

# YT Chronicles - Process Book

## 1. Basic Information

*Project title* : YT Chronicles

*Member details*

- 1) Kalyan Shankar Ragam u1471323@utah.edu
- 2) Vijay Surya Vempati u1472324@utah.edu

*Github Link* - <https://github.com/dataviscourse2024/group-project-yt-chronicles>

## 2. Background and Motivation

This project is driven by my deep interest in exploring the dynamics and patterns of digital media, with a particular focus on YouTube. As a major platform for global content creation and consumption, YouTube has evolved into a hub for influencers, brands, and creators to connect with audiences worldwide. Our motivation stems from a desire to understand the underlying factors that contribute to a channel's success. By visualizing how these factors interplay to influence a channel's performance, we aim to gain insights into user engagement and content virality all over the world and in a specific country.

Given YouTube's global reach, this project provides a unique opportunity to visualize data from channels across various countries and categories. By examining these, we hope to uncover broader trends in global content consumption and identify potential correlations between channel performance and regional factors. This insight can be valuable for aspiring content creators and digital marketers seeking to understand and leverage the nuances of digital media engagement.

## 3. Project Objectives

The primary objectives of this project are to explore and visualize various aspects of YouTube channel metrics to gain insights into content creation and consumption patterns across different countries and categories. Specifically, the project aims to:

- **Analyze Country Variations:** Identify the top 10 YouTubers in each country based on their view counts and subscriber numbers.
- **Identify Recent Trends:** Examine recent changes in subscriber counts and video views over the past 30 days for various channels to uncover current content consumption trends.

- **Determine Popular Content Categories:** Visualize the most popular YouTube channel categories (e.g., gaming, education, lifestyle) in specific countries to gain insights into regional content preferences.

## 4. Data

We have obtained our data from Kaggle. The youtube dataset contains information about the data of top 1000 youtubers and contains fields like rank, id, sub count, view count, category, number of uploads, country, monthly and yearly earnings, date of creation which we believe should be good for this project. The link to the dataset is

<https://www.kaggle.com/datasets/neljiriyewithana/global-youtube-statistics-2023/data>

## 5. Data Processing

For data processing, we ran a python code which helped us get the required columns for the visualization from the original dataset we had. We had to clean and remove the rows that have invalid or garbage data in it. After writing the code accordingly and doing the data processing, we have ended up with the csv we are satisfied with and can use for our visualization. Here is the snapshot of the python code we ran for the data processing and cleaning

```
data_preprocess.py > ...
1  import pandas as pd
2  import os
3
4  # Load the CSV file with the specified encoding
5  df = pd.read_csv("Global YouTube Statistics.csv", encoding='ISO-8859-1')
6
7  # Select the required columns
8  columns_to_keep = ['rank', 'subscribers', 'Title', 'category', 'video views', 'highest_monthly_earnings', 'Country']
9  processed_df = df[columns_to_keep]
10
11 # Print the number of rows before filtering
12 rows_before = processed_df.shape[0]
13 print(f"Number of rows before filtering: {rows_before}")
14
15 # Drop rows with specified conditions
16 filtered_df = processed_df[
17     (processed_df['video views'] > 0) &
18     (processed_df['highest_monthly_earnings'] > 0) &
19     (processed_df['category'].notna()) & (processed_df['category'] != '') &
20     (processed_df['Country'].notna()) & (processed_df['Country'] != '')
21 ]
22
23 # Clean the 'Title' column by removing unwanted characters
24 filtered_df['Title'] = filtered_df['Title'].str.replace(r'\[.*?\]', '', regex=True) # Remove anything in brackets
25 filtered_df['Title'] = filtered_df['Title'].str.replace(r'Ã~Â¿Ã%', '', regex=True) # Remove specific unwanted characters
26 filtered_df['Title'] = filtered_df['Title'].str.strip() # Strip whitespace from the ends
27
```

