A
Mini Project Report on

# YouTube Video Transcribe Summarizer

Submitted in partial fulfilment of the requirements
for the degree of
BACHELOR OF ENGINEERING
IN
**Computer Science & Engineering**
Artificial Intelligence & Machine Learning


by


Ved Sawant (24206001)
Prajwal Dhanawade (24206007)
Gaurav Kshirsagar (24206009)
Umesh Phulare (24206002)


Under the guidance of
**Prof. Priyanka Patil**



**Department of Computer Science & Engineering**
**(Artificial Intelligence & Machine Learning)**
**A. P. Shah Institute of Technology**
**G. B. Road, Kasarvadavali, Thane (W)-400615**
**University Of Mumbai**
**2024-2025**

# A. P. SHAH INSTITUTE OF TECHNOLOGY

# CERTIFICATE

This is to certify that the project entitled "**YouTube Video Transcribe Summarizer"** is a bona fide work of Prajwal Dhanawade (24206007), Gaurav Kshirsagar (24206009), Umesh Phulare (24206002), Ved Sawant (24206001) submitted to the University of Mumbai in partial fulfillment of the requirement for the award of **Bachelor of Engineering** in **Computer Science & Engineering (Artificial Intelligence & Machine Learning).**

_____                                      _____
Prof. Priyanka Patil                                        Dr. Jaya Gupta
Mini Project Guide                                        Head of Department

# A. P. SHAH INSTITUTE OF TECHNOLOGY

# Project Report Approval

This Mini project report entitled "**YouTube Video Transcribe Summarizer**" by **Prajwal Dhanawade, Gaurav Kshirsagar, Umesh Phulare and Ved Sawant** is approved for the degree of *Bachelor of Engineering* in *Computer Science &Engineering*, **(AIML)** *2024-25*.

External Examiner: _____

Internal Examiner: _____

Place: APSIT, Thane

Date:

## Declaration

We declare that this written submission represents my ideas in my own words and whereothers' ideas or words have been included, I have adequately cited and referenced the original sources. I also declare that I have adhered to all principles of academic honesty and integrity and have not misrepresented or fabricated or falsified any idea/data/fact/source in my submission. I understand that any violation of the above will be cause for disciplinary action by the Institute and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been taken when needed.

| Prajwal Dhanawade | Gaurav Kshirsagar | Umesh Phulare | Ved Sawant |
|---|---|---|---|
| (24206007) | (24206009) | (24206002) | (24206001) |

# ABSTRACT

**YouTube Video Transcription and Summarization Web App**

In the digital learning era, students frequently rely on YouTube for academic content, but manually transcribing video lectures is time-consuming and inefficient. This project introduces an AI-powered tool that automates the transcription and summarization process, converting YouTube video lectures into detailed study notes. By leveraging AI technologies such as NLP and summarization tools, the system enhances accessibility, reduces dependency on manual note-taking, and allows students to focus better on learning.

The project utilizes Python, Streamlit, the YouTube API, and libraries like Hugging Face to summarize the transcripts efficiently. This automation addresses challenges such as errors in manual transcription, difficulties faced by non-native speakers and learners with hearing impairments, and the time-consuming nature of traditional methods. By transforming passive video consumption into an active learning experience, this tool makes educational content more structured and accessible for a diverse range of learners.

**Keywords**: YouTube lecture transcription, AI-powered summarization, Study notes automation, Accessible learning.

# **INDEX**

# CHAPTER 1

# INTRODUCTION

# 1. INTRODUCTION

The rapid growth of online education has made YouTube a vital source of learning, but manually transcribing and summarizing educational videos is time-consuming and inefficient. The AI-powered YouTube Lecture Transcription and Summarization Tool is designed to assist students in converting video lectures into well-structured study notes, reducing manual effort and enhancing learning efficiency. By leveraging AI-driven transcription and summarization, this tool simplifies the process of note-taking, making it easier for students to focus on understanding concepts rather than spending hours writing notes.

The purpose of this project is to provide students with an accessible and automated solution for transforming YouTube lectures into structured study material. As digital learning continues to grow, this tool helps students save time by generating concise and accurate notes from video transcripts. Unlike traditional methods, which require frequent pausing and rewinding, this AI-powered solution ensures an effortless and error-free way of summarizing lectures.

This tool stands out by offering a range of features beyond basic transcription, including AI-based summarization, key point extraction, and structured note generation. It is particularly beneficial for learners who face difficulties in manual note-taking, such as non-native speakers and individuals with hearing impairments, ensuring that educational content remains accessible to all.

In addition to automated transcription, this tool provides valuable functionalities such as summarization, real-time processing, and an intuitive interface, making it suitable for students of all disciplines. By simply inputting a YouTube lecture URL, users can receive a structured summary, enabling them to review and retain knowledge effectively. Whether the goal is to revise efficiently, enhance accessibility, or simplify note-taking, this AI-powered tool ensures an optimized learning experience.

# CHAPTER 2
# LITERATURE SURVEY

# 2. LITERATURE SURVEY

## 2.1-HISTORY

The evolution of automated transcription and summarization tools has been closely linked to advancements in artificial intelligence (AI) and natural language processing (NLP). Early efforts in transcription relied on manual processes, requiring individuals to listen to audio and type out spoken content, which was both time-consuming and inefficient. The introduction of speech recognition technology in the late 20th century laid the groundwork for automated transcription systems.

In the early 2000s, basic speech-to-text software emerged, allowing users to convert spoken language into written text. However, these systems often struggled with accuracy, particularly in handling diverse accents, domain-specific terminology, and background noise. Early versions of tools like IBM ViaVoice and Dragon NaturallySpeaking pioneered commercial speech recognition, but their widespread adoption was limited due to high error rates and the need for extensive training.

