# Youtube Video Trending

**APAN 5450 Group 7:**
*Xiaoting Teng*
*Emily Xiao*
*Seline Yin*
*Keqi Yu*
*Chang Yuan*

# TABLE OF CONTENTS

# Background & Business Plan

**Provide Recommendations for Youtube Channel:**

❏ Knowing what is the best publishing time for videos could help businesses target advertisements specific to viewer country and taste

❏ Explore the key factors that influence viewers' preferences and content based on variables such as

  a. Likes & Dislikes
  b. Publish time
  c. Views count
  d. Comment Counts

# Data Source & Description

❏ The dataset is a version of Youtube Trending Videos Statistics from Kaggle which includings Youtube videos that are most popular on a daily basis. The Size of Data is 2.79 GB in total, Including Json files and csv across 11 regions around the world. And based on the main business plan, our team plans to use 6 regions (1.5 GB) to do our analysis, respectively USA, Canada, France, Russia, South Korea and Japan.

❏ Data Link:
https://www.kaggle.com/datasets/rsrishav/youtube-trending-video-dataset

# Team Responsibilities

## Xiaoting
VPC Setup and EC2 instance setup with key pairs

## Seline
Upload and restore data in S3

## Keqi
Clean data in AWS Glue

## Chang
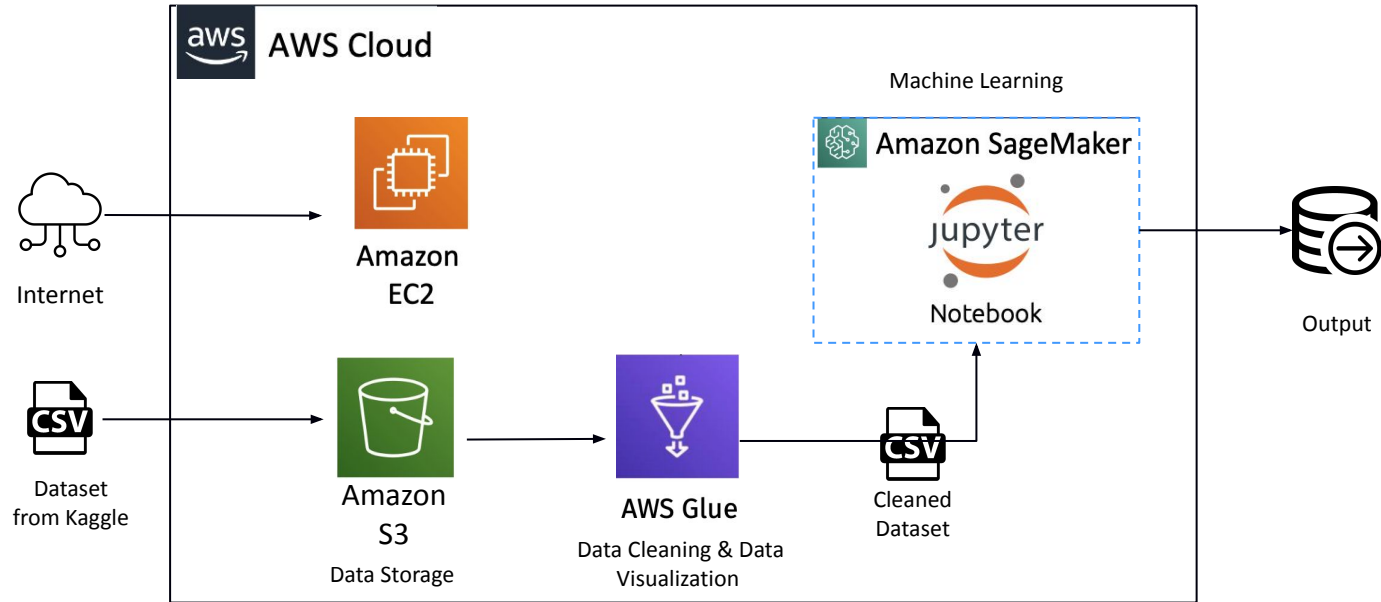Integrate Data from S3 and configuration

## Emily
Classification in AWS SageMaker

# Could Architecture Description

- ❏ **Amazon EC2 (future use)** – The reason why we believe EC2 can be used in future is after we trained our model to solve YouTuber's need and easy to access. Flask of EC2 might be used with EC2.