We did run into a few problems using this code which I'm gonna discuss more about in the visualization design section. But for now, this code did a great job in securing most of the data we wanted perfectly. Here is the snapshot of the original dataset we used for this visualization.

```
data_preprocess.py Global YouTube Statistics.csv index.html JS scripts jupyter_notebook_config.py 1
Global YouTube Statistics.csv
1 rank,youtuber,subscribers,video_views,category,Title,uploads,Country,Abbreviation,channel_type,video_views_rank,country_rank,channel_type_rank,video_views_for_the_last_30_days,lowest_monthl
2 1,T-Series,245000000,2.28E+11,Music,T-Series,20082,India,TN,Music,1,1,1,2258000000,564600,9000000,6800000,108400000,2000000,2006,Mar,13,28.1,1366417754,5.36,471031528,20.593684,78.96280
3 2,YouTube Movies,170000000,0,Film & Animation,youtubemovies,1,United States,US,Games,4055159,7670,7423,12,0,0.05,0.04,0.58,nan,2006,Mar,5,88.2,328239523,14.7,270663028,37.09024,-95.712891
4 3,MrBeast,166000000,28368841870,Entertainment,MrBeast,741,United States,US,Entertainment,48,1,1,1348000000,337000,5400000,64700000,8000000,2012,Feb,20,88.2,328239523,14.7,270663028,
5 4,Cocomelon - Nursery Rhymes,162000000,1.64E+11,Education,Cocomelon - Nursery Rhymes,966,United States,US,Education,2,2,1,1975000000,493800,7900000,94800000,1000000,2006,Sep,1,88.2,
6 5,SET India,159000000,1.48E+11,Shows,SET India,116536,India,IN,Entertainment,3,2,2,1824000000,455900,7300000,5500000,87500000,1000000,2006,Sep,20,28.1,1366417754,5.36,471031528,20.593684,78
7 6,Music,119000000,0,nan,Music,0,nan,nan,Music,4057944,nan,nan,nan,0,0,0,0,nan,2013,Sep,24,nan,nan,nan,nan,nan,nan
8 7, Kids Diana Show,112000000,93247040539,People & Blogs, Kids Diana Show,1111,United States,US,Entertainment,5,3,3,731674000,182900,2900000,2200000,35100000,nan,2015,May,12,88.2,
9 8,PewDiePie,111000000,29058044447,Gaming,PewDiePie,4716,Japan,JP,Entertainment,44,1,4,39184000,9800,156700,117600,1900000,nan,2010,Apr,29,63.2,126226568,2.29,115782416,36.204824,138.252924
10 9,Like Nastya,106000000,90479060027,People & Blogs,Like Nastya Vlog,493,Russia,RU,People,630,5,25,48947000,12200,195800,146800,2300000,100000,2016,Jan,14,81.9,144373535,4.59,107683889,61.52
11 10,Vlad and Niki,98900000,77180169894,Entertainment,Vlad and Niki,574,United States,US,Entertainment,8,5,6,580574000,145100,2300000,1700000,27900000,600000,2018,Apr,23,88.2,328239523,14.7,2
12 11,Zee Music Company,96700000,57856289381,Music,Zee Music Company,8548,India,IN,Music,12,3,2,803613000,200900,3200000,2400000,38600000,1100000,2014,Mar,12,28.1,1366417754,5.36,471031528,20.
13 12,WWE,96000000,77428473662,Sports,WWE,70127,United States,US,Sports,7,6,1,714614000,178700,2900000,2100000,34300000,600000,2007,May,11,88.2,328239523,14.7,270663028,37.09024,-95.712891
14 13,Gaming,93600000,0,nan,Gaming,0,nan,nan,Games,4057944,nan,1,nan,0,0,0,0,nan,2013,Dec,15,nan,nan,nan,nan,nan,nan
15 14,BLACKPINK,89800000,32144597566,People & Blogs,BLACKPINK,543,South Korea,KR,Music,32,1,3,498930000,124700,2000000,1500000,23900000,700000,2016,Jun,29,94.3,51709098,4.15,42106719,35.907757,
16 15,Goldmines,86900000,24118220500,Film & Animation,goldmines,1,nan,nan,Music,4056562,nan,5663,18,0,0,0,0,0.05,0.86,nan,2006,Aug,15,nan,nan,nan,nan,nan,nan
17 16,Sony SAB,83000000,1.01E+11,Shows,Sony SAB,71270,India,IN,Entertainment,4,5,7,1657000000,414200,6600000,5000000,79600000,1100000,2007,Aug,4,28.1,1366417754,5.36,471031528,20.593684,78.962
18 17,5-Minute Crafts,80100000,26236700200,Howto & Style,5-Minute Crafts,2,0,1,United Kingdom,GB,Entertainment,4057901,4707,6781,1,0,0,0,0.05,nan,2020,Jul,27,60,66834405,3.85,55008316,55.37808
19 18,BANGTANTV,75600000,20826993957,Music,BANGTANTV,2281,South Korea,KR,Music,112,2,4,168290000,42100,673200,504900,8100000,400000,2012,Dec,17,94.3,51709098,4.15,42106719,35.907757,127.766922
20 19,Sports,75000000,0,nan,sports,3,United States,US,Entertainment,3898122,6266,5395,16,0,0.06,0.05,0.77,nan,2006,Jan,30,88.2,328239523,14.7,270663028,37.09024,-95.712891
21 20,Justin Bieber,71600000,30608119724,Music,Justin Bieber,249,Canada,CA,Music,38,1,6,176326000,44100,705300,529000,8500000,100000,2007,Jan,15,68.9,36991981,5.56,30628482,56.130366,-106.3467
22 21,HYBE LABELS,71300000,28634566938,Music,HYBE LABELS,1337,South Korea,KR,Music,46,3,5,598173000,149500,2400000,1800000,28700000,900000,2008,Jun,4,94.3,51709098,4.15,42106719,35.907757,127.
23 22,Zee TV,70500000,73139054467,Entertainment,Zee TV,129204,India,IN,Entertainment,9,6,8,1707000000,426800,6800000,5100000,81900000,900000,2005,Dec,11,28.1,1366417754,5.36,471031528,20.59368
24 23,Pinkfong Baby Shark - Kids' Songs & Stories,68200000,38843229963,Education,Pinkfong Baby Shark - Kids' Songs & Stories,2865,United States,US,Education,23,8,2,473387000,118300,1900000,146
25 24,Canal KondZilla,66500000,36775585925,Music,Canal KondZilla,2572,Brazil,BR,Music,25,1,7,447223000,0,0,0,0,nan,2012,Mar,21,51.3,212559417,12.08,183241641,-14.235004,-51.92528
26 25,ChuChu TV Nursery Rhymes & Kids Songs,65900000,45757850229,Education,ChuChu TV Nursery Rhymes & Kids Songs,633,India,IN,Education,18,7,3,420292000,105100,1700000,1300000,20200000,500000,
27 26,Shemaroo Filmi Gaane,65600000,28648024439,Music,Shemaroo Filmi Gaane,8502,India,TN,Music,47,8,8,254961000,63700,1000000,764900,12700000,400000,2010,Jun,11,28.1,1366417754,5.36,471031528,
```