With the rise of cloud computing and machine learning in the 2010s, AI-powered transcription tools saw significant improvements. Services like Google Speech-to-Text, IBM Watson Speech Services, and Microsoft Azure Speech began offering real-time, high-accuracy transcription capabilities. These tools leveraged deep learning models trained on vast datasets, allowing them to handle multiple languages and accents more effectively.

The introduction of YouTube's automatic captions in 2009 marked a major milestone in making online video content more accessible. However, early versions often had errors due to the limitations of the underlying speech recognition models. Over time, YouTube refined its automatic captioning system by incorporating deep learning-based improvements, increasing accuracy and usability.

In recent years, the integration of AI with NLP models such as OpenAI's Whisper, Gemini Pro, and BERT has further enhanced the capabilities of transcription and summarization systems. These advancements have enabled the generation of structured, summarized content from long-form lectures, making learning more efficient and accessible. Today, AI-driven transcription tools are widely used across education, media, and business sectors, transforming how information is processed and consumed.

## 2.2-LITERATURE REVIEW

**YouTube Video Summarizer using NLP: A Review [IJPE (Dec 2023)]**

This review paper delves into the emerging realm of YouTube video summarization utilizing Natural Language Processing (NLP) techniques, a critical area of research with increasing prominence in our multimedia-rich digital age. The paper commences with a broad overview of the field, elaborating on the need for automated video summarization tools to navigate and condense the massive, ever-growing sea of YouTube content. Further, we systematically scrutinize the role and implementation of NLP methods in extracting meaningful textual data from videos, focusing on video transcripts, closed captions, user comments, and associated metadata. Subsequent sections dissect seminal and recent works, studying various NLP techniques such as text summarization, sentiment analysis, topic modeling, and deep learning architectures employed in this context. The paper also focuses on the various metrics used for evaluation and shows datasets generally used to assess the performance of these summarization systems. Finally, we identify current challenges and potential future directions for research in the area, acknowledging the evolving landscape of online video platforms and AI technologies. This review aims to provide researchers and practitioners with an encompassing perspective on the pivotal role of NLP in enabling more efficient, accurate, and intuitive navigation of YouTube content ultimately shaping our digital consumption experiences.

**YouTube Transcript Summarizer To Summarize the content of YouTube.[ IRJET (April 2022)]**

Watching long YouTube videos is very time-consuming and boring. Nowadays YouTube is an essential aspect of providing news and information. It is also considered a second teacher to the students, educational videos are the most viewed videos on YouTube today. In this project, we have tried to provide a quick, precise, and informative summary of a video. Many techniques are already discovered but they only provide test summarization. We have tried to get the summary of a video basically a YouTube video. For this project, we have used a hugging face transformer to summarize the content of a YouTube video along with that we have used python API to get the subtitle of a given video. After that our model will perform text summarization on it and display the summary to the user so that people can save their precious time reading the summary.

**Video Summarization Using Using Fully Convolutional Sequence Networks. [ ECCV 2018]**

This paper addresses the problem of video summarization. Given an input video, the goal is to select a subset of the frames to create a summary video that optimally captures the important information of the input video. With the large amount of videos available online, video

summarization provides a useful tool that assists video search, retrieval, browsing, etc. In this paper, we formulate video summarization as a sequence labelling problem. Unlike existing approaches that use re current models, we propose fully convolutional sequence models to solve video summarization. We firstly establish a novel connection between semantic segmentation and video summarization, and then adapt popular semantic segmentation networks for video summarization. Extensive experiments and analysis on two benchmark datasets demonstrate the effectiveness of our models.

**Implications of Using AI in Translation Systems [AJRESS 2021]**

This review paper provides an overview of the use of artificial intelligence (AI) in Translation Studies (TS), covering statistical machine translation, rule-based machine translation, neural machine translation, and hybrid machine translation. It explores the advantages and limitations of each model, as well as their applications in translation. Additionally, it discusses various techniques for evaluating the effectiveness of AI models in translation, along with their advantages and limitations, such as handling figurative language (e.g., idioms, metaphors) and cultural nuances. The review also delves into research directions for improving AI-based translation, elaborates on the ethical and social implications of AI in translation, and discusses the representation of AI in other disciplines such as literature and arts. Finally, the impact of AI as well as the opportunities and challenges that it could create for translators, such as professional challenges, data privacy, bias, and fairness matters were briefly discussed. By summarizing the main findings, and lessons learnt in AI-based translation, some recommendations regarding the current and future direction of using AI in translation were formulated.

**Turning Whisper into Real-Time Transcription System [NICT 2023]**

Whisper is one of the recent state-of-the-art multilingual speech recognition and translation models, however, it is not designed for real time transcription. In this paper, we build on top of Whisper and create Whisper-Streaming, an implementation of real-time speech transcription and translation of Whisper-like models. Whisper-Streaming uses local agreement policy with self-adaptive latency to enable streaming transcription. We show that Whisper Streaming achieves high quality and 3.3 seconds latency on unsegmented long-form speech transcription test set, and we demonstrate its robustness and practical usability as a component in live transcription service at a multilingual.