- ❏ **Amazon Virtual Private Cloud (Amazon VPC)** – AWS VPC uses security groups as a firewall to control traffic at the instance level, while it also uses network access control lists as a firewall to control traffic at the subnet level. VPC provides much more granular control over security.

- ❏ **Amazon Simple Storage Service (S3)** – Our group use S3 because it can assist the team store data at the lowest cost, backup, and restore data, as well as providing great monitoring to ensure data security, S3 is essential in this research. Also sine S3 offers rich security controls, which means it benefits from a data center and network architecture built to meet the requirements of most security-sensitive organizations. There is also need to set up the own securities plans such as take it own control and permission for the system.

- ❏ **Amazon Glue** – Since AWS Glue is a serverless data integration service that makes data preparation simpler, faster, and cheaper. Thus, our group use it to help data cleaning part and visualization.

- ❏ **Amazon Sagemaker** – Prepare, build, train, and deploy high-quality machine learning models quickly by bringing together a broad set of capabilities purpose-built for machine learning. So our group use Sagemarker for machine learning to predict audience behaviors and count views.

# Security Plan

## IAM

Identity and Access Management since we cannot create the role right now, we use the default LabRole
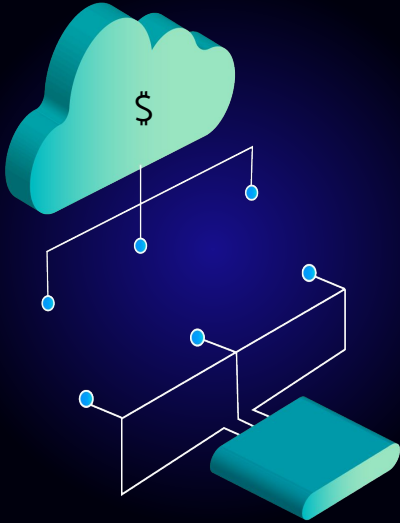
## VPC

Make public subnets so the group can use various platforms and devices to access, analyze, and execute analysis on the data.

## NAT Gateway

Permit proxy access to application servers for a group of servers in a public subnet using the internet

# Cost Analysis

## AWS VPC

★ Total NAT Gateway usage and data processing cost (Monthly): 32.94 USD
★ Total PrivateLink endpoints and data processing cost (Monthly): 7.32 USD

## AWS EC2 Instance

★ Amazon EC2 Instance Savings Plans instances (Monthly): 2.48 USD
★ Amazon EC2 instances (Upfront): 29.78 USD
★ Total Monthly cost: 3.28 USD

## AWS S3

★ Total Monthly cost: 0.5 USD
★ AWS Glue Monthly 8.86USD
★ Sagemaker Monthly ml.t3.medium 3.00USD

**Total Cost:**
★ Upfront: 29.78 USD
★ Monthly: 55.10 USD

# Success Criteria

## Quantitative

**Cost-effectiveness**
The cost of implementing technologies for applicable data and storage is saved for further analysis, which the analysts can directly analyze by using this prepared AWS technologies

**Data Size**
The data is securely and successfully imported into the Cloud ready to analyze and predict variables of interest

## Qualitative

**Data Quality**
Assume raw data from data resources are up-to-date and well-collected, the processed data should be reliable, and the search result from processed data should be accurate.

**Meaningful**
The analysis will provide meaning to social reality by understanding how an individual subjectively perceived.

# Implementation

| VPC ID | State | DNS hostnames | DNS resolution |
|---|---|---|---|
| 📋 vpc-0db3769dbbe01bc25 | ✓ Available | Enabled | Enabled |

