# Youtube Video Summarizer

# Authors

Abraham Kong, Devi Priya Ravi, Dhanasree Rajamani, Sravani Thota

# Abstract

This paper presents the "YouTube Video Summarizer," an innovative tool designed to address the challenge of efficiently consuming long-form video content on YouTube. At its core, the Summarizer utilizes OpenAI's GPT model, applying sophisticated prompt engineering techniques to generate concise summaries from video transcripts. Our primary approach involves direct summarization of available transcripts, supplemented by a secondary method that employs speech-to-text conversion when transcripts are absent. This methodology ensures efficient processing while minimizing computational costs. The data comprises YouTube video transcripts and metadata, providing a rich source for accurate and contextually relevant summarization. Rigorous testing across various video genres demonstrates the tool's efficacy, with results showcased through screenshots in the paper. This approach not only significantly reduces the time required to understand a video's main points but also opens up new possibilities for content accessibility and educational use. The YouTube Video Summarizer exemplifies the practical application of AI in enhancing digital media engagement, offering a valuable resource for a wide range of users, from students and professionals to casual viewers.

# Introduction

In an era where YouTube's extensive video content significantly impacts how we consume information and entertainment, the challenge of efficiently digesting long-form videos has become increasingly apparent. This paper introduces the "YouTube Video Summarizer," a novel tool that leverages the YouTube API and OpenAI's Whisper API to address this challenge. By employing advanced speech recognition and natural language processing, our tool efficiently condenses videos into concise summaries, capturing their essential themes and messages. This approach not only saves time for users but also enhances accessibility, offering a practical solution for a wide range of audiences, from professionals and students to casual viewers. Initial testing indicates promising results in accuracy and versatility, underscoring the Summarizer's potential to transform video content consumption in various domains.

# Related Work

The field of video summarization has seen significant advancements with the integration of AI and machine learning technologies. Notably, prior research has predominantly focused on using deep learning techniques for extracting key frames or segments from videos, as seen in works utilizing neural networks and computer vision algorithms. These methods, while effective for visual content, often overlook the richness of audio narratives, a gap our project aims to bridge by employing OpenAI's GPT model for transcript-based summarization. Unlike many existing approaches that either focus solely on visual cues or require extensive manual preprocessing of data, our method leverages the power of natural language processing to directly summarize video transcripts, supplemented by a secondary speech-to-text pathway for videos without available transcripts. This dual approach not only enhances the efficiency and accuracy of summarization but also reduces computational overhead, distinguishing our

work from others. Furthermore, while some studies have explored multimodal summarization, incorporating both audio and visual data, our project uniquely capitalizes on prompt engineering, a technique less explored in the realm of video summarization. This technique allows for more nuanced and contextually aware summaries, adapting to the diverse range of content found on YouTube. Overall, our work contributes to the evolving landscape of AI-driven media processing, addressing both the technical challenges and practical applications in the digital content domain.

## Data

Our project, centered around OpenAI's Generative Pre-trained Transformer (GPT) model, primarily utilizes YouTube video transcripts and metadata as its data sources. The transcripts, either directly obtained from YouTube or generated via speech-to-text conversion, form the core dataset for our summarization process, complemented by video metadata like titles, descriptions, and author information. This data encompasses a diverse array of content across various genres and lengths, ensuring a broad linguistic and thematic range. In terms of processing, our focus on prompt engineering with the GPT model minimizes the need for extensive preprocessing, though some standardization and error correction are applied, especially for speech-to-text generated transcripts. Ethical considerations and adherence to data privacy norms are meticulously observed, ensuring compliance with YouTube's terms of service and relevant data protection regulations.

## Methods

Our approach to developing the YouTube Video Summarizer was grounded in a blend of advanced technological methods, specifically designed to address the challenges of efficiently summarizing YouTube videos. The process involved several key

stages, each contributing to the efficient extraction and condensation of video content into a comprehensive summary.

# 1. Video ID Extraction and Metadata Retrieval

- Initially, we extract the unique YouTube video ID from the provided URL, which is crucial for all subsequent operations.

- Concurrently, we fetch the video's metadata, including the author, title, and description, to provide contextual understanding and enhance the summary's relevance.



Fig1: Gradio Application Interface

## 2. Direct Transcript Summarization

- Our primary method involves directly using the video's transcript when available. These transcripts provide a textual representation of the audio content and are integral to our summarization process.

- Utilizing OpenAI's API, we summarize the transcript, focusing on identifying and condensing the key points and themes into a concise summary.