**Summary of literature review in tabular form:**

| Paper Title | Authors | Journal/Conference | Techniques Used | Key Findings/Results |
|---|---|---|---|---|
| YouTube Video Summarizer using NLP: A Review | Yogendra Singh, Rishu Kumar, Soumya Kabdal, Prashant Upadhyay | IJPE (Dec 2023) | NLP methods like BERT, TF-IDF, and LSA; Keyframe-based techniques; Sentiment analysis; Topic modeling | Focuses on the role of NLP in extracting meaningful data from video transcripts, captions, and metadata. Explores various evaluation metrics and future research directions. |
| YouTube Transcript Summarizer To Summarize the content of YouTube | Sourav Biswas, Atul Kumar Patel | IRJET (April 2022) | Hugging Face Transformers; DistilBERT; Extractive and Abstractive Summarization | Uses NLP models like DistilBERT for concise YouTube video summarization. Demonstrates improved efficiency in summarizing educational content and professional meetings. |
| Video Summarization Using Fully Convolutional Sequence Networks | Mrigank Rochan, Linwei Ye, and Yang Wang | ECCV (2018) | Deep learning, Graph-theoretic methods, LSTM, BERT | Introduces a semantic graph-based method for improved sentence correlation and summarization. Also explores personalized summarization with user-preferred templates. |
| Implications of Using AI in Translation Systems | Mansour Amini, Latha Ravindran, Kam-Fong Lee | Asian Journal of Research in Education and Social Sciences (AJRESS - 2021) | Transformer models for language translation; NLP for text comprehension | Explores AI's impact on translation accuracy, performance, and linguistic challenges in multilingual translation tasks. |
| Turning Whisper into Real-Time Transcription System | Dominik Machácek, Raj Dabre, Ondrej Bojar | National Institute of Information and Communications Technology (NICT), 2023 | AI-driven summarization; NLP text processing techniques | Explores efficient summarization for large volumes of video content with minimal information loss. Offers insights into enhancing summarization quality. |

# CHAPTER 3

# PROBLEM STATEMENT

# 3. PROBLEM STATEMENT

The growing reliance on YouTube as a primary source of educational content presents significant challenges for students who need to convert video lectures into structured study material. Manually transcribing and summarizing lectures is a time-consuming and inefficient process that requires constant pausing, rewinding, and note-taking, often resulting in incomplete or inaccurate information. This approach is particularly problematic for students managing multiple subjects, as it significantly reduces study efficiency.

Additionally, accessibility remains a major concern, as learners with hearing impairments or non-native speakers struggle to comprehend video content effectively. Existing solutions, such as YouTube's automatic captions, often lack accuracy and do not provide structured summaries, making it difficult for students to extract key concepts from lengthy lectures. The absence of an automated and structured approach to transcription and summarization limits students' ability to review and retain information efficiently.

Therefore, there is a critical need for an AI-powered tool that automates the transcription and summarization of YouTube lectures. Such a system would enhance learning by providing students with accurate, structured study notes, eliminating the need for manual effort, and improving accessibility. By leveraging advanced AI and NLP technologies, this tool would enable users to efficiently process educational content, transforming passive video consumption into an active and effective learning experience.

# CHAPTER 4

# EXPERIMENTAL SETUP

# 4. EXPERIMENTAL SETUP

This project is a web-based application designed to automate the process of YouTube video transcription and summarization using advanced AI models and text processing techniques. It leverages multiple frameworks and tools to ensure accuracy, efficiency, and an interactive user experience.

## 4.1 Hardware Setup

- **Processor**: Minimum Quad-core CPU with support for parallel processing to handle Whisper AI transcription efficiently.
- **GPU** (Optional but Recommended): NVIDIA CUDA-enabled GPU for faster Whisper AI inference.
- **RAM**: Minimum 8GB RAM (16GB recommended for optimal performance with large transcripts).
- **Storage**: At least 10GB free space for audio processing and temporary files.

## 4.2 Software Setup

- **Operating System:** Compatible with **Windows**, **Linux**, or **MacOS**.
- **Python Environment:**
    - Python **3.10** or higher.
    - Required libraries:
        a) **yt-dlp** – For downloading video/audio.
        b) **whisper** – For audio transcription.
        c) **transformers** – For text summarization.
        d) **google-generativeai** – For text refinement.
        e) **streamlit** – For frontend UI.
        f) **FPDF** – For PDF generation.
- **Visual Studio Code (IDE):** Used for writing, debugging, and testing code.
- **Live Server Extension:** For real-time code updates during frontend development.
- **Google Chrome:** Used for testing, debugging, and inspecting web application behaviour.
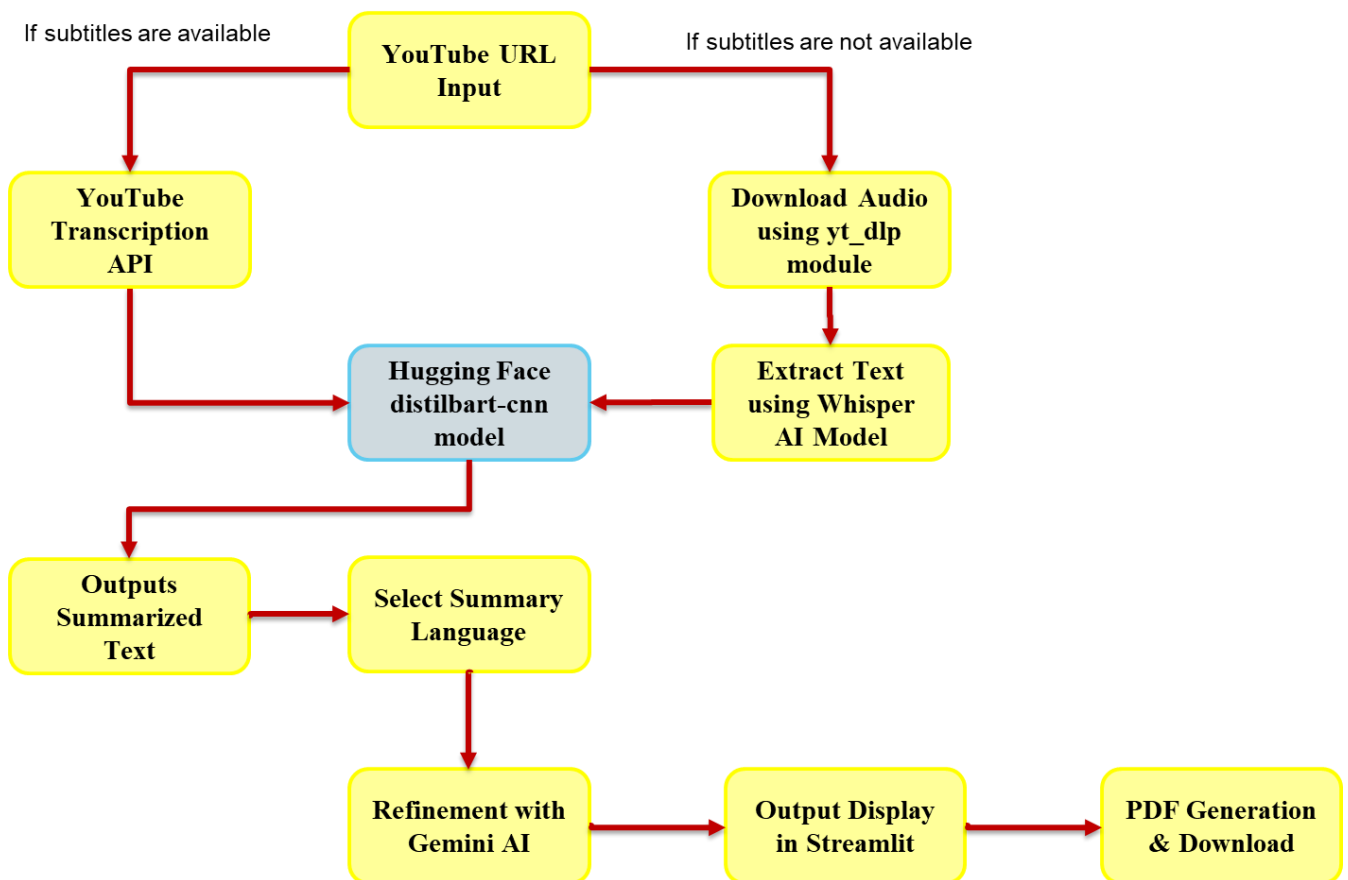
# CHAPTER 5
# PROPOSED SYSTEM
# &
# IMPLEMENTATION

