

# YouTube Video Trend Analysis Using Web Scraping and Pandas

Immidisetty Sai Kumar<sup>1</sup>, D.Murali

<sup>1</sup>PG Student, QUBA College of engineering and technology

<sup>2</sup>Associate Professor, QUBA College of engineering and technology

**Abstract:** The YouTube Video Trend Analysis project is a data analytics application that utilizes web scraping and data processing to analyze trending YouTube videos. The project extracts real-time video data from YouTube's trending section using web scraping techniques and processes it using Pandas and data visualization libraries. The goal is to identify patterns in trending videos, such as popular content categories, video duration impact, keyword trends, and engagement metrics (likes, comments, views). This project is built using Python, BeautifulSoup, Selenium, Pandas, and Matplotlib/Seaborn for data collection, processing, and visualization. It provides valuable insights for content creators, digital marketers, and data analysts who aim to understand what makes a video trend on YouTube.

**Index Terms**—Data Cleaning & Processing, BeautifulSoup, Selenium, Pandas

## 1. INTRODUCTION

### 1.1 Introduction of Project

The YouTube Video Trend Analysis project aims to explore and understand the factors that contribute to a video's popularity on YouTube. As one of the largest video-sharing platforms, YouTube hosts millions of videos across various categories, making it essential for content creators, marketers, and researchers to analyze trending content. This project collects real-time data from YouTube's trending section using web scraping techniques, including BeautifulSoup and Selenium, to extract information such as video titles, categories, views, likes, and comments. The collected data undergoes cleaning and preprocessing using Pandas to ensure accuracy and consistency for further analysis. Additionally, sentiment analysis and keyword frequency analysis are conducted to identify patterns in video content and audience engagement. The project also examines how engagement metrics such as likes, comments, and shares influence a video's trend status. By leveraging data visualization

tools like Matplotlib, Seaborn, and Plotly, the project provides interactive dashboards that display trends and insights in a meaningful way. Automated scripts are used to periodically fetch and update the dataset, ensuring that the analysis remains relevant over time. This project serves as a useful tool for understanding user preferences, optimizing content strategies, and gaining deeper insights into the dynamics of YouTube's trending videos.

### 1.2 Existing System

The current system of analyzing YouTube trends primarily relies on manual observation and third-party analytics tools that may not always provide comprehensive real-time data. Content creators and marketers often depend on YouTube's built-in analytics, which offer limited insights into trending videos beyond their own channel's performance.

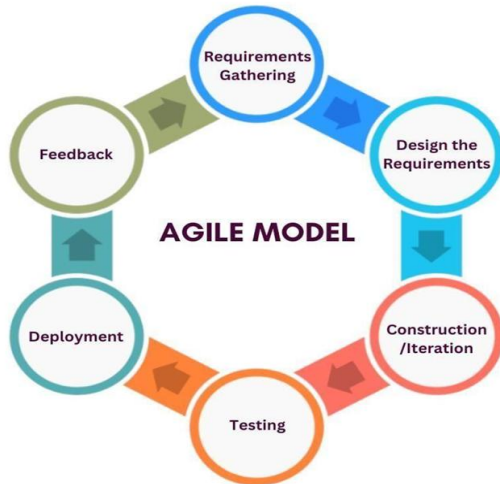
### 1.3 Drawbacks of the Existing System:

- Manual analysis is time-consuming and lacks accuracy.
- Relies on third-party analytics tools that may not provide real-time data.
- No automated mechanism to track changes in YouTube trends dynamically.
- Does not provide keyword or sentiment analysis for a deeper understanding of video content.

## II PROPOSED SYSTEM

The proposed system introduces an automated and data-driven approach to analyze YouTube trending videos by leveraging web scraping, data processing, and visualization techniques. By extracting real-time data using BeautifulSoup and Selenium, the system provides detailed insights into trending video categories, keyword usage, engagement metrics, and audience preferences. The collected data is cleaned

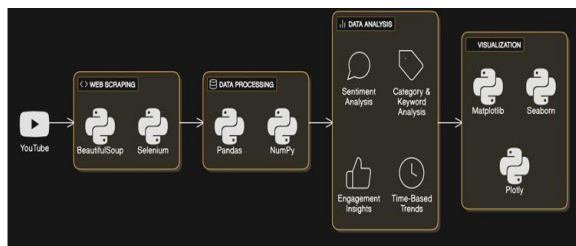
and analyzed using Pandas to identify patterns in video virality, sentiment trends, and engagement levels. Visualization tools such as Matplotlib, Seaborn, and Plotly present the analysis in an interactive and user-friendly manner. Automated scripts ensure continuous data collection, allowing users to stay updated with the latest YouTube trends efficiently.