This data set contained some garbage and invalid data as you can see from the snapshot. The dataset is also quite extensive and contains about 28 columns of information which we won't be using for our visualization project. We picked the columns which we feel we will use for our project but we can always go back and get more columns to work with if we want to add more visualizations apart from the planned ones. Here is the snapshot of the processed data after deciding what we want and running the code.

```
data > data.csv
1 rank,subscribers,Title,category,video_views,highest_monthly_earnings,Country
2 1,245000000,T-Series,Music,2.28E+11,9000000,India
3 3,166000000,MrBeast,Entertainment,28368841870,5400000,United States
4 4,162000000,Cocomelon - Nursery Rhymes,Education,1.64E+11,7900000,United States
5 5,159000000,SET India,Shows,1.48E+11,7300000,India
6 7,112000000,ýýý Kids Diana Show,People & Blogs,93247040539,2900000,United States
7 8,111000000,PewDiePie,Gaming,29058044447,156700,Japan
8 9,106000000,Like Nastya Vlog,People & Blogs,90479060027,195800,Russia
9 10,98900000,Vlad and Niki,Entertainment,77180169894,2300000,United States
10 11,96700000,Zee Music Company,Music,57856289381,3200000,India
11 12,96000000,WWE,Sports,77428473662,2900000,United States
12 14,89800000,BLACKPINK,People & Blogs,32144597566,2000000,South Korea
13 16,83000000,Sony SAB,Shows,1.01E+11,6600000,India
14 18,75600000,BANGTANTV,Music,20826993957,673200,South Korea
15 20,71600000,Justin Bieber,Music,30608119724,705300,Canada
16 21,71300000,HYBE LABELS,Music,28634566938,2400000,South Korea
17 22,70500000,Zee TV,Entertainment,73139054467,6800000,India
18 23,68200000,Pinkfong Baby Shark - Kids' Songs & Stories,Education,38843229963,1900000,United States
19 25,65900000,ChuChu TV Nursery Rhymes & Kids Songs,Education,45757850229,1700000,India
20 26,65600000,Shemaroo Filmi Gaane,Music,28648024439,1000000,India
21 27,64600000,Colors TV,Shows,61510906457,4800000,India
22 28,61000000,T- SERIES BHAKTI SAGAR,Music,29533230328,0.04,India
23 29,59500000,Dude Perfect,Sports,16241549158,564800,United States
24 30,59500000,Movieclips,Film & Animation,59316472754,458700,United States
25 31,59300000,Tips Official,Music,33431802698,1700000,India
26 32,58400000,El Reino Infantil,Music,57271630846,2400000,Argentina
27 33,58000000,Wave Music,Music,40602020242,028100,India
```

## 6. Proposed Visualization Design

We aim to keep the page minimalistic to ensure a clean, user-friendly experience. Instead of overwhelming the user with multiple pages, our design will focus on a single, dynamic page that adapts based on user input.

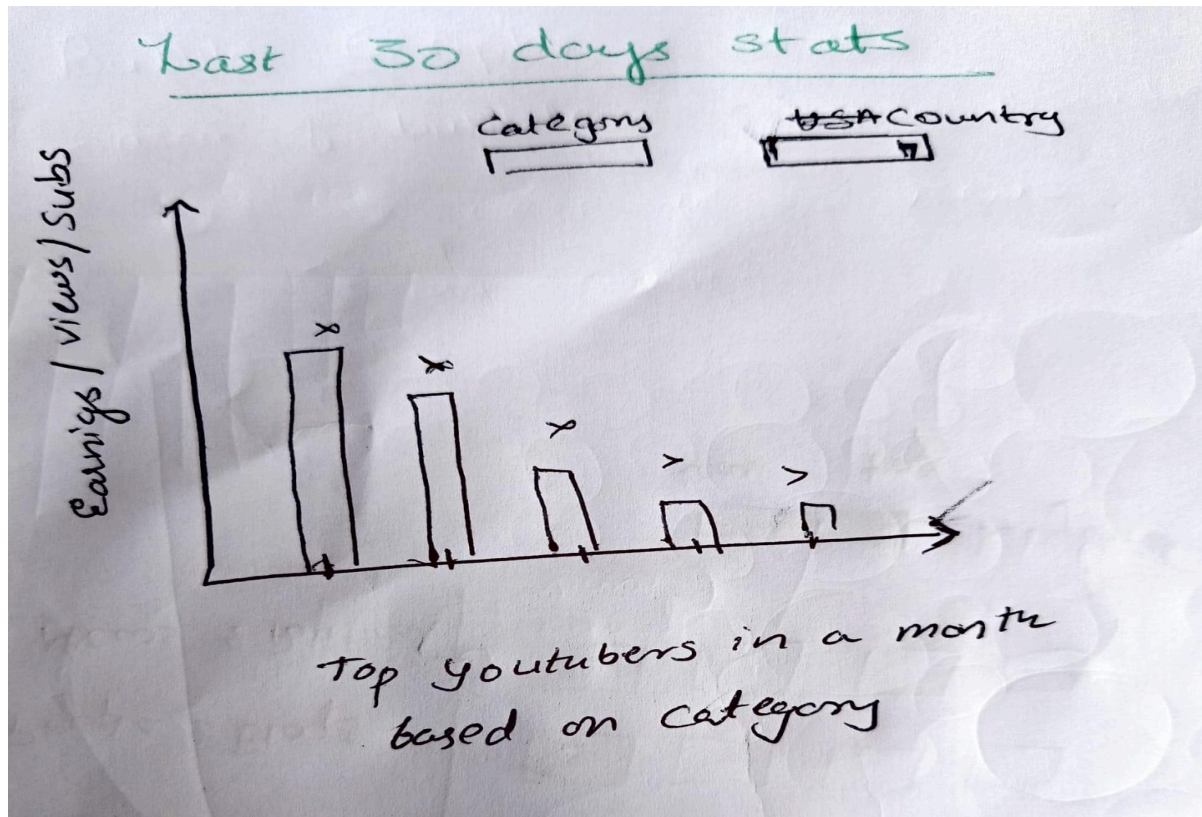
- 1) Upon visiting the site, users will be directed to the main page, where they will be greeted with a world map and a bar graph displaying the top 10 YouTubers in the world (by default). This graph can be dynamically sorted by either views or subscriber count and by countries as well. When a user selects a specific country on the map, the graph will update to show the top 10 YouTubers from that country, with the ability to sort them by views or subscriber count, as desired.