# 5. PROPOSED SYSTEM & IMPLEMENTATION

This block diagram outlines the workflow of the YouTube Video Summarizer project, showing the key components and their interactions.

## 5.1 DIAGRAMS OF PROPOSED SYSTEM

If subtitles are available

**YouTube URL Input**

If subtitles are not available

**YouTube Transcription API**

**Download Audio using yt_dlp module**

**Hugging Face distilbart-cnn model**

**Extract Text using Whisper AI Model**

**Outputs Summarized Text**

**Select Summary Language**

**Refinement with Gemini AI**

**Output Display in Streamlit**

**PDF Generation & Download**

## 5.2 DESCRIPTION OF DIAGRAMS.

- **YouTube URL Input**
  - The system begins with the user providing a YouTube video URL as input.
  - This initiates two parallel processes — one for extracting the transcript via the YouTube Transcription API and another for audio download and transcription using Whisper AI.

- **YouTube Transcription API**
  - If the video has subtitles available, the system fetches the transcript directly using the YouTube Transcription API.
  - This ensures faster and more accurate text retrieval for supported videos.

- **Download Audio using yt_dlp Module**
  - If subtitles are not available, the system downloads the video's audio using the yt_dlp module.

- **Extract Text using Whisper AI Model**
  - The downloaded audio is processed using the Whisper AI model, which performs speech-to-text conversion to extract the spoken content.
  - This step ensures transcription even for videos without subtitles.

- **Hugging Face distilbart-cnn Model (Summarization)**
  - The extracted text (from either the YouTube Transcription API or Whisper AI) is passed to the Hugging Face distilbart-cnn model for generating a concise and meaningful summary.
  - This model is optimized for producing clear, structured summaries.

- **Outputs Summarized Text**
  - The summarized text is presented as a direct output.
  - Users can optionally refine this summary further.

- **Select Summary Language**
  - Users can select their preferred language for refined output (e.g., English, Hindi, Marathi).

- **Refinement with Gemini AI**
  - The summarized content is refined using Gemini AI, improving clarity, coherence, and structure.
  - This step enhances the quality of the final content.

- **Output Display in Streamlit**
  - The refined summary is displayed interactively using Streamlit, offering a user-friendly interface for improved readability.

- **PDF Generation & Download**
  - The final summary is formatted into a well-structured PDF document for users to

download.
- o This ensures easy sharing and offline reference for users.


Overall, the diagrams illustrate the hierarchical structure of user access, the basic data flow within the system & the overall structure of the database.

