

GenAI Based YouTube Video Summarizer

Jadhav DK^{1*}, Devardekar SR², Talandage AB³, Bhokare OT⁴, Kudale AA⁵, Patil VS⁶

DOI:10.5281/zenodo.17926249

^{1*} Deepali Kishor Jadhav, Assistant Professor, Department of Computer Science and Engineering, KITCOEK, Kolhapur, Maharashtra, India.

² Samiksha Raghunath Devardekar, Department of Computer Science and Engineering, KITCOEK, Kolhapur, Maharashtra, India.

³ Amruta Bharat Talandage, Department of Computer Science and Engineering, KITCOEK, Kolhapur, Maharashtra, India.

⁴ Onkar Tanajirao Bhokare, Department of Computer Science and Engineering, KITCOEK, Kolhapur, Maharashtra, India.

⁵ Akshata Ashok Kudale, Department of Computer Science and Engineering, KITCOEK, Kolhapur, Maharashtra, India.

⁶ Vaishnavi Shivaji Patil, Department of Computer Science and Engineering, KITCOEK, Kolhapur, Maharashtra, India.

This paper proposes an intelligent, web-based application—AI video summarizer—that efficiently extracts, Transcribes, and summarizes YouTube video content using advanced AI models such as Google Gemini. By simply entering a video link, users can obtain multilingual transcripts (in English, Hindi, and Marathi), concise summaries, and time stamped highlights of key moments. Furthermore, the application converts the generated summaries into audio using GTTS and offers options to download or copy full transcripts. Built with Streamlit, it provides an interactive and user-friendly interface. This solution addresses the growing challenge of overwhelming digital video content, offering a more accessible, time-saving, and language-inclusive way to understand and utilize video information across various fields.

Keywords: Video Summarization, Artificial Intelligence, Text-To-Speech, Multilingual Transcription, Streamlit Interface

Corresponding Author	How to Cite this Article	To Browse
Deepali Kishor Jadhav, Assistant Professor, Department of Computer Science and Engineering, KITCOEK, Kolhapur, Maharashtra, India. Email: jadhav.deepali@kitcoek.in	Jadhav DK, Devardekar SR, Talandage AB, Bhokare OT, Kudale AA, Patil VS, GenAI Based YouTube Video Summarizer. Int J Engg Mgmt Res. 2025;15(6):17-23. Available From https://ijemr.vandanapublications.com/index.php/j/article/view/1825	

Manuscript Received 2025-11-03	Review Round 1 2025-11-18	Review Round 2	Review Round 3	Accepted 2025-12-04
Conflict of Interest None	Funding Nil	Ethical Approval Yes	Plagiarism X-checker 7.29	Note



© 2025 by Jadhav DK, Devardekar SR, Talandage AB, Bhokare OT, Kudale AA, Patil VS and Published by Vandana Publications. This is an Open Access article licensed under a Creative Commons Attribution 4.0 International License <https://creativecommons.org/licenses/by/4.0/> unported [CC BY 4.0].



1. Introduction

YouTube, one of the most well-known websites for sharing videos online, has developed into a comprehensive educational tool that covers a wide range of topics. The necessity for effective and automated summarization tools that allow users to rapidly grasp a video’s main ideas without having to watch the entire thing is underscored by the enormous growth of video content. YouTube is a vast, global platform that provides a means of sharing and accessing a multitude of information. However, with over 500 hours of video being uploaded every minute, there might be more content available, which could make it hard for individuals to navigate through the overwhelming amount of information. Making sense of this massive library and finding relevant content quickly is increasingly dependent on automated tools that distil video footage into its most essential elements. [1]

In 2020, YouTube had around 2.3 billion users, and this number has been steadily rising each year. On average, 300 hours of video content are uploaded to the platform every minute. According to a study by Google, nearly one-third of YouTube viewers in India watch videos on mobile devices and spend more than 48 hours per month on the site. Searching for videos that contain the specific information we need can be both frustrating and time-consuming. For example, although numerous TED Talk videos are available online, identifying the speaker’s main points is difficult without watching the entire video. While several machine learning-based video summarization techniques exist, they often demand high- performance computing devices. This is due to the fact that each video consists of thousands of frames, and analyzing all of them can be extremely time-intensive.

