

# **ContextualVid: Video Translation, Summarization, and Content Generation**

## **FIELD OF INNOVATION:**

The area of innovation is in the smooth integration of cutting-edge technologies like machine learning, speech recognition, and natural language processing to produce a versatile platform that gets beyond conventional language boundaries for the access and comprehension of video information. Our technology ensures accessibility for a worldwide audience by utilizing speech-to-text and text-to-speech algorithms to provide real-time transcription and translation capabilities. Moreover, the use of advanced summarizing methods enables the extraction of significant insights from films, facilitating users' rapid assimilation of crucial information. The innovation is further enhanced by adding a contextual generator function that uses cutting-edge natural language processing to deliver pertinent context gleaned from the video's content and promotes comprehension and deeper engagement. Our platform transforms the way people interact with video content by fusing breakthrough technologies in a way that makes it easily accessible, understandable, and captivating in a variety of languages and cultural contexts.

## **OBJECTIVE:**

The goal of the idea is to provide a flexible and intuitive platform that improves accessibility, breaks down barriers based on language, and encourages more in-depth interaction with video material. The invention aims to give users a complete toolkit for navigating and understanding video content in their preferred language by integrating advanced language processing technologies, such as transcription, translation, summarization, multilingual voiceover generation, text generation, and contextual understanding. Furthermore, the platform places a high priority on security, privacy, scalability, and flexibility to satisfy users' changing demands and guarantee a smooth and enjoyable viewing experience in a variety of situations and audiences.

## **SCOPE OF THE INVENTION:**

The invention's scope includes a full range of accessibility and language processing functionalities designed for video material, such as summarization, translation, transcription, multilingual voiceover creation, and contextual understanding. This innovation seeks to improve comprehension, promote inclusion for audiences throughout the world, and overcome language barriers through the integration of state-of-the-art technologies and user-friendly design. Strong security measures are also put in place to protect user data and guarantee

privacy, and scalability and flexibility are prioritized to meet future improvements and changing user requirements regarding the accessibility and comprehension of video material.

### **DESCRIPTION:**

Our project intends to change the accessibility and comprehension of video material by providing a holistically integrated platform that streamlines the transcription, translation, summarization, text generation, and multilingual voiceover production processes. Our website offers customers a user-friendly interface that allows them to easily post videos or link to current material using cutting-edge technologies such as text-to-speech and speech-to-text. After an upload, our system uses advanced algorithms to produce accurate transcriptions in the user's language of choice, guaranteeing consistency and precision. To improve their comprehension and engagement, viewers may also generate more context based on the themes and subjects of the movie to go deeper into the information. Users don't need to watch lengthy videos to rapidly understand the main ideas and insights because of the availability of brief descriptions. It also expands the content of the video which makes the user understand the content wisely. Our platform's multilingual voiceovers, which allow users to hear the information in their native tongue, are one of its best features. This promotes diversity and breaks down barriers caused by language barriers. Our platform prioritizes security, privacy, and user-friendliness to make video content globally accessible and understandable, appealing to a wide spectrum of viewers with different language preferences and proficiency levels.

### **NOVELTY:**

Our project is novel because it takes a complete approach to removing language barriers and improving accessibility for video material by integrating state-of-the-art technology seamlessly. Our idea integrates translation, summarization, multilingual voiceover creation, transcription, text generation, and contextual understanding into a single, unified platform, in contrast to other systems that might only offer one area of language translation or transcription. In addition to meeting the various demands of users, this all-encompassing strategy guarantees a more inclusive and engaging watching experience. Our platform is unique in that it bridges language gaps and promotes worldwide connection and comprehension through video content, setting it apart from other solutions with its emphasis on user-friendliness, security, and privacy.

### **BRIEF EXPLANATION ABOUT THE INVENTION:**

The basic components required for building this application are:

- User Interface (UI) and User Experience (UX) Design
- Text-to-Speech (TTS) Technology.
- Speech-to-Text (STT) Technology.
- Multilingual Translation.
- Generative Pre-Trained Transformer (GPT).

- Video Scrapping.
- Text Summarization.
- Text generation.
- Audio extractor.
- Python libraries.

## **WORKING:**

This website will be divided into Seven parts. This will get the video or video URL as input and provide the transcription of the video, translated audio as the user wishes, Summary of the video, and text expansion of the video. This will be useful to the user for a better understanding of the video content.

### **A. Video Scrapping:**

This website allows users to download videos either by providing URLs or directly uploading videos to the platform. When a user submits a URL, the system starts the download by utilizing the supplied URL to obtain the video material immediately. The system downloads the provided video files straight from external sources for videos that are posted to the site. Title, duration, and available formats are not among the metadata that are extracted throughout the procedure. Rather, it just downloads the video file in its original format. The system downloads the video material with default settings; it does not provide a way to choose preferred video formats or quality levels. Although the website manages downloads from URLs and user uploads with efficiency, it lacks advanced information extraction or format selection tools, giving consumers a simple download experience.

### **B. Extraction of Audio:**

The video that was scrapped will act as an input here. A process to extract audio from a video file. To do this, the system will read the audio portion of the video and save it independently. Through this extraction process, users can retrieve only the audio material from the original video file by separating the audio data from the video data. Usually, the extracted audio is saved in MP3 format.

### **C. Transcribing the Audio:**

The extracted audio will be used here to transcribe the audio. It begins by setting up logging for debugging purposes. Next, for computational efficiency, a transcription model is initialized with certain characteristics. With a predetermined beam size, a beam search method is utilized to transcribe the audio content using the initialized model. Data about the identified language and the probability that goes along with it are retrieved and shown. Finally, the transcribed text segments are processed and written into a text file. In general, the procedure for setting up and applying a transcription model to transform audio information into text.