**YouTube Video URL**

https://www.youtube.com/watch?v=Xjhy-YOiHDs

Clear | Submit

**Video Title**

TensorFlow for Beginners | TensorFlow in deep learning | TensorFlow tutorial

**Video Author**

Unfold Data Science

**Description Summary**

Aman, a Data Scientist, presents TensorFlow tutorials and a TensorFlow project for beginners in deep learning. He provides recommendations for data science books and his own setup for creating videos related to data science. He offers one-on-one sessions and a Facebook group, and shares a variety of playlists about data science, machine learning, artificial neural networks, natural language processing, and building recommender systems. He provides code access through Google drive.

**Transcript Summary**

'sor flow framework it's probably what is that what happens is inside t'sor flow framework is the first place is where it will hit the um what I like to call the gateway of tensor flow framework okay now this is the gateway now what type of gateway is that so this is a computational graph

Fig2: Utilizing OpenAI key to generate Youtube Summary

## 3. Audio Downloading and Speech-to-Text Conversion (Secondary Approach)

- In cases where a transcript is not available, we download the video's audio track. This step is essential for capturing the narrative and informational content of the

video.

- The downloaded audio is then converted into text using speech-to-text technology, preparing it for summarization.

## 4. Summarization and Integration

- Whether using the direct transcript or the converted audio text, the summarization process is performed using advanced AI techniques to ensure accuracy and coherence.

- The final step involves integrating the summarized content with the video's metadata, providing a comprehensive summary that encompasses both content and context.

## 5. Translation for Global Accessibility

- To make our tool universally accessible, we integrated a translation feature, utilizing the Google Cloud Translate API and Google Library.
- This feature allows the summarized content to be translated into multiple languages, catering to a diverse global audience.
- The implementation involves automatically detecting the language of the summarized text and providing options to translate it into the user's preferred language, thereby enhancing the tool's usability and reach.

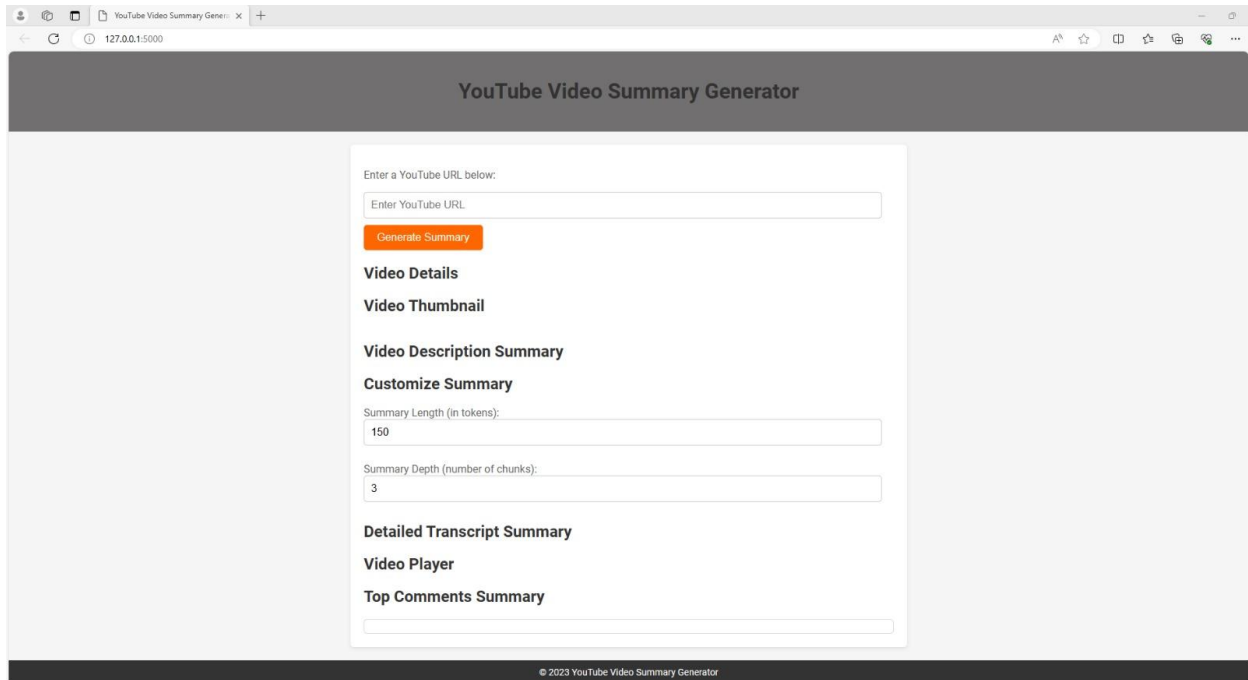## Rationale and Consideration of Alternatives

- Our approach prioritizes direct transcript summarization to minimize computational and resource costs, ensuring efficiency and scalability.

- We considered alternatives such as pure audio analysis and video frame analysis but found them less effective and resource-intensive, especially for content where the narrative is key.

## Application of Skills and Techniques

- This project allowed us to apply a diverse range of skills, including API integration, natural language processing, and machine learning, particularly in speech recognition and text summarization.