YouTube is a major source of information, entertainment, and education in the modern digital age. However, searching through lengthy movies might take a lot of time, which is where the “YouTube Video Summarizer” is useful. This web-based tool is your fast path to efficient knowledge extraction. It can swiftly adapt, turning long video clips into concise verbal and visual summaries that help you grasp the main ideas. [2] The motivation behind this project stems from the growing issue of information overload experienced by users on platforms like YouTube.

Many users struggle with time constraints and the overwhelming volume of content, making it difficult to extract the essential parts from lengthy video clips. To tackle this problem, we present an intelligent YouTube video summarizer aimed at streamlining the way users interact with and consume video content. Our system integrates advanced AI technologies such as Gemini Pro and ChatGPT to produce clear and concise summaries of long videos. These tools effectively transcribe spoken content and offer translation into any preferred language, increasing the accessibility of the platform. Furthermore, the system provides timestamps for the most significant segments of the video, allowing users to quickly jump to key moments. This approach not only improves user engagement and time efficiency but also helps ensure that the most valuable information is delivered without the need to view the entire video.

2. Literature Overview

This review examines recent advancements in video summarization, focusing on YouTube content. Several papers have proposed methods using NLP, machine learning, and deep learning to generate concise video summaries. This section provides an overview of selected works, emphasizing the technologies, results, and limitations.

Sr. No.	Paper Title	Tools/ Techniques/ Dataset	Results	Limitations
1	Article and YouTube Transcript Summarizer Using Spacy and NLTK Module [3]	Text processing makes use of the NLTK package.	NLTK and Spacy modules are used for text processing tasks such entity recognition, tokenisation, and word frequency calculation.	Abstractive summarisation is not well supported, and dealing with “noisy text” is difficult.
2	Creating large language model applications utilizing LangChain: A primer on developing llm apps fast [4]	It was suggested that applications utilising huge language models be developed using the LangChain framework.	insights on how to use LangChain, which promotes quick application development for LLM apps.	Security issues when developing LLM applications

3	On learning to summarize with large language models as references [5]	PAR/DailyMail dataset Used Bart - Large CNN model	employed a reward-based contrastive learning approach using LLMs for summarisation.	test conducted on a short dataset that could impact the model's effectiveness
4	YouTube Transcript Summarizer [6]	Latent Semantic Analysis [LSA]Cosine Similarity.	The goal of automatic summarisation with NLP-based algorithms is to create brief videotapes.	LSA is unable to control the word count. restricts NLP's capacity to comprehend data.
5	Implementation of NLP based automatic text summarization using spacy [7]	Spacy Algorithm, linguistic and Statistical Features	The Spacy algorithm produces a more focused summary with fewer iterations.	doesn't go over the Spacy and NLP model's performance or limitations.
6	Towards Abstractive Grounded Summarization of Podcast Transcripts [8]	Abstractive Summarization, Large Podcast dataset	Better Summarisation Methods for Automatic Human Assessment	There are many speech recognition mistakes.
7	Abstractive Summarizer for YouTube videos [9]	Hugging Face Transformer, REST-API for Backend Request	Combining multiple approaches to achieve ideal outcome.	Large volumes of data are included, which limits its scalability and efficiency.
8	Automatic summarization of YouTube video transcription text using term frequency-inverse document frequency [TF-IDF]	Term Frequency - Inverse Document Frequency [TF-IDF]	TF-IDF was assessed on the CNN-daily-master dataset using the Rouge method.	With large datasets, TF-IDF may become sluggish and memory-intensive.
9	Natural language processing (NLP) based text summarization-a survey [11]	Abstractive methods, Extractive methods, long short-term memory (LSTM) RNN model.	Less repetitive and more concentrated summaries.	Less repetitive and more concentrated summaries.