## 5.3 IMPLEMENTATION

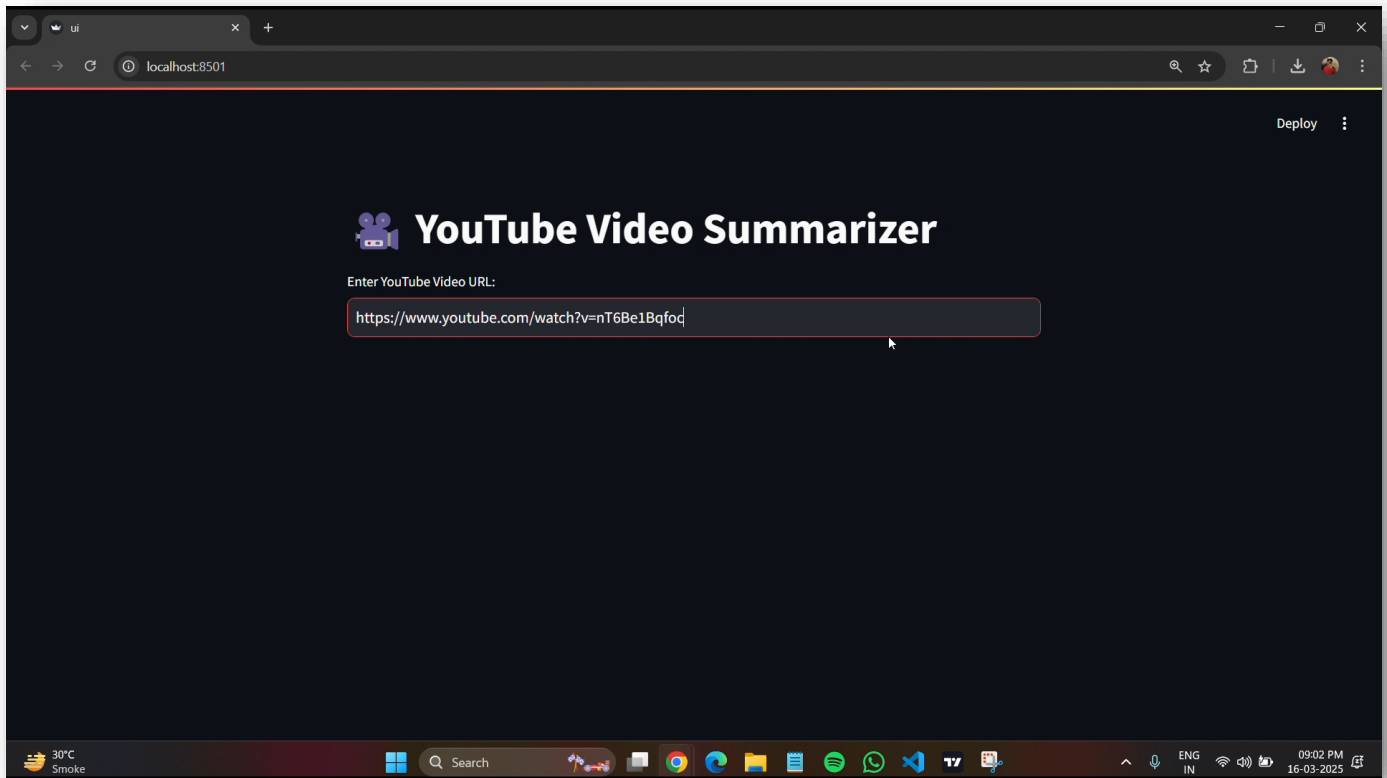Implementation of proposed system is included here as screenshots:
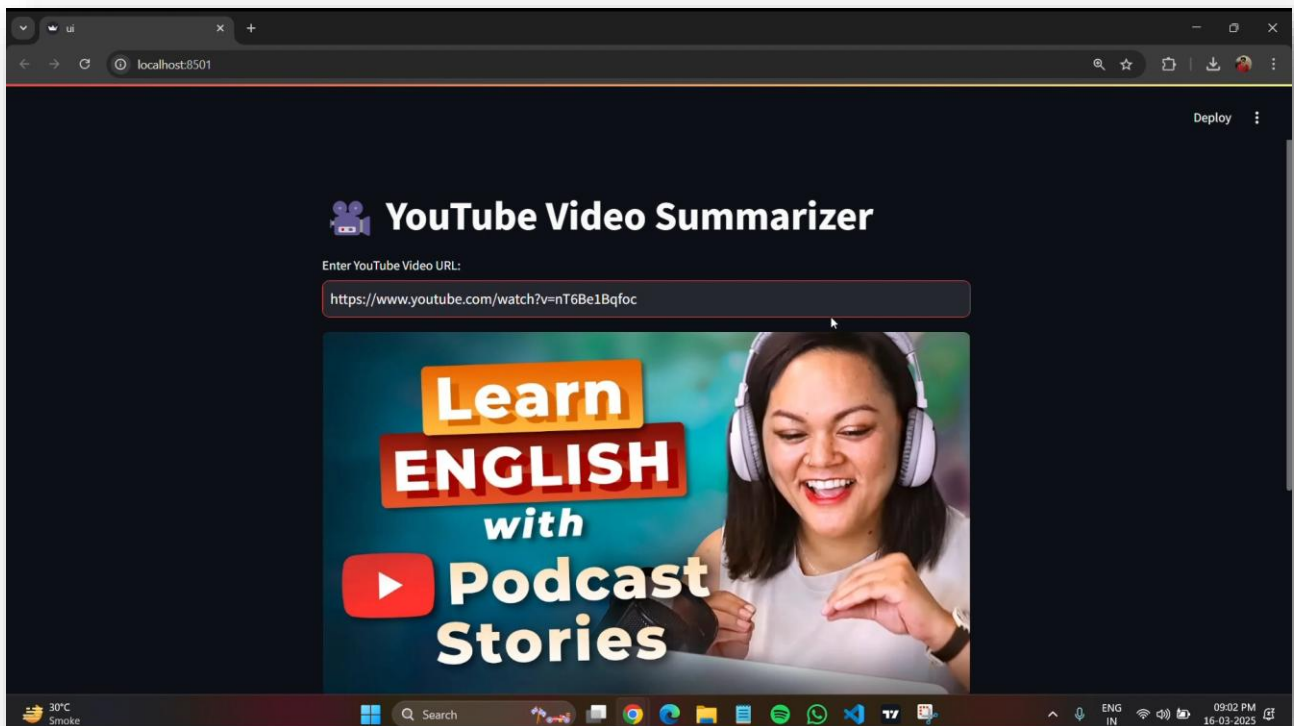