The world map can be used as a heatmap where the country's color intensity is proportional to the category (views/subscribers/earnings).

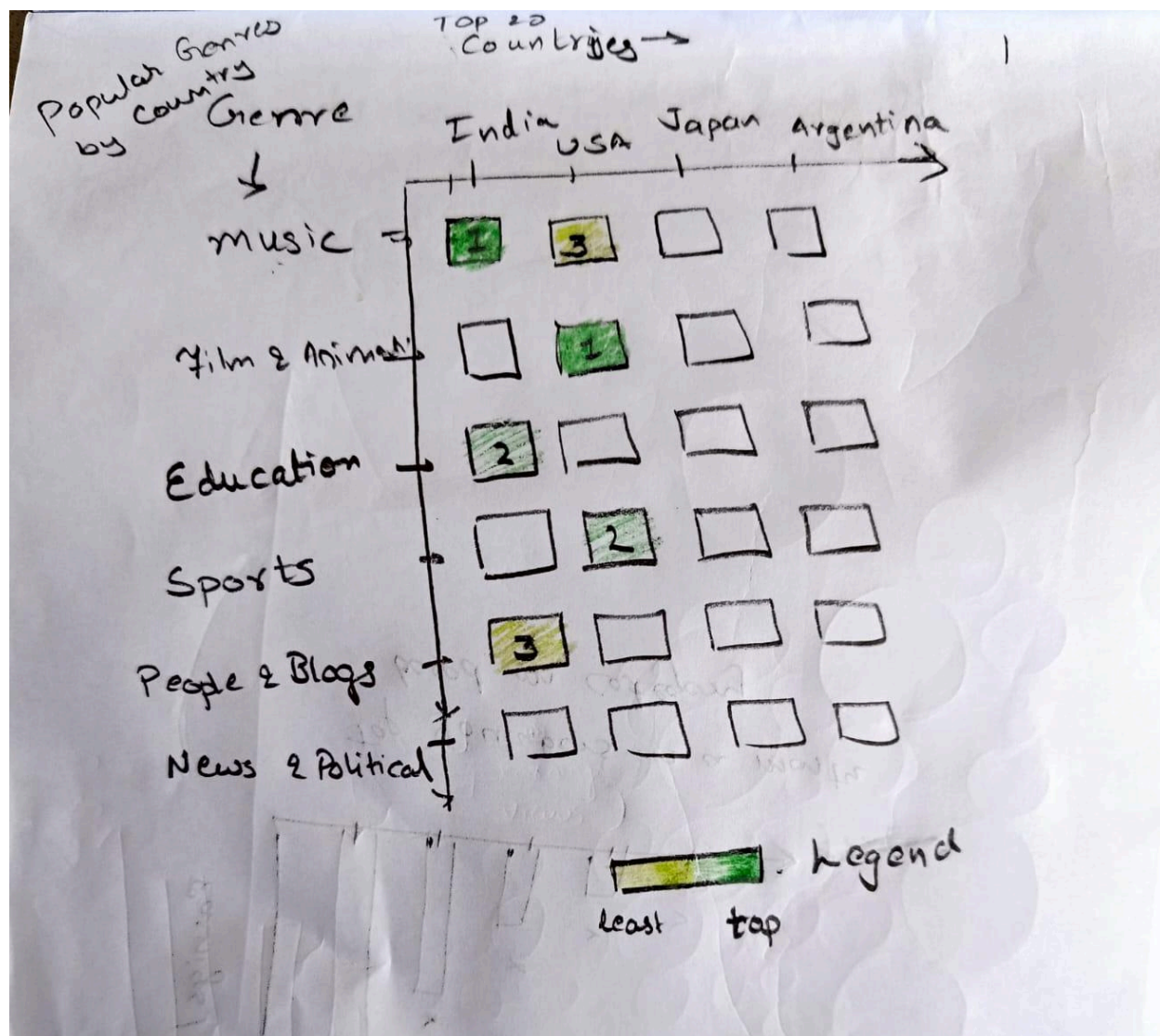


- 2) Below the main graph, we will display a second graph focused on recent trends, highlighting changes in subscriber counts and video views over the past 30 days. This will help uncover current content consumption trends either globally (by default) or for the selected country on the world map.