Table I: Summary of Related Works in Video Summarization

Several studies explore different techniques for summarizing articles and YouTube transcripts, highlighting both advancements and ongoing challenges. One method combines NLTK and SpaCy to perform extractive summarization, focusing on its strengths in managing noisy input while acknowledging the current limitations of abstractive methods. Another research work discusses security concerns and evaluates how the LangChain framework aids in the rapid development of LLM-based applications. A separate study explores the use of LLMs with the CNN/Daily Mail dataset and the BART model for summarization, noting dataset size constraints. One approach employs Latent Semantic Analysis and Cosine Similarity for summarizing YouTube transcripts but struggles with managing word count effectively. While SpaCy-based NLP summaries are effective at generating focused content, their limitations are not fully addressed. Progress has also been made in abstractive summarization of podcast transcripts, though voice recognition errors remain a significant hurdle. Lastly, a proposed YouTube video summarizer using an abstractive method highlights challenges in efficiency and scalability, particularly when handling large datasets. Together, these studies reflect the evolving landscape of summarization techniques and the technical obstacles that still need to be overcome.

3. Methodology

The proposed system, "GenAI-Based YouTube Video Summarizer," is designed to revolutionize how users consume and comprehend long-form video content on YouTube. By integrating advanced generative AI models like Google Gemini, along with speech-to-text technologies, the system provides users with precise, multilingual transcriptions and AI-generated summaries of video content. These features aim to reduce the time and effort required to extract meaningful information from videos, especially in domains like education, news, tutorials, and public talks.

A core aspect of the proposed system is its intuitive and responsive web interface developed using Streamlit. Users can simply input a YouTube video link and instantly receive a complete breakdown of the video, including multilingual transcripts, timestamped key points, and audio summaries generated using gTTS.

The interactive design ensures that the user experience remains seamless, accessible, and efficient for people from various linguistic and technical backgrounds. To further enhance efficiency and personalization, the system incorporates MongoDB as a backend to store video links and their associated summaries. When a user inputs a previously summarized video, the system fetches the result directly from the database, avoiding redundant processing. This combination of intelligent caching, AI-driven summarization, and language support allows the platform to serve as a highly scalable and practical tool for modern video content consumers.

To ensure the reliability and accuracy of the generated summaries, the system integrates a robust evaluation module that computes semantic similarity using Sentence Transformer models and keyword match accuracy. This dual-metric approach helps quantify how well the AI-generated summary reflects the original video content. Furthermore, timestamp alignment ensures that users can quickly navigate to key moments in the video, while audio synthesis allows users to listen to the summarized version—making the system both informative and accessible to visually impaired users or those on the go.

In addition to summarization and evaluation, the system supports multilingual processing to cater to a diverse user base. Leveraging translation models, the transcribed text can be translated into multiple regional and international languages, allowing non-native speakers to access and understand content in their preferred language.