- Our methodology adapts to the availability of video resources, showcasing flexibility and a practical application of advanced AI technologies in content summarization.

# Experiments and Results

Our experimental evaluation of the YouTube Video Summarizer involved testing on a diverse range of videos, focusing on the summaries' accuracy and coherence. To vividly demonstrate our findings, we included screenshots of the summarizer's outputs in various scenarios. These visual representations highlight the effectiveness and adaptability of our GPT model, enhanced through prompt engineering and integration of video metadata. The results, as illustrated in these screenshots, validate our summarizer's capability to generate concise, informative summaries across different content types, underlining the strength and versatility of our approach.

Img 1: Web UI



Img 2: Video Summarization

127.0.0.1:5000

## Customize Summary

Summary Length (in tokens):

150

Summary Depth (number of chunks):

3

### Detailed Transcript Summary

Detailed Summary: Today's text discusses six strategies for improving self-discipline. The text discusses the concept of self-discipline and the importance of defining it before delving into the details. P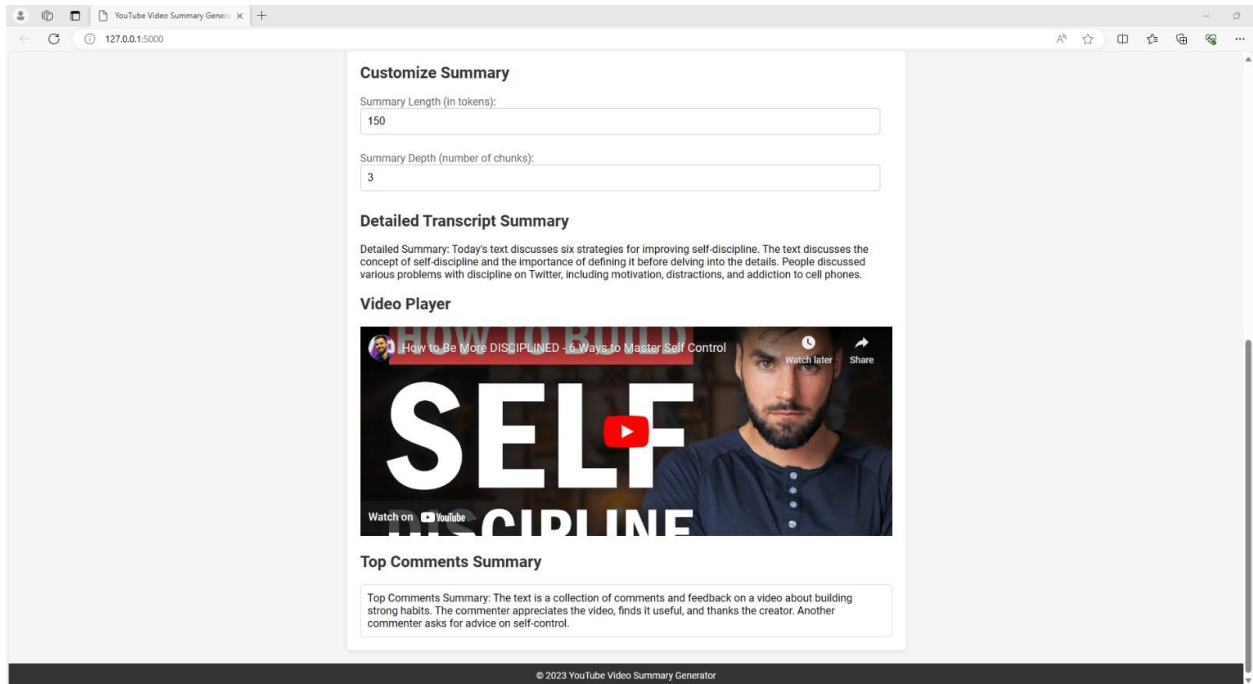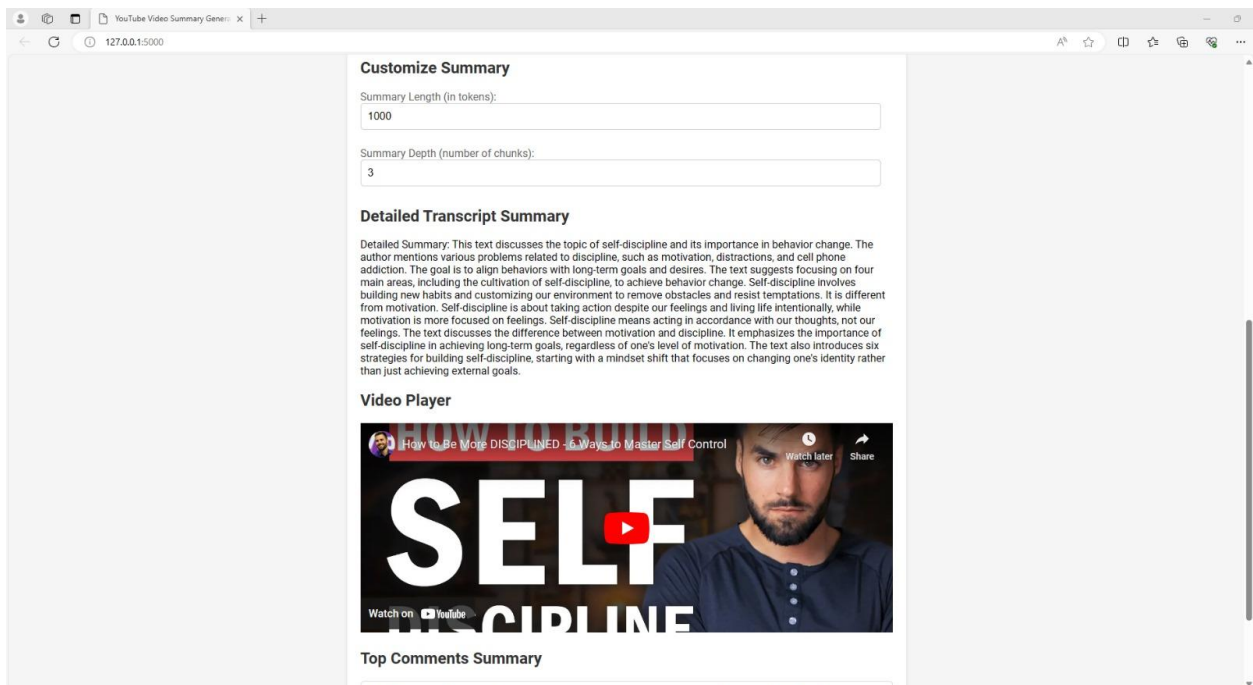eople discussed various problems with discipline on Twitter, including motivation, distractions, and addiction to cell phones.