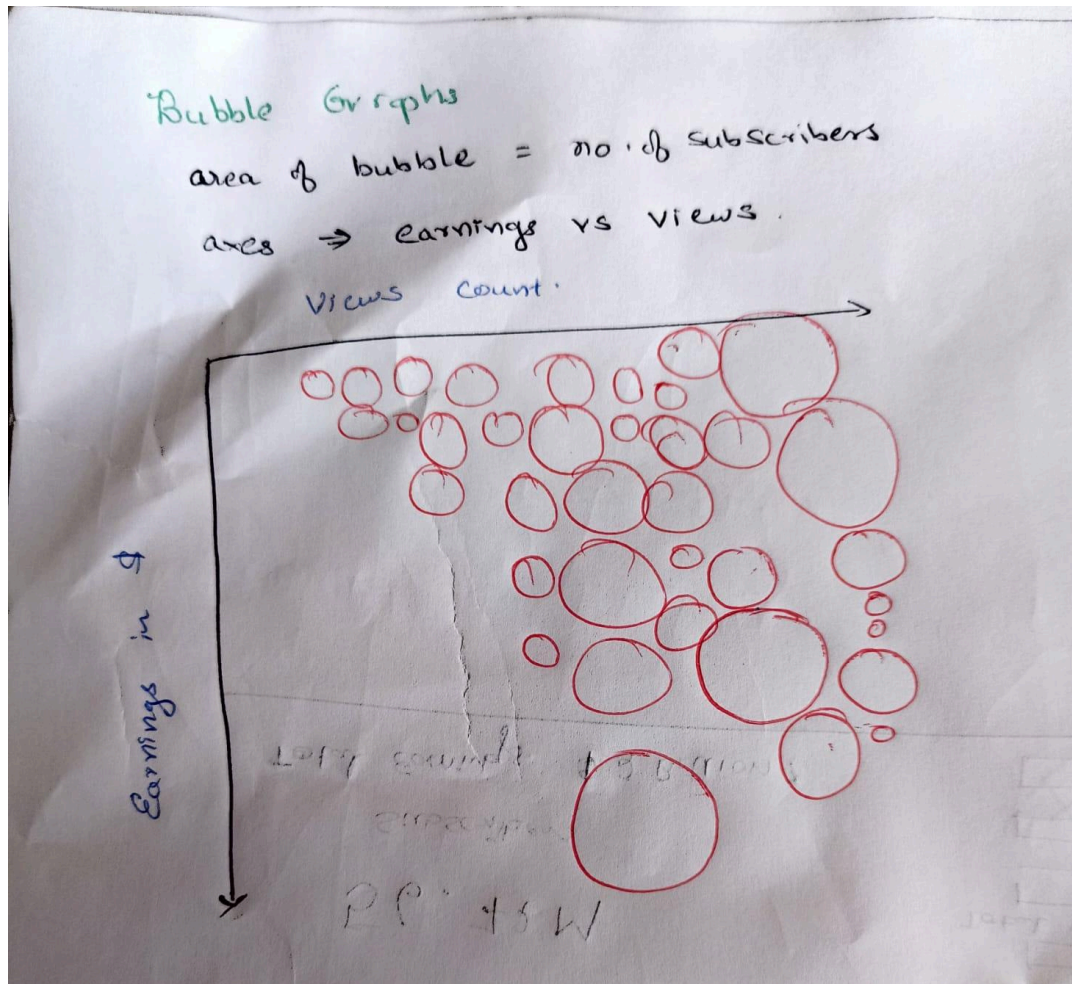




- 3) At the bottom of the page, we will feature a heat map that shows the popularity of different YouTube channel categories across countries. This visualization will reveal which types of content creators are most popular in specific regions.



**Optional :** We aim to implement a bubble graph where the size of the bubble indicates the subscriber count, axes represent Earnings vs View counts.



## 7. Proposed Must Have Features

- Interactive heatmaps and charts for visualizing country-specific YouTube stats.  
The world map should function as a heatmap where the color intensity reflects a selected metric (views, subscribers, or earnings).
- Filter options to select categories, countries  
By default, the bar graph displays the top 10 YouTubers globally.  
The graph should allow sorting by views or subscriber count.  
When a country is selected on the map, the graph should update to show the top 10 YouTubers from that country.
- Recent Trends Graph  
Display changes in subscriber counts and video views over the past 30 days.  
By default, show global trends, but when a country is selected, it should update to reflect trends for that specific country.
- Category Popularity Heatmap



A heatmap at the bottom of the page that visualizes the popularity of different YouTube categories (e.g., Gaming, Education, Music) across countries.