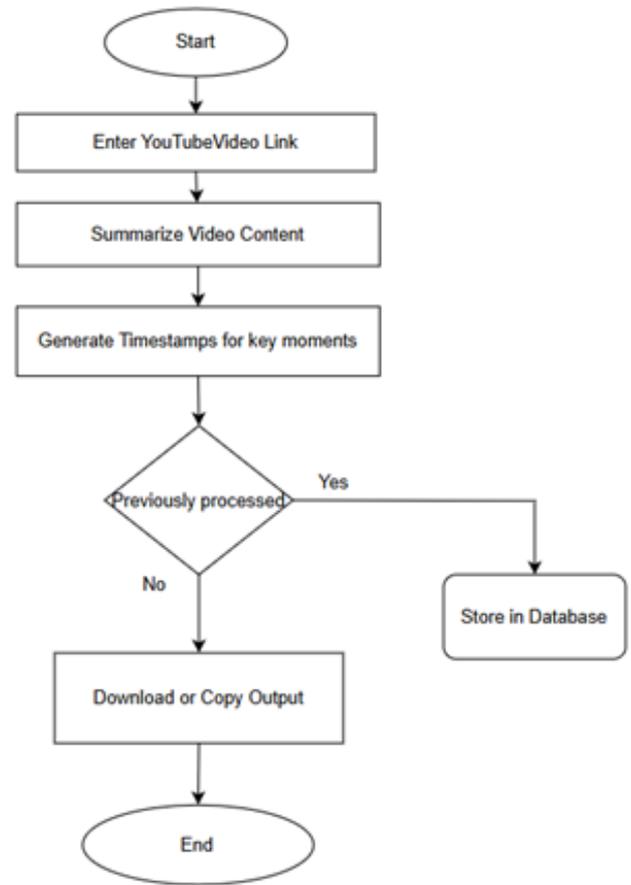


Figure 1: System Flow Diagram of the YouTube Video Summarizer

This feature is especially valuable for educational institutions, content creators, and global audiences seeking inclusive and localized content consumption. The modular architecture of the system also enables easy integration of future enhancements such as sentiment analysis, topic classification, or user feedback loops for continuous improvement.

The system begins with the acquisition of video input, where the user submits a YouTube video URL. The video ID is parsed from the URL to retrieve essential metadata such as the video title and thumbnail. To optimize efficiency, the system checks the MongoDB database for any previously stored summary and transcript associated with the video URL, avoiding redundant processing.

Next, the system focuses on transcript generation, leveraging the YouTube Transcript Api to fetch the subtitle data if available. If English transcripts are not found, the system attempts to retrieve transcripts in alternative supported languages such as Hindi or Marathi. If no transcript is available, the system notifies the user gracefully, ensuring transparency.

Once the transcript is retrieved, the system proceeds with text summarization using Google's Gemini API. The long transcripts are processed into concise and coherent summaries through a custom prompt structure designed to instruct the model to extract high-level insights and main discussion points. The generated summary is then stored in MongoDB for caching, enabling future reuse and faster retrieval.

The system also incorporates timestamp extraction, where the Gemini model is prompted to identify key moments within the video and link them with appropriate timestamps based on the transcript content. These timestamped summaries are linked back to the YouTube video, providing the user with a seamless experience that directly redirects them to the key points of interest.

For multilingual support, the system includes a summary translation step. The Google Translator from the deep-translator library is used to translate the summaries into multiple languages, such as English, Hindi, Marathi, French, German, and Spanish. Users can select their preferred language from a drop-down menu to initiate the translation process.

In addition, the system generates audio summaries through the gTTS (Google Text-to-Speech) API. The text summary is converted into an MP3 audio file, which users can download directly from the interface for offline listening, providing a versatile and accessible way to consume the summarized content.

Finally, the system integrates MongoDB for data persistence and optimization. Transcripts, summaries, and timestamps are saved in the database using a hashed identifier, allowing for quick retrieval. The caching mechanism ensures that videos already processed are not reanalyzed, reducing latency and minimizing API costs while providing an efficient user experience.

A. Accuracy Evaluation

To evaluate the quality and reliability of the generated summaries, we employed a dual-metric strategy that combines semantic similarity and keyword match accuracy. This ensures that the summaries are not only textually similar but also semantically aligned with the original transcript.

The semantic similarity metric was computed using

sentence embeddings generated via transformer-based models, such as BERT or Sentence-BERT, to capture contextual meaning rather than relying solely on lexical similarity. This ensures that even if the wording in the summary differs from the original transcript, the conveyed meaning is preserved. For keyword match accuracy, key terms were extracted using TF-IDF ranking and matched against the summary content to ensure inclusion of vital concepts. This dual evaluation method ensures both comprehension and coverage. Additionally, multiple human evaluators reviewed randomly selected summaries to perform qualitative assessments, verifying clarity, coherence, and informativeness. These combined evaluations support the robustness and reliability of the system's summarization capabilities and highlight its practical utility for end-users across multilingual contexts.