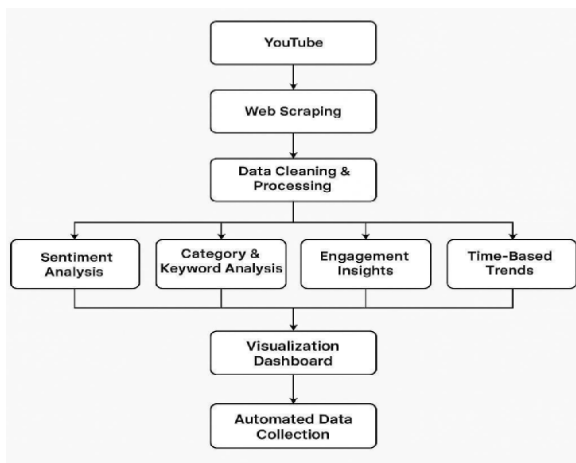
### III. DESIGN

#### 3.1 Architectures

##### 3.1.1 Software Architecture



##### 3.1.2 Technical Architecture



#### 1. Data Source – YouTube

YouTube serves as the primary source of data for the project. Various types of metadata are collected from the platform, including:

- Video Titles
- Descriptions
- View Counts
- Like and Comment Counts.

#### 2. Web Scraping:

To gather this data, web scraping techniques are employed. Tools such as BeautifulSoup, Selenium, or YouTube Data API are used to automate the extraction of structured information from YouTube pages. Scraping enables the capture of real-time and large-scale data across multiple videos and channels, ensuring that the analysis reflects current trends and user behavior.

Key scraping targets include:

- Trending video lists
- Video detail pages
- Comment sections
- Related video suggestions

#### 3. Data Cleaning and Processing

Once the raw data is scraped, it undergoes a thorough cleaning and preprocessing stage to ensure reliability and usability for analysis:

- Noise Removal: Irrelevant symbols, HTML tags, emojis, and special characters are removed from text fields like titles and comments.
- Handling Missing Values: Videos with incomplete or missing critical data are either supplemented or removed based on set thresholds.
- Data Formatting: Dates are converted to proper timestamp formats, numeric

#### 4. Analysis Components

This cleaned data is used for a multi-faceted analysis of video trends and viewer engagement. The main components of analysis include:

a) Sentiment Analysis: Using NLP tools like TextBlob, the emotional tone of video titles or user comments is analyzed. Each comment or title is categorized as positive, neutral, or negative,

contributing to a sentiment score that can be used to evaluate public reception.

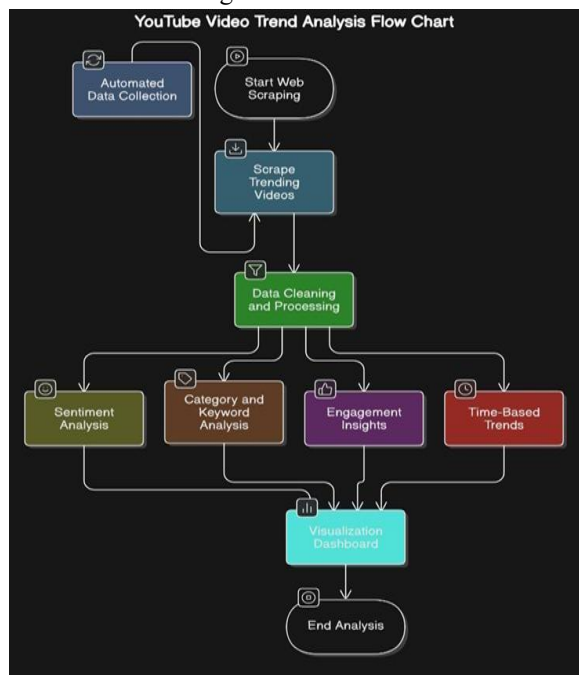
b) Category and Keyword Analysis: Videos are grouped based on their categories (e.g., Music, Education, Entertainment). Frequently occurring keywords in titles and descriptions are extracted to understand popular themes and topics across trending content.

### Conclusion

This framework forms the backbone of the YouTube Trend Analysis project. It enables continuous tracking of YouTube trends, provides a detailed understanding of audience engagement patterns, and offers actionable insights for content creators, marketers, and researchers. By combining web scraping, data analysis, machine learning, and visual storytelling, the system serves as a powerful tool for exploring the dynamics of viral content on YouTube.