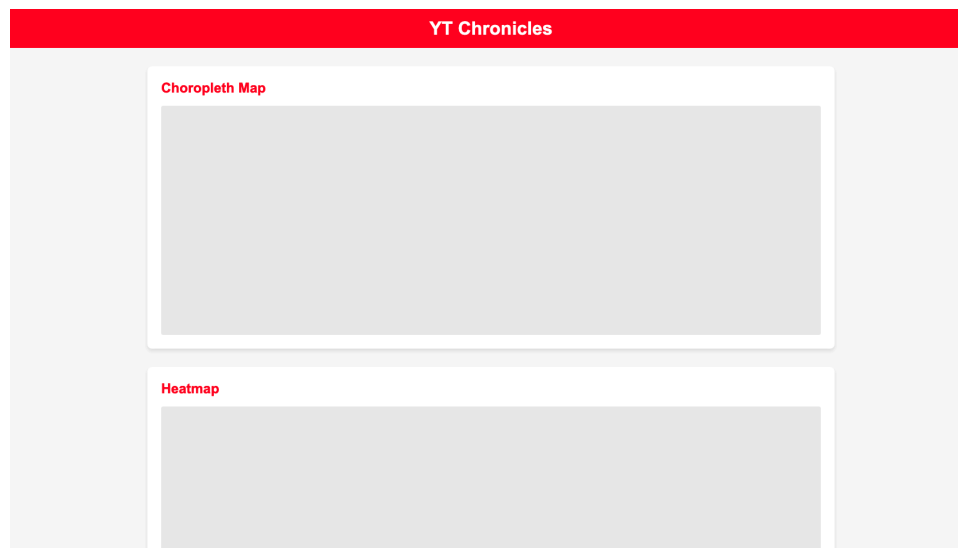
## 8. Proposed Optional Features

Here are some additional features we believe would enhance our project, though we recognize they may be challenging to implement within our given timeframe. Nevertheless, we will make an effort to incorporate them to further enrich the project

- We aim to create individual profiles for the top YouTubers, including their profile picture and key information displayed in a dashboard format, allowing users to quickly learn more about them.
- We want to add a feature that randomly recommends a YouTuber to watch, based on the country and category the user selects.
- We plan to leverage unused data, such as population, unemployment rate, and gross tertiary education enrollment, to find correlations with YouTube metrics and create insightful visualizations that reveal deeper patterns.
- Have a dark mode option for better user experience.

## 9. Implementation Process Stages

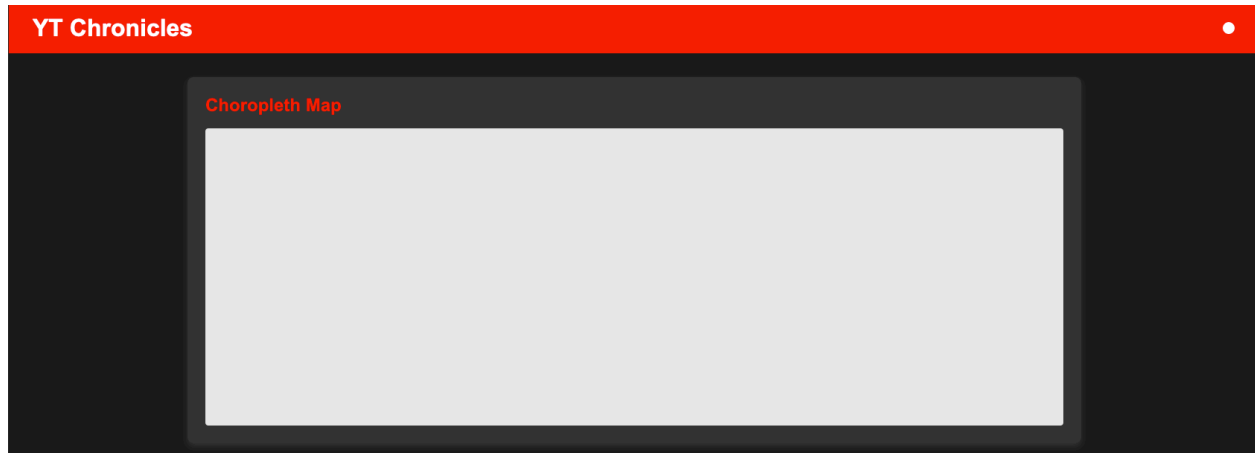
To start off, we wanted to have a skeleton website to work with and add our visualizations to which we could refine going on and rearrange our visualizations in. For the starting skeleton website design, we decided to keep it simple and go with the colors Red, White and Black primarily as those colors represent Youtube as a website. We plan to keep this color scheme going forward playing with their hues to make it visually appealing. Here is the screenshot of the initial skeleton website we made to work with



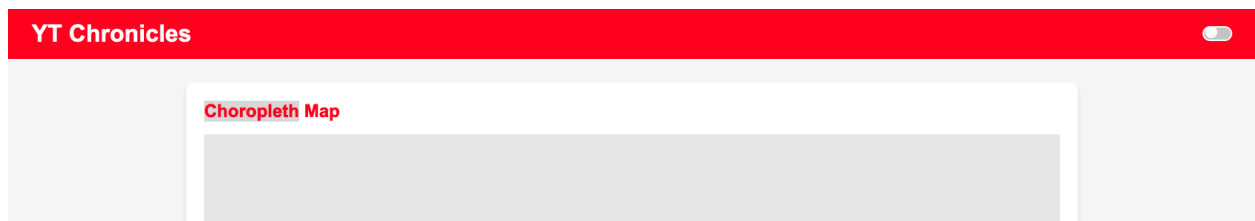


We have added names of the visualization we planned to and allocated space for them too. This is merely for our convenience for now and the website will be redesigned without those names and blocks on them in the future to make it visually appealing overall. For now, this is the canvas we are planning to go ahead and work with.