*Figure 1: Main page*
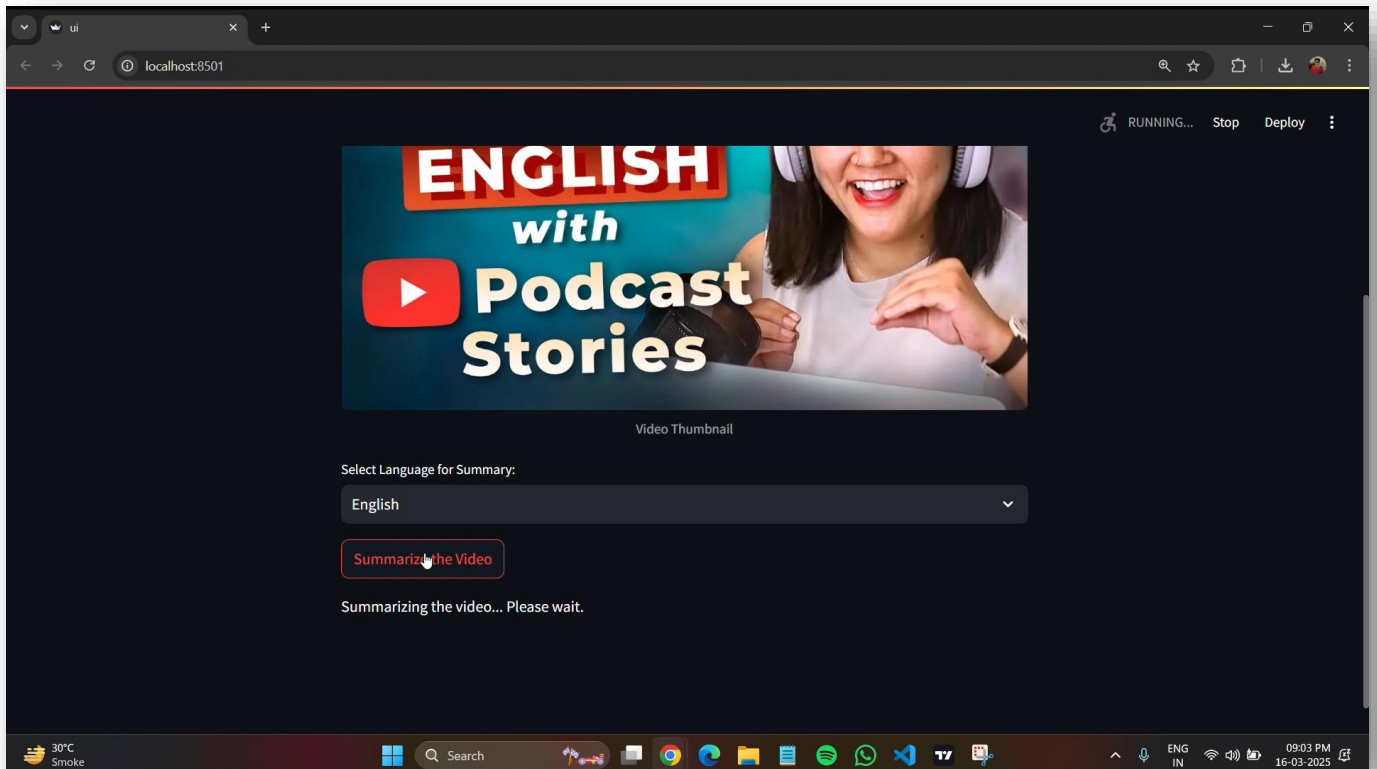


*Figure 2: YouTube Video Thumbnail*

19

*Figure 3: Preferred Summarization Language Selection*
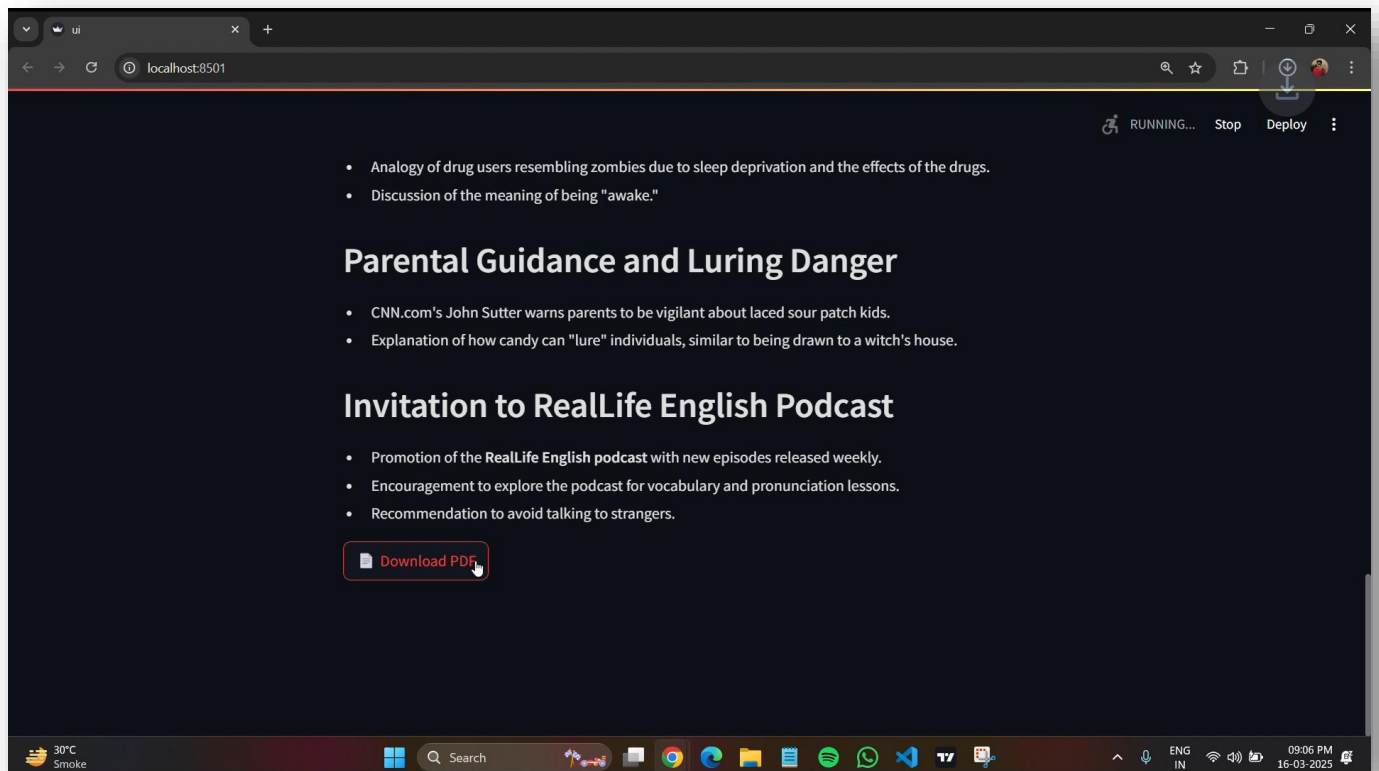