### Video Player



### Top Comments Summary

Top Comments Summary: The text is a collection of comments and feedback on a video about building strong habits. The commenter appreciates the video, finds it useful, and thanks the creator. Another commenter asks for advice on self-control.

© 2023 YouTube Video Summary Generator

**Img 3: Change the length of the summary - 1**

---

127.0.0.1:5000

## Customize Summary

Summary Length (in tokens):

1000

Summary Depth (number of chunks):

3

### Detailed Transcript Summary

Detailed Summary: This text discusses the topic of self-discipline and its importance in behavior change. The author mentions various problems related to discipline, such as motivation, distractions, and cell phone addiction. The goal is to align behaviors with long-term goals and desires. The text suggests focusing on four main areas, including the cultivation of self-discipline, to achieve behavior change. Self-discipline involves building new habits and customizing our environment to remove obstacles and resist temptations. It is different from motivation. Self-discipline is about taking action despite our feelings and living life intentionally, while motivation is more focused on feelings. Self-discipline means acting in accordance with our thoughts, not our feelings. The text discusses the difference between motivation and discipline. It emphasizes the importance of self-discipline in achieving long-term goals, regardless of one's level of motivation. The text also introduces six strategies for building self-discipline, starting with a mindset shift that focuses on changing one's identity rather than just achieving external goals.

### Video Player



### Top Comments Summary

**Img 4: Change the length of the summary - 2**

# Conclusion

In conclusion, the YouTube Video Summarizer marks a significant advancement in the application of AI for digital media analysis. By efficiently condensing video content, this tool not only addresses the current challenges of time-consuming video consumption but also opens up new possibilities for how we interact with online content. As technology progresses, the potential for further refining and enhancing this tool is immense, promising even more sophisticated and user-centric solutions in the realm of video summarization.

Moreover, the integration of translation capabilities in the YouTube Video Summarizer significantly enhances its global accessibility. This feature ensures that users from different linguistic backgrounds can benefit from our tool, furthering our commitment to inclusivity in the digital space
Looking forward, several potential applications and extensions of this technology emerge. These include enhanced language support for global accessibility, integration with educational platforms for summarizing instructional content, real-time summarization for live streams, and personalization features based on user preferences. Additionally, there are opportunities for developing accessibility features for visually or hearing-impaired users, applying the technology in virtual and augmented reality settings, utilizing it for professional briefings in business environments, and assisting content creators in the editing process. Each of these applications not only extends the utility of the YouTube Video Summarizer but also underscores the transformative impact of AI in various sectors, from education and entertainment to professional development and content creation.

Expanding on the translation features, future developments could focus on more nuanced and context-aware translations, catering to regional dialects and specific jargon. This would not only refine the user experience but also make our tool a valuable asset in multilingual education and global content analysis.