Our next goal is to add a dark mode feature to the website so we promptly did that and also made a few changes in the placement of text in the website. It looked something like this during this stage.



The issue we ran into here is that the dark mode button in the top right merges with the header background so we have to add an outline for the button to make it more visible like this.



After that is done, it is now time to work on the actual visualizations and we decided to start with the World Map visualization as we believe it is the most crucial and integral visualization in the whole webpage and making it work would mean other visualizations that are going to be integrated to it in the future are going to work fine too.

To start with our world map visualization, i first visited the website <https://geojson-maps.kyd.au/> which has been really helpful to me in downloading the world map in a .geojson format. After I have downloaded it, it is time to map it properly, give it color and arrange it into our website.

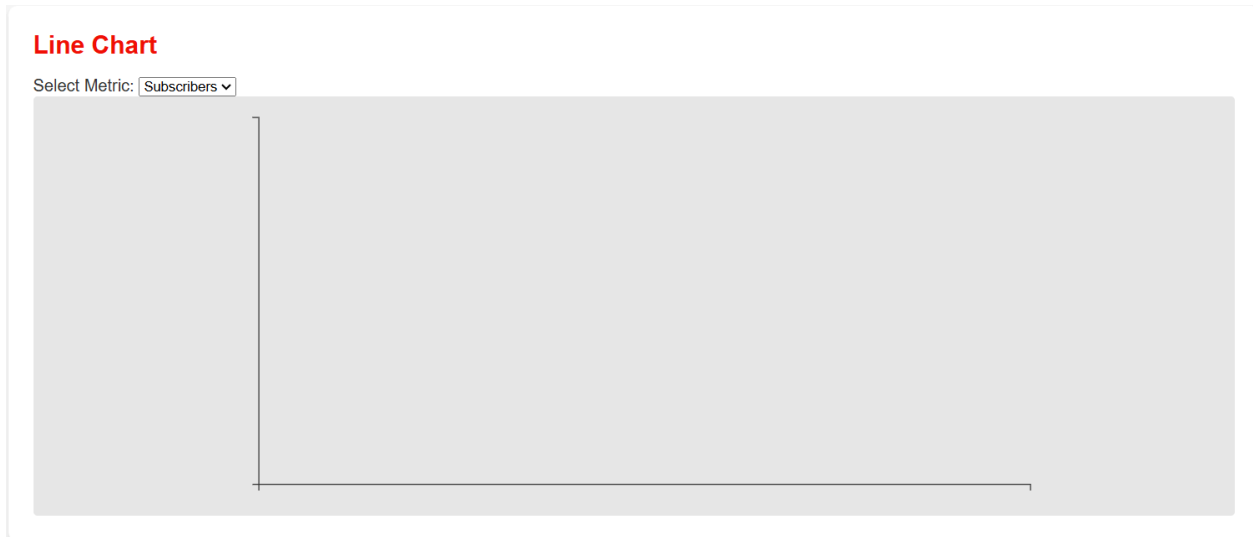


After our initial steps the visualization looked something like this. Though it is not perfect as the visualization exceeded the space allocated to it, we are happy that the visualization is working perfectly. We also added a selected country option to see what country we have selected at the moment. The next step obviously is to get our visualization into the required space. After further coding and tweaking this is what we ended up with.

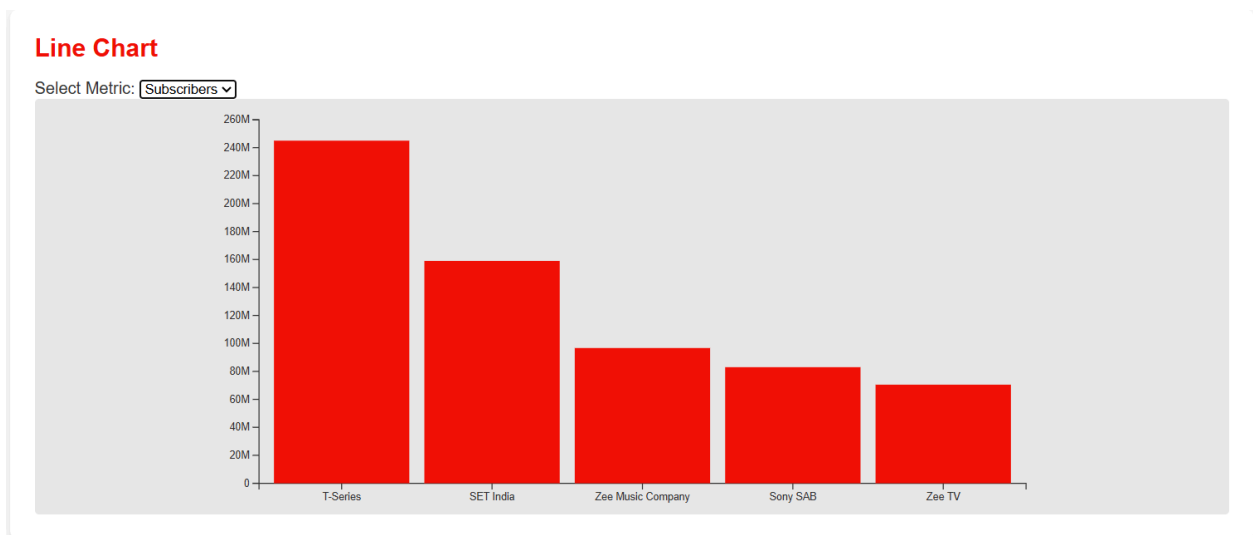


We were able to get the map properly into the block and also added a reset selection button to get back to none. This will be helpful to us when we start visualizing our line chart as setting the option to none would give us the top youtube channels in a set category in the entire world.