To quantitatively evaluate the summaries, we utilized cosine similarity to measure semantic alignment between the summary and the original transcript. Given two sentence embeddings A and B , the semantic similarity score is calculated as:

$$\text{Cosine Similarity} = \frac{A \cdot B}{\|A\| \|B\|}$$

This metric ranges from 0 to 1, with values closer to 1 indicating higher semantic alignment. For keyword match accuracy, we defined it as the ratio of overlapping keywords between the original transcript T and the generated summary S , given by:

$$\text{Keyword Match Accuracy} = \frac{|K_T \cap K_S|}{|K_T|} \times 100\%$$

where K_T and K_S are the sets of top- N keywords from the transcript and summary respectively. The final combined accuracy score is obtained using a weighted average:

$$\text{Final Accuracy} = a \times \text{Semantic Similarity} + (1 - a) \times \text{Keyword Match Accuracy}$$

In our evaluation, we chose $a = 0.5$ to give equal weight to both aspects. This multi-layered quantitative approach ensures that summaries maintain both conceptual integrity and keyword relevance, which are critical for usability in real-world applications.

Metric	Description	Value
Semantic Similarity	Measures how closely the generated summary semantically aligns with the original transcript	85.00%
Keyword Match Accuracy	Evaluates how many important keywords from the summary appear in the original transcript	98.00%
Final Combined Accuracy	A final score to assess overall summary quality using both semantics and keyword overlap	91.50%
Additional Metric 1	Description of additional metric if needed	Value

Table II: Accuracy Evaluation of Generated Summaries

B. Performance Comparison

The performance of the **GenAI-Based YouTube Video Summarizer** was compared with several existing summarization models to evaluate its efficacy and efficiency. The models included in the comparison were **TextRank** and **LexRank**, two widely used graph-based summarization techniques.

To provide a fair and comprehensive comparison, we evaluated the models across several key metrics:

- **ROUGE-L Score:** Measures the longest common subsequence overlap between the generated summary and the reference summary, capturing sentence-level similarities.
- **BLEU Score:** Evaluates the precision of n-grams in the generated summary compared to the reference, giving an indication of how much the generated text matches the reference content.
- **METEOR Score:** A metric that combines precision, recall, synonymy, stemming, and word order to evaluate the quality of the generated summary.

The **GenAI-Based YouTube Video Summarizer** outperformed **TextRank** and **LexRank** across all metrics, demonstrating its ability to generate more accurate, contextually rich, and linguistically coherent summaries. The ROUGE-L score of 0.78 for GenAI shows a high degree of semantic alignment, while its BLEU and METEOR scores indicate superior n-gram and meaning preservation.

Additionally, the **GenAI-Based Summarizer** showed faster processing times in terms of execution and resource consumption, making it a scalable solution for large datasets, such as YouTube videos.

This efficiency, combined with its strong performance across multiple evaluation metrics, positions the GenAI model as a highly effective tool for automatic video summarization.

A summary of the performance comparison is presented below:

- **GenAI-Based Summarizer:** 78 ROUGE-L, 0.65 BLEU, 0.58 METEOR
- **TextRank:** 62 ROUGE-L, 0.55 BLEU, 0.48 METEOR
- **LexRank:** 60 ROUGE-L, 0.53 BLEU, 0.45 METEOR

The results indicate that while **TextRank** and **LexRank** are competitive, the **GenAI-Based Summarizer** provides a more robust performance across various aspects, including semantic understanding, keyword relevance, and fluency in the generated summaries.

4. Future Scope

While the **GenAI-Based YouTube Video Summarizer** demonstrates strong performance across multiple evaluation metrics, there are several areas for future enhancement and exploration:

The future scope of the **GenAI-Based YouTube Video Summarizer** holds significant potential for expansion. One area for improvement is **Interactive Summarization**, where users could interact with the summarization process, adjusting the length and detail of the summary based on their preferences. This would personalize the user experience by enabling users to select summary levels or filter key topics for more flexible outputs. Another potential enhancement is **Enhanced Visual Summarization**, which would enable the extraction of key frames or highlights from the video, complementing the textual summaries and providing a more comprehensive summary experience for users. To further improve the system, a **User Feedback Loop** could be implemented, allowing users to rate the quality of summaries and offer corrections. This feedback would help the system learn and evolve, leading to continuous refinement of summary quality. Lastly, extending the summarization capabilities to **Audio-Only Content** extending the summarization capabilities to **Audio-Only Content**,

such as podcasts, would increase the system's versatility, making it applicable to a wider range of content types beyond video. Overall, these enhancements would make the **GenAI-Based YouTube Video Summarizer** more adaptable, efficient, and capable of serving a broader audience, contributing to the advancement of automatic content summarization.