*Figure 4: Video's Summarized Content*
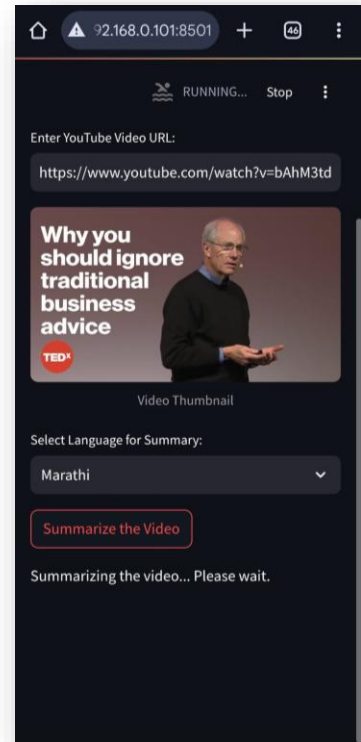
*Figure 5: Downloading Summarized Content PDF*
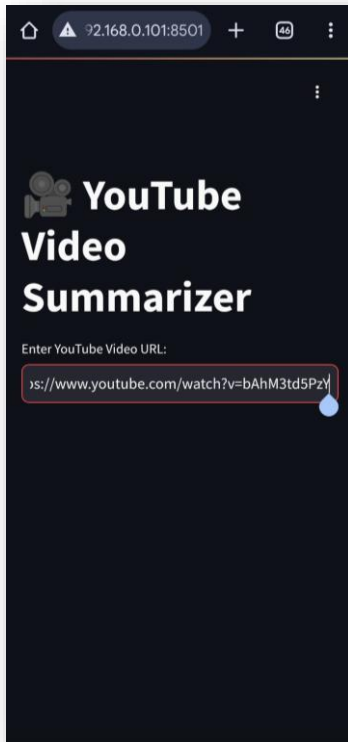


*Figure 6: summary.pdf*
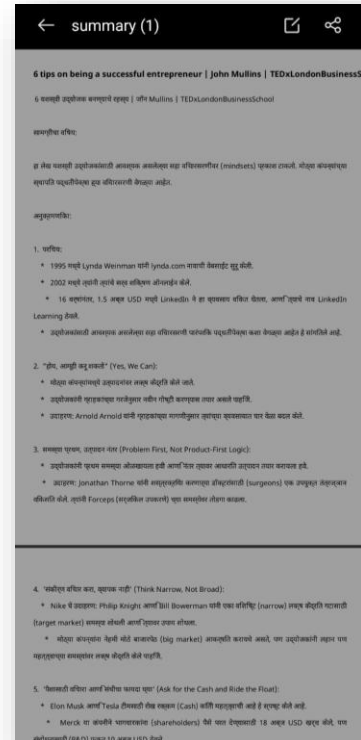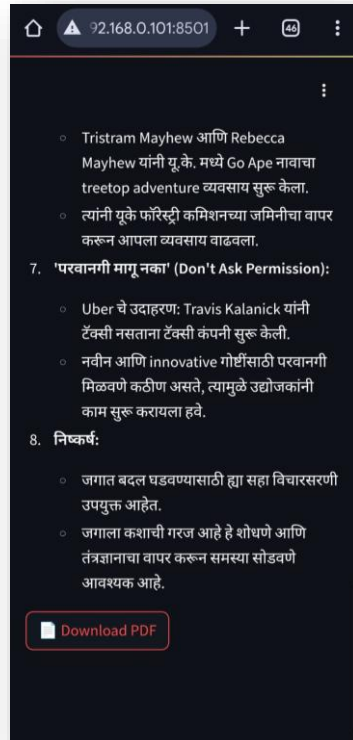
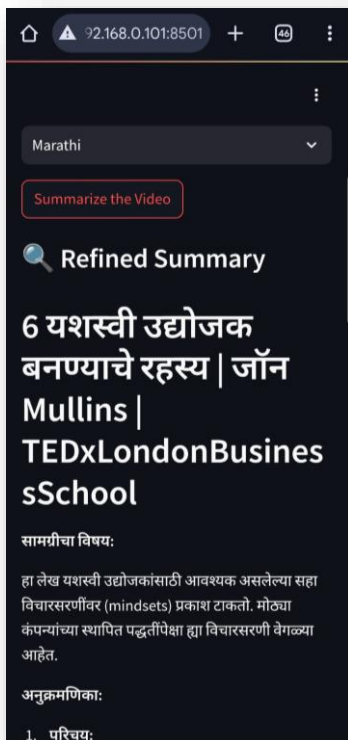*Figure 7: Running Same Website in Mobile*



*Figure 8: Summarized Content in downloadable PDF file*

## 5.4 Applications

**AI-Powered YouTube Lecture Transcription and Summarization Tool**

1. **Lecture Transcription:**
   The tool provides detailed and accurate transcriptions of YouTube lectures, enabling users to access a written version of the content. This feature benefits students, researchers, and professionals by allowing them to read and reference key points effectively.

