

YouTube Comment Analyzer

Ms. Shreeya Pandagale¹, Ms. Vaishnavi Sankpal², Ms. Arya Deshmukh³ Mrs. Namrata Madavi⁴

Students, Department of Computer Technology^{1,2,3}

Lecturer, Department of Computer Technology⁴

Bharati Vidyapeeth Institute of Technology, Navi Mumbai

Abstract: *The exponential growth of user-generated content on digital platforms such as YouTube has resulted in an enormous volume of textual comments expressing public opinions, emotions, and reactions. Analyzing such large-scale unstructured data manually is inefficient, time-consuming, and impractical. Sentiment Analysis, a sub-field of Natural Language Processing (NLP), enables automated extraction of emotional polarity from textual content.*

This paper presents a Transformer-based YouTube Comment Sentiment Analyzer that extracts comments using the YouTube Data API, preprocesses noisy textual data, and classifies sentiments into Positive, Negative, and Neutral categories using a fine-tuned RoBERTa model. The proposed system integrates advanced NLP preprocessing techniques such as tokenization, stop word removal, lemmatization, spam filtering, and normalization.

The system generates statistical summaries and graphical sentiment distributions to support decision-making. Experimental observations indicate that transformer-based models significantly improve contextual understanding compared to traditional machine learning techniques. The proposed solution demonstrates scalability, efficiency, and real-world applicability in social media analytics.

Keywords: Sentiment Analysis, Natural Language Processing, RoBERTa, Transformer Models, YouTube Data API, Social Media Analytics

I. INTRODUCTION

The rapid advancement of digital communication platforms has transformed the way individuals express opinions and interact online. Among these platforms, YouTube stands as one of the most widely used video-sharing platforms globally.

Millions of videos are uploaded daily, and each video receives hundreds or thousands of comments from viewers.

These comments serve as valuable indicators of audience perception, satisfaction levels, and emotional responses.

However, the massive volume of comments generated daily makes manual analysis infeasible. Businesses, content creators, and researchers require automated systems to extract meaningful insights from textual data efficiently. Sentiment analysis, also known as opinion mining, enables classification of text into emotional categories such as positive, negative, and neutral.

Earlier sentiment analysis techniques relied on traditional machine learning algorithms such as Naïve Bayes, Support Vector Machines (SVM), and Logistic Regression. These methods required manual feature engineering using techniques like Bag-of-Words (BoW) and TF-IDF. Although effective to some extent, they lacked contextual understanding.

Recent advancements in deep learning, particularly transformer-based models like RoBERTa, have significantly enhanced text classification accuracy by capturing semantic and contextual relationships through attention mechanisms.

This paper proposes a YouTube Comment Sentiment Analysis System leveraging a transformer-based architecture to perform accurate and scalable sentiment classification.



Problem being addressed:

YouTube comments are highly unstructured and noisy in nature. They frequently include:

- o Informal language and slang
- o Emojis and special characters
- o URLs and spam content
- o Repeated characters
- o Mixed-language text
- o Sarcasm and ambiguous tone

II. LITERATURE SURVEY

Early research in sentiment analysis primarily focused on supervised machine learning techniques. Pang and Lee demonstrated that Naïve Bayes and SVM could classify movie reviews with moderate accuracy. However, these models struggled with contextual variations and sarcasm. Deep learning approaches such as Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM) networks improved sequential modeling of text. Nevertheless, these models faced limitations in handling long-range dependencies and required substantial computational resources. Transformer-based models such as BERT and RoBERTa introduced self-attention mechanisms that significantly enhanced contextual text representation. RoBERTa, a robustly optimized version of BERT, achieved state-of-the-art performance in multiple NLP tasks due to improved training strategies and larger datasets. Several studies have analysed sentiment on platforms like Twitter and Facebook. However, limited research integrates transformer-based models specifically for YouTube comment analytics with structured visualization output.

