**Data source - Data Catalog**
AWS Glue Data Catalog

**Data source - Data Catalog**
AWS Glue Data Catalog

**Transform - Join**
Join

**Data target - S3 bucket**
Amazon S3

There are two catalogs one for each parquet files (csv and json converted). The transformation includes joining both schemas on category_id and id.

The data inserted in my New Final S3 bucket will look like this:

# region=ca/

**Objects** | Properties

## Objects (16)

Objects are the fundamental entities stored in Amazon S3. You can use **Amazon S3 inventory** [↗] to get a list of all objects in your bucket. For others

| C | 🗗 Copy S3 URI | 🗗 Copy URL | ⬇ Download | Open ↗ | Delete | Actions ▼ |

🔍 Find objects by prefix

| | Name ▲ | Type ▽ | Last modified |
|---|---|---|---|
| ☐ | 📁 category_id=1/ | Folder | - |
| ☐ | 📁 category_id=10/ | Folder | - |
| ☐ | 📁 category_id=15/ | Folder | - |
| ☐ | 📁 category_id=17/ | Folder | - |
| ☐ | 📁 category_id=19/ | Folder | - |
| ☐ | 📁 category_id=2/ | Folder | - |
| ☐ | 📁 category_id=20/ | Folder | - |
| ☐ | 📁 category_id=22/ | Folder | - |
| ☐ | 📁 category_id=23/ | Folder | - |

Visualizing the data in Final S3 Bucket:
Using AWS Quicksight, I built a dashboard which helps deliver business insights and give visual representation. I connect Athena as my data source and now Quicksight can automatically query in Athena which has catalog on my final youtube analytics data.

Dashboard:

**Views by Categories**

Music
Entertainment
People & Blogs
Science & Tech...
Film & Animation
Comedy
Comedy 6,458 (0%)
Sports
Education

**Total Comments**

282,617,342

**Total Likes**

2,685,842,999

**Views by Region**

ca
gb
us
10,242,91...

**Top 10 Liked Youtube Channels**

SHOWING TOP 10 IN CHANNEL_TITLE

600M
400M
200M
0

PewDiePie · xxxtentacion · SpaceX · Primitive T... · CaseyNeistat · Vevo · zefrank1 · bill wurtz · MAMAMOO · Simone Giertz

**Likes by Categories**

Film & Animation
Gaming
Howto & Style
News & Politics
Sports
Pets & Animals
Travel & Events
Education
Autos & Vehicles
Nonprofits & Act...
Shows

0   25M   50M   75M   100M   125M   150M

**Trend by Categories**

● Comedy

6M
4M
2M
0

**Tags Word Cloud**

SHOWING TOP 50 IN TAGS

Dota Dota 2 ESL ONE Katowice VP MVP
Louise Pentland
Linkin Park
Gwen Stefani Blake Shelton You Make It F...
Sia Snowman Holiday
CBS 2 News Morning
the simpsons
Emirates first class
klpolish
SpaceX Heavy Landing
markiplier
Trap Adventure 2
Jordan Peterson
SHANtell martin
Bruno Mars and Cardi B Finesse R&B...
SATIRE
Oblivion
Sia Ho Ho Holiday
Lonzo Ball
viralhog
Imagine Dragons Thunder Evolve Ellen DeG...
Sia Candy Cane Lane Holiday
Frasier Crane (Fictional Character)
11-8-17
bill wurtz
Trixie Mattel Rupaul All Stars One Stone...
Beach House Lemon Glow
Sia Santa's Coming For Us Holiday
Star Wars The Last Jedi Plot Holes
Conan O'Brien Conan Conan (TV Series) TR...

**Categories-Region analysis**

| Categories | Region | | | |
|---|---|---|---|---|
| | us | gb | ca | Total |
| | Views | Views | Views | Views |
| ⊞ Autos & ... | 1,844,590 | | 41,705,590 | 43,550,180 |
| ⊞ Comedy | 2,208,997,330 | 2,097,358,080 | 2,121,049,270 | 6,427,404,680 |
| ⊞ Education | 35,557,900 | 5,552,000 | 31,189,820 | 72,299,720 |
| ⊞ Entertainment | 7,859,609,290 | 7,155,895,470 | 4,503,943,060 | 19,519,447,820 |
| ⊞ Film & ... | 2,586,576,850 | 3,825,907,600 | 1,096,942,650 | 7,509,427,100 |
| ⊞ Gaming | 324,885,230 | 620,955,890 | 467,492,430 | 1,413,333,550 |
| ⊞ Howto & Style | 307,259,030 | 76,260,430 | 85,252,740 | 468,772,200 |
| ⊞ Music | 4,877,982,510 | 23,339,602,030 | 2,209,526,320 | 30,427,110,860 |
| ⊞ News & ... | 405,864,100 | 227,468,070 | 579,222,680 | 1,212,554,850 |
| ⊞ Nonprofits ... | 169,118 | 831,138 | 1,888,235 | 2,888,491 |
| ⊞ People & ... | 5,266,856,560 | 8,507,862,810 | 4,198,445,700 | 17,973,165,070 |
| ⊞ Pets & ... | 321,068,050 | 48,515,230 | 239,218,980 | 608,802,260 |