2. **Summarized Notes Generation:**
   The tool generates concise and comprehensive summaries of lengthy lecture videos. These summaries capture the essence of the lecture, saving users valuable time while ensuring they understand the key concepts.

3. **Keyword Extraction for Quick Reference:**
   By extracting important keywords from the lecture, the tool allows users to navigate the content efficiently. This feature is especially useful for searching specific topics within long lectures.

4. **Interactive Learning Interface:**
   The tool offers a user-friendly interface that allows users to interact with transcripts and summaries. Features like keyword search and note-taking enhance the learning experience, making the tool suitable for academic and professional use.

5. **Multilingual Support:**
   The transcription tool supports multiple languages, enabling users to transcribe and summarize lectures in their preferred language. This inclusivity broadens the accessibility of educational content worldwide.

6. **Customizable Export Options:**
   Users can save and export transcriptions and summaries in various formats, such as PDF, Word, or plain text. This flexibility ensures that the content can be easily integrated into personal notes or shared with peers.

7. **Enhanced Accessibility for Disabled Users:**
   The transcription feature makes video lectures accessible to individuals with hearing impairments. Summarized notes also support learners with cognitive challenges, ensuring inclusivity in education.

8. **Efficient Time Management:**
   By summarizing and highlighting critical points, the tool helps users save time by focusing only on relevant content. This is especially beneficial for students preparing for exams or professionals conducting research.

9. **Research and Collaboration:**
   The tool aids researchers by providing structured content that can be cited or discussed. Teams can use the tool to collaboratively analyze and interpret lecture material, fostering productive group work.

10. **Cost-Effective Educational Resource:**
    By offering a comprehensive solution for lecture transcription and summarization, the tool eliminates the need for manual note-taking or hiring transcription services, making it a cost-effective choice for students and educators.

# CHAPTER 6
# CONCLUSION

# 6. CONCLUSION

## 6.1 Conclusion

The AI-Powered YouTube Transcript Summarization System has successfully met its objectives by providing users with an efficient tool to transcribe, summarize, and navigate video content seamlessly. The system offers highly accurate transcription, concise summaries, and intuitive keyword-based navigation, ensuring a user-friendly experience. The tool empowers users by enhancing learning efficiency and accessibility, particularly for students, educators, and professionals seeking to maximize their productivity with video resources. Overall, the system effectively addresses the challenges it was designed to solve, as outlined in the project specifications.

## 6.2. Limitations of System

A primary limitation of the system is its reliance on robust internet connectivity for real-time transcription and summarization processes, which may hinder users in areas with poor network coverage. Additionally, the system's performance could be affected when processing videos with low audio quality or heavily accented speech.

## 6.3 Future Scope of the Project

The project holds significant potential for future enhancements aimed at improving its accessibility, efficiency, and overall user experience. One of the key directions is the implementation of enhanced multilingual support, allowing the system to cater to a wider audience by accommodating various languages and dialects. Additionally, the integration of visual elements such as charts, graphs, and diagrams will enable more comprehensive content analysis. To cater to different user needs, adaptive summarization capabilities can be introduced, allowing users to choose between brief, moderate, or detailed summaries. Further, the system can be improved by incorporating video context awareness, enabling it to detect scene changes, differentiate between speakers, and understand transitions within content more accurately. A highly impactful feature would be real-time summarization, allowing users to obtain summaries of live streams and webinars instantly. From a performance standpoint, optimizing Whisper AI for better GPU utilization will significantly reduce processing time and enhance efficiency. Lastly, offering a developer-friendly API would facilitate seamless integration with other platforms and applications, expanding the utility and reach of the system.

# REFERENCES

**Research papers:**

[1] Yogendra Singh, Rishu Kumar, Soumya Kabdal, Prashant Upadhyay, "YouTube Video Summarizer using NLP: A Review", International Journal of Engineering Research and Technology (IJERT), vol. 19, no. 12, pp. 817-823, December 2023.

[2] Sourav Biswas, Atul Kumar Patel, "YouTube Transcript Summarizer To Summarize the content of YouTube", International Research Journal of Engineering and Technology (IRJET), Volume: 09 Issue: 04, April 2022.

[3] Mrigank Rochan, Linwei Ye, and Yang Wang, "Video Summarization Using ECCV 2018 Techniques", European Conference on Computer Vision (ECCV), published in 2018.

[4] Eka Wahyu Aditya Shahrinaz Ismail and Noormadinah Allias, "Implications of Using AI in Translation Systems", Asian Journal of Research in Education and Social Sciences e-ISSN: 2682-8502 | Vol. 6, No. 1,740-754, 30 April 2024.

[5] Dominik Machácek , Raj Dabre, Ondrej Bojar, "Turning Whisper into Real-Time Transcription System", National Institute of Information and Communications Technology, 2023.

**URL**

https://www.ijpe-online.com/EN/10.23940/ijpe.23.12.p6.817823

https://www.irjet.net/archives/V9/i4/IRJET-V9I4520.pdf

https://openaccess.thecvf.com/content_ECCV_2018/papers/Mrigank_Rochan_Video_Summarization_Using_ECCV_2018_paper.pdf

https://myjms.mohe.gov.my/index.php/ajress/article/view/24617

https://arxiv.org/abs/2307.14743