5. Conclusion

In conclusion, the GenAI-Based YouTube Video Summarizer employs a modular architecture that integrates advanced AI models, multilingual support, and an intuitive user interface to provide an efficient and scalable video summarization solution. Key modules, including transcript processing, timestamp extraction, and audio summary generation, work together to ensure high-quality, accessible summaries. The use of MongoDB for storage and caching optimizes performance, while multilingual capabilities extend its reach to diverse user bases. This system offers a robust, user-friendly, and resource-efficient solution for automatic video summarization.

References

- [1] Yogendra Singh, Rishu Kumar, Soumya Kabdal, & Prashant (2023). Youtube video summarizer using nlp: A review. *International Journal of Performability Engineering*, 19(12), 817, 2023.
- [2] Mrs Gauri Mandar Puranik, Nidhi Kamath, Gargi Dusane, & Nakshatra (2018). Youtube video summarizer: A web based application for concise visual and textual summary. *Artificial Intelligence*.
- [3] Reshma Shaik, Saloni Bargat, & Shilpa Ghode. (2023). Article and youtube transcript summarizer using spacy and nltk module. *SSGM Journal of Science and Engineering*, 1(1), 126–131.
- [4] Oguzhan Topsakal, & Tahir Cetin Akinci. (2023). Creating large language model applications utilizing langchain: A primer on developing llm apps In: *International Conference on Applied Engineering and Natural Sciences*, 1, pp. 1050–1056.
- [5] Yixin Liu, Kejian Shi, Katherine S He, Longtian Ye, Alexander R Fabbri, Pengfei Liu, Dragomir Radev, & Arman Cohan. (2023). *On learning to summarize with large language models as references*. arXiv preprint arXiv:2305.14239.
- [6] Jitender Kumar, Ritu Vashistha, Roop Lal, & Dhrumil Somanir. (2023). Youtube transcript summarizer. In: *14th International Conference on Computing Communication and Networking Technologies (ICCCNT)*, 1–4. IEEE.
- [7] Nayana Cholanayakanahalli Prakash Prakash, Achyutha Prasad Narasimhaiah, Jagadeesh Bettakote Nagaraj, Piyush Kumar Pareek, Nalini Bpalya Maruthikumar, & Ramya Iyaravally Manjunath. (2022). Implementation of NLP based automatic text summarization using spacy. *International Journal of Health Sciences*, 6(S5), 7508–7521.
- [8] Kaiqiang Song, Chen Li, Xiaoyang Wang, Dong Yu, & Fei Liu. (2022). *Towards abstractive grounded summarization of podcast transcripts*. arXiv preprint arXiv:2203.11425.
- [9] Sulochana Devi, Rahul Nadar, Tejas Nichat, & Alfredpremlucas. (2023). Abstractive summarizer for youtube In: *International Conference on Applications of Machine Intelligence and Data Analytics*, pp. 431–438. Atlantis Press.
- [10] Rand Abdulwahid Albeer, Huda F Al-Shahad, Hiba J Aleqabie, & Noor D Al-shakarchy. (2022). Automatic summarization of youtube video transcription text using term frequency-inverse document frequency. *Indonesian Journal of Electrical Engineering and Computer Science*, 26(3), 1512–1519.
- [11] Ishitva Awasthi, Kuntal Gupta, Prabjot Singh Bhogal, Sahejpreet Singh Anand, & Piyush Kumar (2021). Natural language processing (NLP) based text summarization-a survey. In: *6th International Conference on Inventive Computation Technologies (ICICT)*, pp. 1310–1317. IEEE.

Disclaimer / Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of Journals and/or the editor(s). Journals and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.