III. METHODOLOGY

The proposed YouTube Comment Sentiment Analyzer follows a structured approach consisting of data collection, preprocessing, model training, sentiment classification, and result generation. First, comments are extracted using the YouTube Data API by providing a video link. The system also allows CSV file input for previously stored datasets. This ensures flexibility in analysing both live and offline data.

Since YouTube comments are often unstructured and noisy, preprocessing is performed to improve classification accuracy. This includes removal of URLs and special characters, spam filtering, tokenization, stopword removal, and lemmatization. These steps convert raw textual data into a clean and structured format suitable for model input.

For sentiment classification, a transformer-based model, RoBERTa (Robustly Optimized BERT Pretraining Approach), is used. The model is fine-tuned on a labelled dataset containing positive, negative, and neutral comments. Transformer models utilize self-attention mechanisms to capture contextual relationships between words, enabling better understanding of semantic meaning compared to traditional machine learning techniques.

After training, the model classifies each comment into one of three categories: Positive, Negative, or Neutral. The system then calculates the percentage distribution of each sentiment category and generates a summarized analytical output.

Finally, the results are presented in a structured format to help users understand overall audience sentiment efficiently.

IV. IMPLEMENTATION

1. Backend infrastructure- The backend infrastructure of the YouTube Comment Sentiment Analyzer is responsible for handling data extraction, preprocessing, sentiment classification, and result generation. The system integrates the YouTube Data API to retrieve comments from a given video link. A preprocessing module cleans and normalizes the extracted data by removing unwanted elements and preparing the text for analysis. The RoBERTa transformer-based model is implemented within the backend to classify comments into Positive, Negative, and Neutral categories. The backend also computes statistical summaries such as total comments analysed and sentiment percentage distribution before sending the results to the user interface.



2. User interface design The user interface is designed to be simple, interactive, and user-friendly so that users can easily operate the system without technical knowledge. It allows users to enter or paste a YouTube video link into the input field provided on the interface. After submitting the link, the system processes the video data using sentiment analysis techniques. Once the processing is completed, the interface displays the sentiment results in a well-structured format, such as positive, negative, and neutral categories.

In addition to textual results, the system also presents graphical summaries like charts or graphs to help users visualize the overall audience sentiment.

The design focuses on clarity and accessibility, ensuring that users can easily interpret the analysis results.

This approach helps users quickly understand public opinion and audience reactions to the selected YouTube video without requiring any technical expertise.