## Instance summary for i-07e2bd36fbf10a0c2 (Group 7 Web Server 1) Info

Updated less than a minute ago

[ ↻ ] [ Connect ] [ Instance state ▼ ] [ Actions ▼ ]

| | | |
|---|---|---|
| **Instance ID**<br>📋 i-07e2bd36fbf10a0c2 (Group 7 Web Server 1) | **Public IPv4 address**<br>📋 52.90.154.101 \| open address ⬈ | **Private IPv4 addresses**<br>📋 10.0.2.177 |
| **IPv6 address**<br>– | **Instance state**<br>✓ Running | **Public IPv4 DNS**<br>📋 ec2-52-90-154-101.compute-1.amazonaws.com \|<br>open address ⬈ |
| **Hostname type**<br>IP name: ip-10-0-2-177.ec2.internal | **Private IP DNS name (IPv4 only)**<br>📋 ip-10-0-2-177.ec2.internal | |
| **Answer private resource DNS name**<br>IPv4 (A) | **Instance type**<br>t2.micro | **Elastic IP addresses**<br>– |
| **Auto-assigned IP address**<br>📋 52.90.154.101 [Public IP] | **VPC ID**<br>📋 vpc-0db3769dbbe01bc25 (Group 7 lab-vpc) ⬈ | **AWS Compute Optimizer finding**<br>ⓘ Opt-in to AWS Compute Optimizer for recommendations.<br>\| Learn more ⬈ |
| **IAM Role**<br>– | **Subnet ID**<br>📋 subnet-0e1d91b173e545bd5 (group-7-lab-subnet-public2) ⬈ | **Auto Scaling Group name**<br>– |

Setting up the VPC
Setting up key paris and Launch EC2 Instance

## Access control list (ACL)

Grant basic read/write permissions to other AWS accounts. Learn more 🔗

Edit

ℹ️ **This bucket has the bucket owner enforced setting applied for Object Ownership**
When bucket owner enforced is applied, use bucket policies to control access. Learn more 🔗

| Grantee | Objects | Bucket ACL |
|---|---|---|
| Bucket owner (your AWS account)<br>Canonical ID: 📋 d0911766e6081fa57f20897186b20c1d23088a2c901cac58530d525d43df4b5e | List, Write | Read, Write |
| Everyone (public access)<br>Group: 📋 http://acs.amazonaws.com/groups/global/AllUsers | - | - |
| Authenticated users group (anyone with an AWS account)<br>Group: 📋 http://acs.amazonaws.com/groups/global/AuthenticatedUsers | - | - |
| S3 log delivery group<br>Group: 📋 http://acs.amazonaws.com/groups/s3/LogDelivery | - | - |

Created bucket and upload dataset in S3 with subsequent access.

Also for AWS Glue and AWS Sagemaker

| Job name ▽ | Type | Last modified ▼ | AWS Glue version |
|---|---|---|---|
| group7dataclean | Glue ETL | 12/11/2022, 10:18:33 PM | 3.0 |

Created a job for cleaning data in AWS Glue.

Created domain of AWS Sagemaker with the VPC setted up for connecting notebook and using Machine Learning in Jupyter Notebook.

| Name | Status | Domain ID |
|---|---|---|
| group7 | ⊘ Ready | ⧉ d-6oplofmd2joc |

| Created | Last modified | VPC |
|---|---|---|
| Thu Dec 08 2022 09:48:01 GMT-0500 (Eastern Standard Time) | Thu Dec 08 2022 09:52:53 GMT-0500 (Eastern Standard Time) | vpc-0db3769dbbe01bc25 |

| Authentication method | Execution role | |
|---|---|---|
| AWS Identity and Access Management (IAM) | ⧉ arn:aws:iam::715923724147:role/LabRole | |

# Data Cleaning in AWS Glue

```python
# Connect to S3 bucket
bucket = 'group7bucketnew'
file_key = ['US_youtube_trending_data.csv','CA_youtube_trending_data.csv','FR_youtube_trending_data.csv',
s3uri = []
for i in file_key:
    s3uri.append('s3://{}/{}'.format(bucket, i))
```

```python
data = pd.DataFrame()
for i in s3uri:
    df = pd.read_csv(i)
    df["origin"] = i[21:]
    data = data.append(df)
```