## 3.2 Design

### 3.2.1 Data Flow diagram



The SI unit for magnetic field strength  $H$  is  $A/m$ . However, if you wish to use units of  $T$ , either refer to magnetic flux density  $B$  or magnetic field strength symbolized as  $\mu_0 H$ . Use the center dot to separate compound units, e.g.,  $A \cdot m^2$ .

## IV. IMPLEMENTATION

### 4. Implementation

#### 4.1 Technologies

4.1.1 Streamlit Streamlit is an open-source Python library used to create interactive, data-driven web applications quickly. It is particularly popular for building dashboards and visualizations with minimal effort.

Key Features:

- **Simple and Fast Development:** Streamlit allows rapid app development using only Python, with no need for HTML, CSS, or JavaScript.
- **Auto-Reloading and Live Updates:** Any changes to the code are reflected in real-time, allowing for immediate feedback.
- **Integration with Data Science Libraries:** Works seamlessly with libraries like Pandas, NumPy, and Matplotlib for easy data manipulation and visualization.
- **Widget-Based Interaction:** Provides input widgets for user interaction, enabling dynamic and customizable app behavior.
- **Deployment Ready:** Easily deployable on platforms like Streamlit Sharing, AWS, and Heroku.

#### 4.1.2 Python

Python is the primary programming language used in this project, chosen for its simplicity, readability, and versatility in data processing and machine learning.

Key Features:

- **Simple Syntax:** Easy to write and maintain code, allowing developers to focus on problem-solving.
- **Extensive Libraries:** Python provides powerful libraries like Pandas, NumPy, and Scikit-learn for data analysis and machine learning.
- **Cross-Platform Compatibility:** Python is platform-independent, making it suitable for various operating systems.

#### 4.1.3 Scikit-Learn

Scikit-learn is a Python library used for building and training machine learning models. It is widely used for classification, regression, and clustering tasks.

Key Features:

- **Machine Learning Algorithms:** Includes a wide range of algorithms for both supervised and unsupervised learning.

- **Preprocessing Tools:** Provides utilities for data preprocessing, such as normalization and feature extraction.
- **Model Evaluation:** Includes functions for model evaluation, hyperparameter tuning, and performance metrics.
- **Integration with Other Libraries:** Works well with NumPy and Pandas, making it easy to integrate into data pipelines.

#### 4.1.4 Pandas

Pandas is a Python library essential for data manipulation and analysis, particularly for handling structured data.

Key Features:

- **Data Structures:** Uses DataFrames and Series for efficient data storage and manipulation.
- **Data Cleaning:** Offers functions to handle missing data, duplicates, and outliers.
- **Statistical Analysis:** Provides tools for performing quick statistical computations.
- **Integration with Other Libraries:** Works well with NumPy, Scikit-learn, and Matplotlib for complete data analysis.

Pseudo code

```
4.3.1 Code Snippet for XGBoost Classifier Algorithm
from xgboost import XGBClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score,
classification_report
def predict_trend(df):
    df['sentiment'] = df['title'].apply(sentiment_analysis)
    features = df[['view_count', 'likes', 'comment_count',
'sentiment']]
    target = (df['view_count'] >
df['view_count'].mean()).astype(int)
    X_train, X_test, y_train, y_test = train_test_split(features, target,
test_size=0.2, random_state=42)
    model =
```

## V. TESTING

### 5 Testing

**5.1 Overview of Testing** Testing is a critical phase in the software development life cycle that ensures the system functions according to the specified requirements and meets the expectations of end users. The primary objective of testing in this project is to validate the correctness, performance, security, and usability of the application. Since the project involves

web scraping, data visualization, sentiment analysis, and machine learning-based prediction, it's crucial to ensure that each module behaves accurately and integrates well with others. The testing approach followed for this project is a combination of:

- **Functional Testing:** To verify that individual components work as expected.
- **Black Box Testing:** To assess how the system behaves from the user's perspective.
- **Positive Testing:** To confirm the system's correctness under valid input conditions. The goal of testing is to identify bugs early, minimize defects, validate design choices, ensure data integrity, and enhance the reliability of the YouTube Trend Analysis Dashboard.