3. General Flow Chart

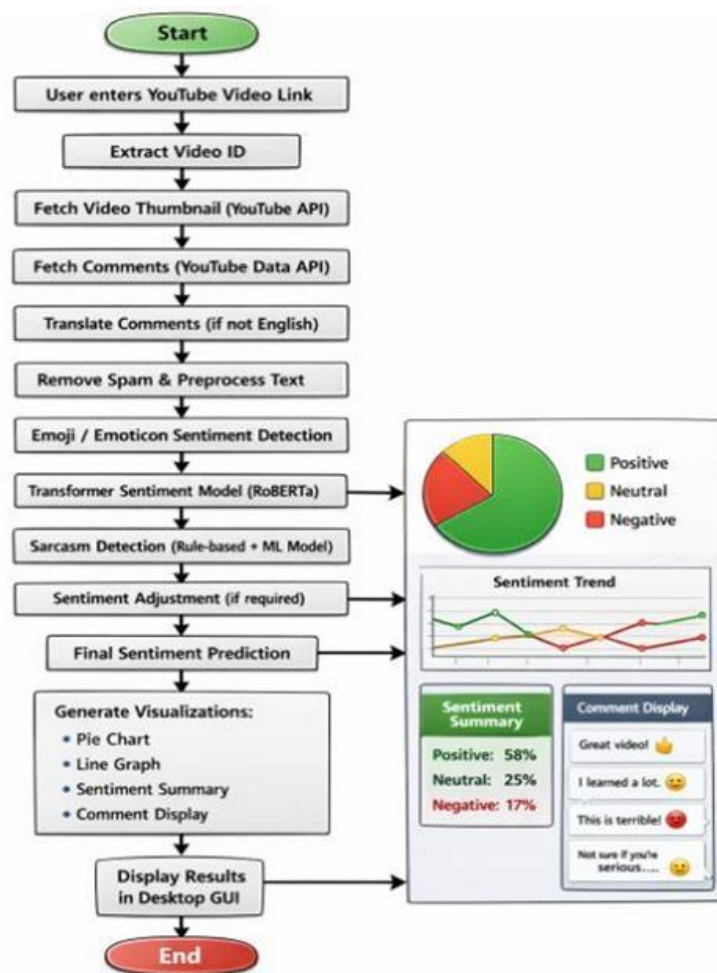


Fig1. Flow system of comment analyzer



4. Overview

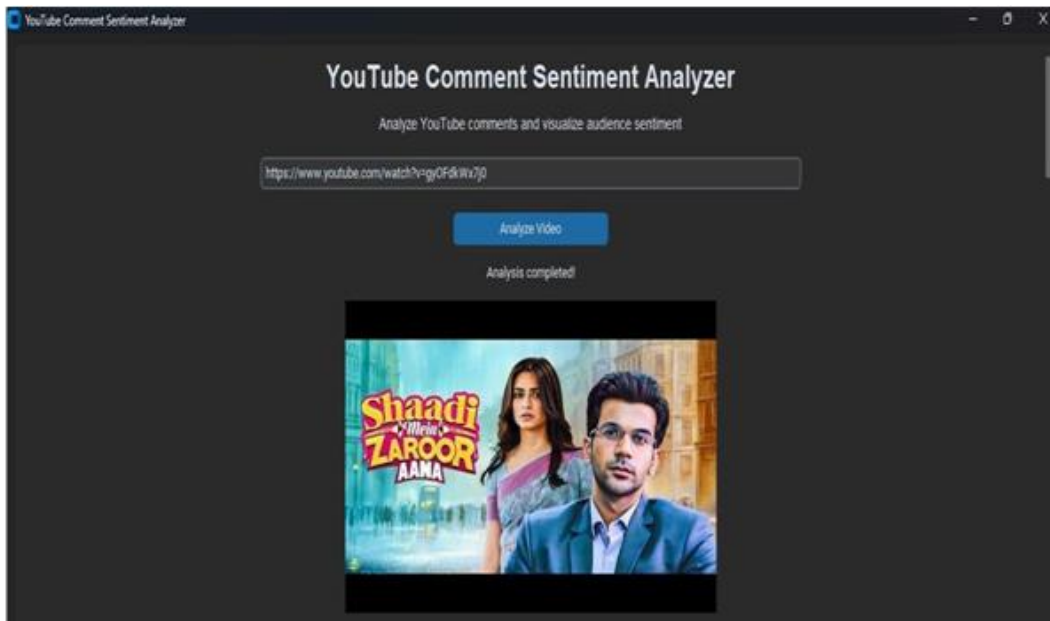


Fig. 2 Input Section Page



Fig.3. Sentiment Graph Page.



- Emoji and Emoticon Awareness.
- Multilingual Comment support.
- Youtube specific Spam Filtering.
- Interactive Desktop-Based user Interface.
- Scalable and Extensible Design.

Weakness:

- High Computational requirements.
- Dependency on internet connectivity.
- Limited Comment Volume per Request.
- May have difficulties in translation of languages.

VI. FUTURE SCOPE

Integration of Advanced Transformer Models

Future work can explore more advanced transformer architectures such as BERT, RoBERTa-large, or domain-specific fine-tuned models to further improve sentiment classification accuracy. Larger models trained on YouTube-specific datasets could better capture informal language and contextual nuances.

Real-Time Sentiment Monitoring: The system can be extended to support real-time sentiment analysis by continuously monitoring new comments using the YouTube Data API. This would allow live sentiment tracking for trending videos, live streams, and ongoing events.