1. Connect to S3 bucket from AWS Glue and merge all the data files together

The original dataset has

**1012039 rows x 17 columns**

**2. Convert Data Type for Date Column to datetime format and convert datetime to date, month, year, time, hour**

**3. Delete abnormal data**

```python
mask = (data.view_count<=0)
df = data.loc[~mask]
```

```python
# Transforming Trending date column to datetime format
df['trendingAt'] = pd.to_datetime(df['trendingAt'], format='%Y-%m-%dT%H:%M:%SZ')
df['publishedAt'] = pd.to_datetime(df['publishedAt'], format='%Y-%m-%dT%H:%M:%SZ')
```

```python
df.insert(loc=3, column='published_date', value=df.publishedAt.dt.date)
df.insert(loc=4, column='published_month', value=df.publishedAt.dt.month_name())
df.insert(loc=5, column='published_day', value=df.publishedAt.dt.day_name())

df.insert(loc=10, column='trending_date', value=df.trendingAt.dt.date)
df.insert(loc=11, column='trending_month', value=df.trendingAt.dt.month_name())
df.insert(loc=12, column='trending_day', value=df.trendingAt.dt.day_name())
```

```python
df.insert(loc=6, column='published_time', value=df.publishedAt.dt.time)
df.insert(loc=7, column='published_hour', value=df.publishedAt.dt.hour)

df.insert(loc=13, column='trending_time', value=df.trendingAt.dt.time)
df.insert(loc=14, column='trending_hour', value=df.trendingAt.dt.hour)
```
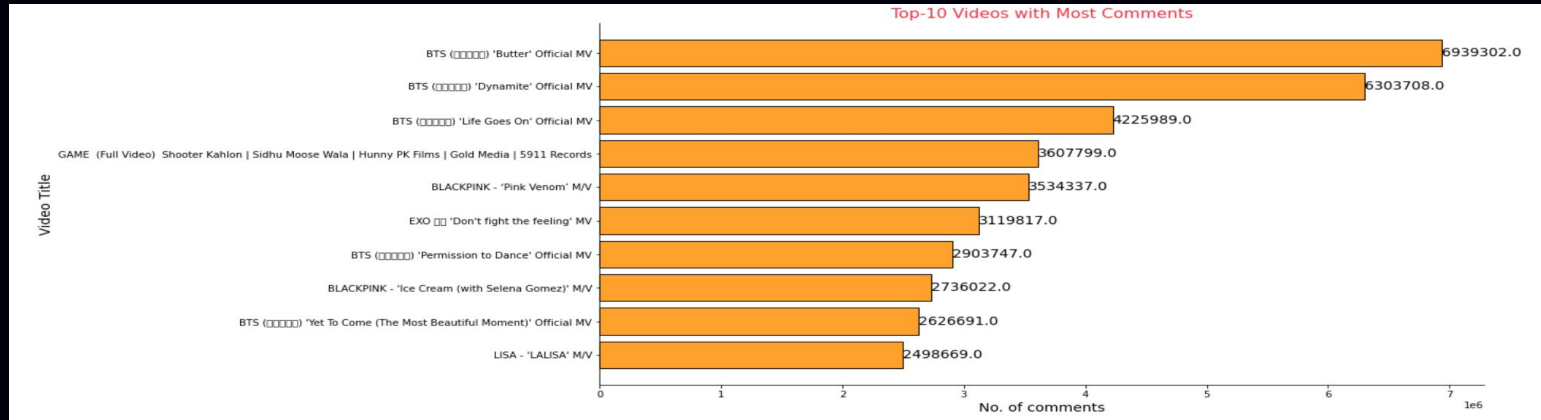
**Finally, after the data cleaning, the cleaned dataset has 1011848 rows x 27 columns**
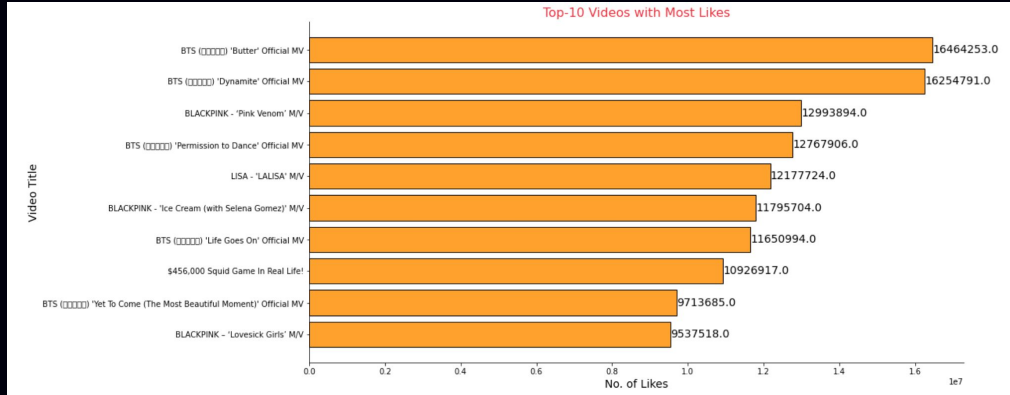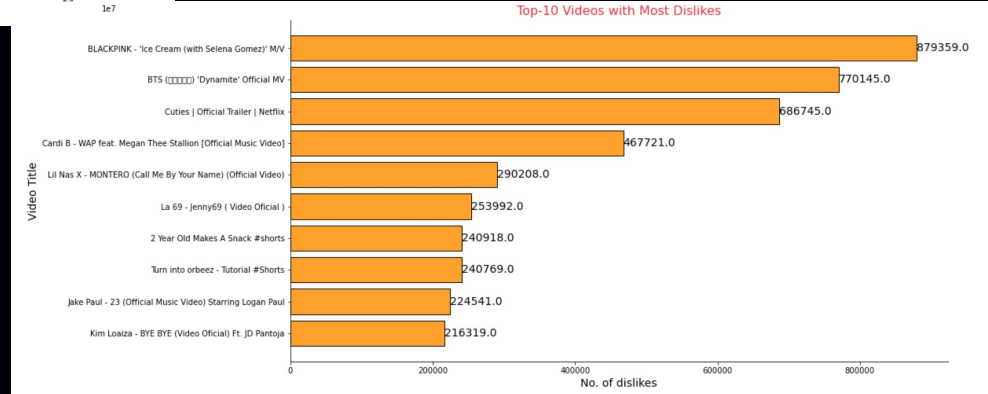
# Implement Results