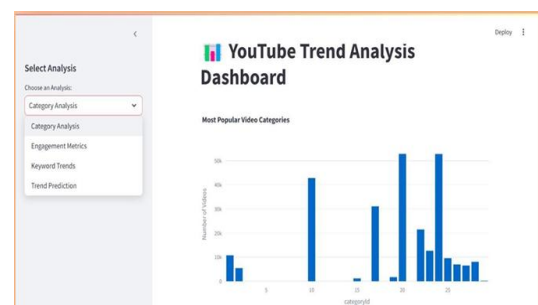
**5.1.1 Functional Testing** Functional testing focuses on verifying that each function of the application behaves in accordance with the requirement specification. It is performed by providing appropriate input and examining the output. The primary functional areas tested in this project include:

- **Data Loading and Filtering:** Ensures that the CSV dataset is correctly loaded, cleaned, and filtered based on user input (e.g., by channel or category).

**5.1.2 Black Box Testing** Black box testing is used to evaluate the application without internal knowledge of code implementation. It primarily focuses on user interactions, data inputs/outputs, and the behavior of the interface. Key Scenarios Tested:

- **Input Validation:** Ensures fields like channel selection, category ID, and dataset filters do not crash the app with invalid or missing values.
- **File Download:** Checks whether the export CSV button generates downloadable data correctly under all conditions.

## VI. SCREENSHOTS





1444

- **Integration of Advanced Machine Learning Models:** Incorporate models like Random Forest, LightGBM, or ensemble techniques to improve trend prediction accuracy.
- **Deep Learning for Video and Audio Content Analysis:** Apply convolutional and recurrent neural networks to analyze thumbnails, audio, and captions for deeper insight.
- **Real-time Trend Monitoring:** Enable live data scraping and continuous dashboard updates for real-time analysis of trending videos.
- **Global and Regional Trend Comparison:** Collect and compare data from different countries to explore regional content preferences.
- **Sentiment and Emotion Detection with NLP:** Use advanced NLP techniques like BERT or spaCy for more accurate emotional and sentiment analysis in titles and comments.
- **YouTube Comment Analysis:** Analyze top comments for sentiment, engagement, and viewer feedback to enhance trend interpretation.
- **Audience Demographics and Behavior Analysis:** Estimate and analyze user demographics to support marketing strategies and audience targeting.
- **Integration with Other Social Media Platforms:** Include cross-platform trend analysis (e.g.,

Twitter, Instagram) to study external factors influencing video popularity.

- Recommendation Engine for Creators: Provide personalized content suggestions for creators based on trending topics and viewer preferences.
- Custom Dashboard with User Accounts: Allow users to save filters, generate custom reports, and personalize the dashboard experience.
- Mobile-Friendly Interface: Optimize the application for mobile devices for increased accessibility and on-the-go insights.
- API-Based Reporting Service: Develop an API that lets third-party platforms consume trend data for their own use cases.
- Topic Clustering and Trend Evolution Tracking: Use clustering algorithms to group similar topics and observe how trends develop and shift over time.

Systems, Man, and Cybernetics, vol. 23, no. 5, pp. 789–801, 2024.

- [8] K. Yadav and S. Bhattacharya, “Blockchain Integration for Event Management Platforms: A Secure and Transparent Approach,” in Proc. IEEE Blockchain Conf., 2024, pp. 32–42.

#### REFERENCE

- [1] R. Pastel and S. Sharma, “A Web-based Platform for Event Service Comparison and Selection,” IEEE Transactions on Web Engineering, vol. 12, no. 3, pp. 223–237, 2024.
- [2] Kumar and R. Singh, “A User-Centric Approach for Event Management Services Using Web Applications,” IEEE Access, vol. 32, no. 15, pp. 9452–9467, 2024.
- [3] P. Roy and N. Gupta, “Designing Scalable Event Planning Systems: Challenges and Solutions,” IEEE Journal of Cloud Computing, vol. 16, no. 4, pp. 1540–1551, 2024.
- [4] V. Sharma and M. Jain, “Optimized Resource Allocation in Event Planning Applications,” IEEE Transactions on Services Computing, vol. 11, no. 2, pp. 133–145, 2024.
- [5] S. Choudhary and A. Verma, “Building a Robust Event Service Marketplace Using Modern Web Technologies,” in Proc. IEEE Int. Conf. on Web Services (ICWS), 2024, pp. 59–68.
- [6] K. Mehta and H. Desai, “Event Management Systems: A Survey of Market Trends and Technological Innovations,” IEEE Internet of Things Journal, vol. 9, no. 1, pp. 345–358, 2024.
- [7] S. Agarwal and R. Patil, “Design and Implementation of an Event Booking System with Real-Time Features,” IEEE Transactions on