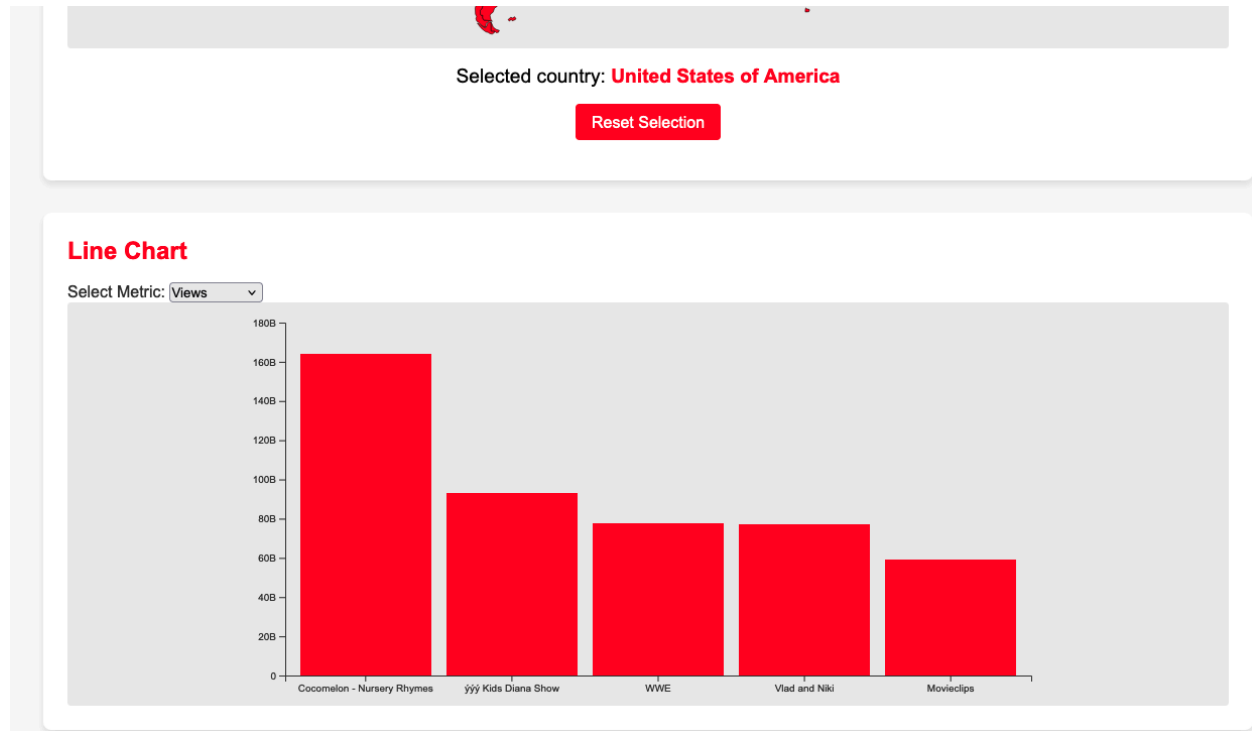
Now, the next step is to work on the line chart and this is where we ran into our first hurdle. The map we made has the location of USA as United States of America while our data set has it as just United States making it not being able to pull the data properly from our dataset showing something like this



But the visualization is coded and implemented perfectly in other cases where the country matches the dataset entry perfectly. So when you select India, it works properly showing something like this



We have worked on our dataset as well as the code to iron out these errors so that the data shows perfectly for all the countries as expected. After making a few changes we are able to output the desired output we need just like this



The line chart is now working perfectly. We have included three metrics to sort and select the youtube channels which are Views, Subscribers and Earnings. The plan was to include top 10 youtube channel per metric on the line chart but on further discussion we have decided to bring the count down to top 5 as there is just not enough data in the dataset we decided to show top 5 channels in most countries other than USA and India.

## 9. Planned Features

As of this milestone, we are able to get the World Map and Line chart visualization working as we intended. The project will continuously be updated in the upcoming milestones with more planned visualizations and other features. Here is a list of visualizations and features that are planned for the final release.

- A heatmap showing the top countries and the popular categories just as shown in our planned visualization design section (Must have)
- Adding a bubble graph comparing earnings to views (optional)
- Redesigning the website and structuring the visualization properly to make it more visually appealing (Must do)
- Make the visualizations be able to convey much more information like making the color of the country red go darker to lighter depending on the subscriber count in a World Map (Must do)
- Adding tasteful transitions and effects to our visualizations (optional)

This section and the entire process book will keep getting updated as we add more visualizations and features to the webpage. Thank you!