Top-10 Videos with Most Views

BTS and BlackPink videos make the most of the top 10 videos with most views and most comments


Top-10 Videos with Most Comments

# Implement Results

Top-10 Videos with Most Likes



Top-10 Videos with Most Likes

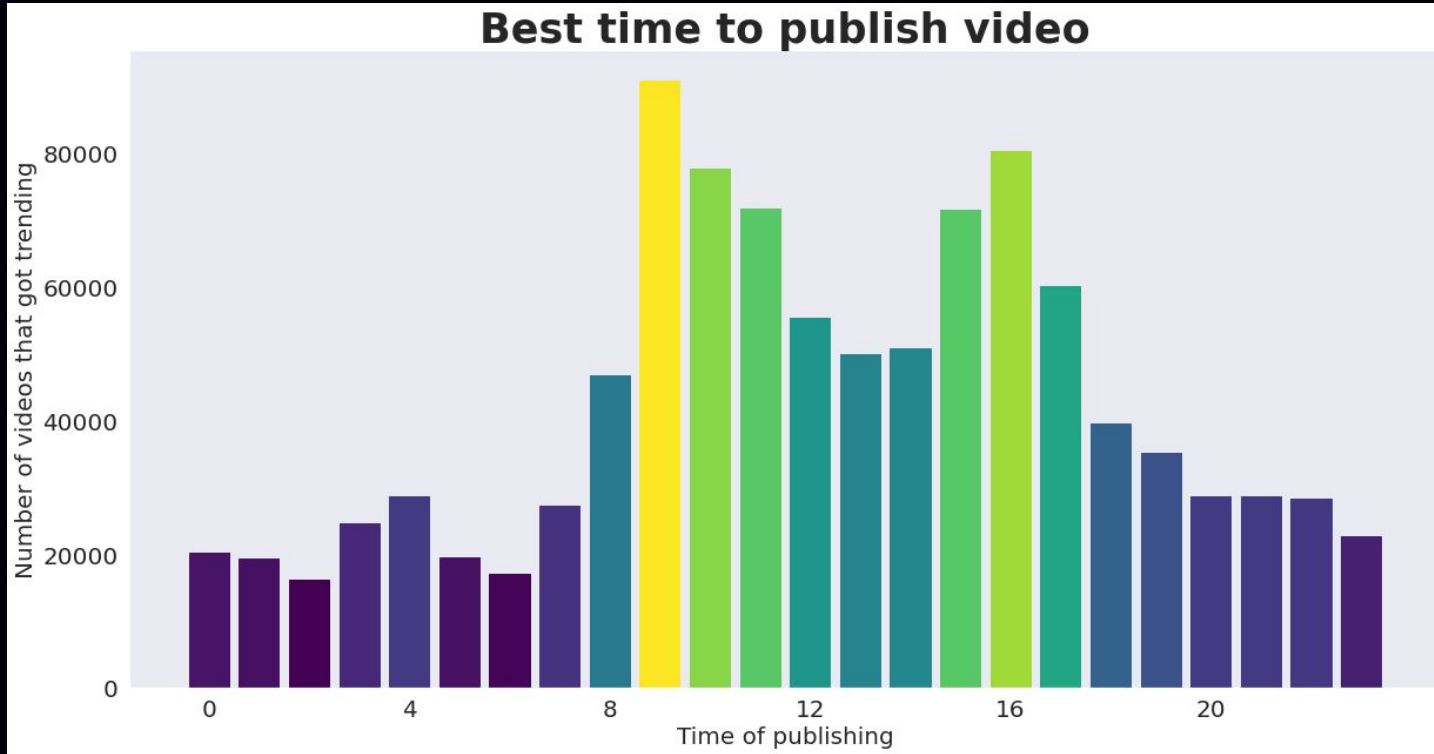| Video Title | No. of Likes |
| --- | --- |
| BTS (방탄소년단) 'Butter' Official MV | 16464253.0 |
| BTS (방탄소년단) 'Dynamite' Official MV | 16254791.0 |
| BLACKPINK - 'Pink Venom' M/V | 12993894.0 |
| BTS (방탄소년단) 'Permission to Dance' Official MV | 12767906.0 |
| LISA - 'LALISA' M/V | 12177724.0 |
| BLACKPINK - 'Ice Cream (with Selena Gomez)' M/V | 11795704.0 |
| BTS (방탄소년단) 'Life Goes On' Official MV | 11650994.0 |
| $456,000 Squid Game In Real Life! | 10926917.0 |
| BTS (방탄소년단) 'Yet To Come (The Most Beautiful Moment)' Official MV | 9713685.0 |
| BLACKPINK – 'Lovesick Girls' M/V | 9537518.0 |

**BTS and BlackPink videos make the most of the top 10 videos with most likes, which coincides with the previous slide**

**The video with most dislikes were also from the music video.**

Top-10 Videos with Most Dislikes

| Video Title | No. of dislikes |
| --- | --- |
| BLACKPINK - 'Ice Cream (with Selena Gomez)' M/V | 879359.0 |
| BTS (방탄소년단) 'Dynamite' Official MV | 770145.0 |