Deployment as a Web Application: Currently implemented as a desktop application, the system can be deployed as a web-based platform using frameworks such as Flask or Django. This would improve accessibility and allow users to analyze videos directly through a browser

Scalability for Large-Scale Data: Future enhancements may include handling thousands of comments instead of a fixed limit (e.g., 100 comments). Implementing pagination and batch processing techniques would enable large-scale sentiment analysis

Improved Sarcasm and Context Detection: Sarcasm detection can be further enhanced using advanced contextual models trained specifically on irony-rich datasets. This would reduce incorrect sentiment flipping and improve reliability.

Emotion-Level Classification: Instead of limiting classification to positive, negative, and neutral, future versions can classify comments into more detailed emotional categories such as happiness, anger, sadness, surprise, and frustration.

Engagement and Analytics Integration: The system can be expanded to analyze correlations between sentiment and video engagement metrics such as likes, replies, view count, and subscriber growth, providing deeper insights for content creators and marketers.

Offline Model Optimization: Future improvements may include optimizing the transformer model for lightweight deployment using techniques such as model distillation or quantization, reducing computational requirements and improving execution speed.

Multilingual Native Model Support: Instead of translating comments into English, future implementations can use multilingual transformer models to directly classify comments in different languages, improving linguistic accuracy.

Optional Closing Line (You can include this): Overall, the proposed system provides a strong foundation for intelligent sentiment analysis of YouTube comments and can be extended into a scalable, real-time, and production-ready analytical platform



VII. CONCLUSION

The YouTube Comment Sentiment Analyzer successfully demonstrates the application of transformer-based models for analyzing user-generated content on social media platforms. The system effectively extracts, preprocesses, and classifies YouTube comments into Positive, Negative, and Neutral categories using contextual sentiment analysis.

By integrating a structured backend infrastructure and a user-friendly interface, the system provides clear and meaningful insights into audience reactions. The implementation highlights the advantages of transformer models in understanding contextual relationships within textual data, resulting in improved sentiment classification accuracy compared to traditional approaches.

Although the current system operates in batch mode, it provides a scalable and reliable framework for large-scale comment analysis. With future enhancements such as real-time processing and live chat sentiment analysis, the system has strong potential for broader applications in content evaluation, brand monitoring, and audience engagement analysis.

Overall, the project establishes an efficient and practical solution for automated sentiment analysis of YouTube comments

REFERENCES

- [1] Yinhan Liu, Myle Ott, Naman Goyal, et al., “RoBERTa: A Robustly Optimized BERT Pretraining Approach,” arXiv preprint arXiv:1907.11692, 2019.
- [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” Proceedings of NAACL-HLT, 2019.
- [3] Ashish Vaswani, Noam Shazeer, Niki Parmar, et al., “Attention Is All You Need,” Advances in Neural Information Processing Systems (NeurIPS), 2017
- [4] Alec Radford, et al., “Language Models are Unsupervised Multitask Learners,” OpenAI Technical Report, 2019
- [5] Bo Pang and Lillian Lee, “Opinion Mining and Sentiment Analysis,” Foundations and Trends in Information Retrieval, vol. 2, no. 1–2, pp. 1–135, 2008.
- [6] Bing Liu, “Sentiment Analysis and Opinion Mining,” Synthesis Lectures on Human Language Technologies, 2012.
- [7] Steven Bird, Ewan Klein, and Edward Loper, Natural Language Processing with Python, O’Reilly Media, 2009.
- [8] Thomas Wolf, Lysandre Debut, Victor Sanh, et al., “Transformers: State-of-the-Art Natural Language Processing,” Proceedings of EMNLP: System Demonstrations, 2020.
- [9] Google Developers, “YouTube Data API Documentation,” Available: <https://developers.google.com/youtube>
- [10] Alec Go, Richa Bhayani, and Lei Huang, “Twitter Sentiment Classification using Distant Supervision,” Stanford University Technical Report, 2009.