| Cuties | Official Trailer | Netflix | 686745.0 |
| Cardi B - WAP feat. Megan Thee Stallion [Official Music Video] | 467721.0 |
| Lil Nas X - MONTERO (Call Me By Your Name) (Official Video) | 290208.0 |
| La 69 - Jenny69 ( Video Oficial ) | 253992.0 |
| 2 Year Old Makes A Snack #shorts | 240918.0 |
| Turn into orbeez - Tutorial #Shorts | 240769.0 |
| Jake Paul - 23 (Official Music Video) Starring Logan Paul | 224541.0 |
| Kim Loaiza - BYE BYE (Video Oficial) Ft. JD Pantoja | 216319.0 |

Top-10 Videos with Most Dislikes

# Results



**This graph can help YouTubers to pick the best time to publish their videos.**

# Classification Results

| | | |
|---|---|---|
| **Data Split** | Training Set: 0.8 | Testing Set: 0.2 |
| **Model Used** | Support Vector Machine Classification | Logistic Regression Classification |
| **RMSE Value** | 0.1906 | 0.2527 |
| **Usage** | Using our model, YouTubers can predict their **View Counts** for their videos and know how many audience would view their videos | |

# Future Plan

1.  Add AWS Key Management Service(KMS) encryption key for additional security
2.  Deploy Flask on EC2 for our product development
3.  Create IAM roles to achieve specific permissions for future development and maintenance
4.  Use AWS RDS database

# THANKS!