#### **D. Translation:**

The transcribed file, summarized text, and generated text will be the input in the process. A procedure for translating text into languages. To translate text segments into the intended target language, it makes use of a translation service. Text segments are read from an input file, each segment is translated separately, and the translated text segments are saved in an output file. The target language code, which indicates the preferred language for translation, must be entered by the user. Users can access the translated content to gain a better understanding of the video's content by saving the translated parts in a separate file. The translated summary and text generation will also be translated.

#### **E. Text To Speech:**

The translated files, translated summary, and translated generated text will be converted into speech. It is a procedure for employing a text-to-speech (TTS) engine to transform text into speech. It reads content from an input file that most likely has text that has been translated into other languages. After that, the text is processed to get rid of any extraneous formatting and newline characters. The processed text and the specified language parameter are then passed to the TTS engine. The engine produces a voice in the designated language that corresponds to the text that has been supplied. Lastly, the produced speech is stored at the designated output file location as an audio file, usually in MP3 format. All in all, translated text segments may be turned into audible speech, offering a different way to access and listen to text content. The speech audio of translated files, translated summary, and translated generated text will be stored in separate MP3 format.

#### **F. Text Summarization:**

In this process, the transcribed text will be summarized. A technique for automatically condensing text information that has been taken from a file. It starts by reading an input file, most likely including a lengthy text document. Then calculates the number of sentences to include in the summary, typically set to 25% of the total number of sentences in the input text. It then initializes a tokenizer and parser to handle the English text. It then applies the LexRank algorithm, which is used as a summarizer, to provide a text summary. The most significant sentences from the input text are chosen for the summary based on their relevance and significance. The generated summary sentences are printed to offer a clear overview of the primary ideas and pertinent details found in the original text content. All things considered, the ability to automatically summarize text material makes it simpler to understand and extract information from long publications. The summarized text will move to the translation part. The summary will be translated into the desired language and it will be moved to the text-to-speech part to convert the translated summary into speech and store it in a separate MP3 format.

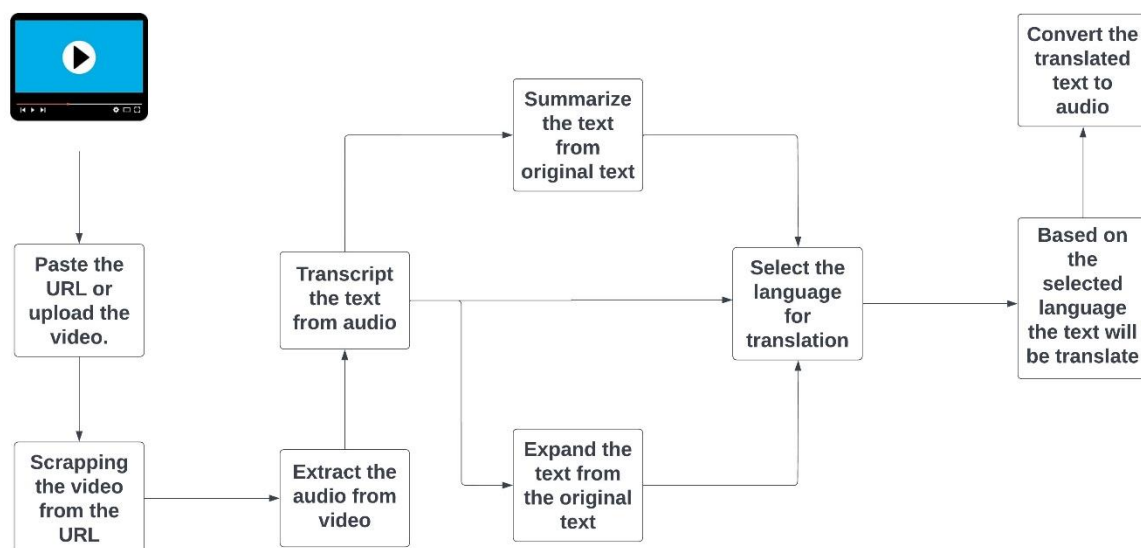
#### **G. Text Generation:**

In text generation, it uses transcribed text as input. A method of producing text using a GPT-2 language model that has already been trained. The GPT-2 tokenizer and model are first loaded using pre-trained weights. The task of the tokenizer is to transform the input text into numerical tokens that the model can process. On the other hand, the model is a neural network that has been trained to anticipate the following word in a string of text tokens. reads the contents of an input file, mostly containing a prompt or context for text generation, after the tokenizer and

model have been loaded. The tokenizer is then used to tokenize the input text, resulting in a series of numeric token IDs. This tokenized input sequence is passed to the GPT-2 model, which generates a sequence of text tokens representing the continuation of the input text.

To ensure diverse and coherent text generation, the model's generation process is regulated by several parameters, such as the maximum length of the generated text, the number of beams employed in beam search, and restrictions on repetitive n-grams. Lastly, all unique tokens are eliminated and the resultant text tokens are decoded back into text that can be read by humans using the tokenizer. Based on the given input prompt or context, the generated output is represented by the text that results. Similar to the summarization, the generated text will be translated into the desired language of the user and it will be moved to the text-to-speech part to convert the translated generated text into speech and store it in a separate MP3 format.

### **Flowchart of the process:**



### **SUMMARY OF THE INVENTION:**

The creation is a feature-rich web platform intended to improve accessibility to video content and overcome linguistic obstacles. The platform includes several features including transcription, translation, summarization, multilingual voiceover production, text generation, and contextual understanding. It integrates cutting-edge technologies like speech recognition, natural language processing, and machine learning. In contrast to current solutions that frequently concentrate on certain facets of language processing, our platform offers a comprehensive strategy, guaranteeing smooth communication across linguistic barriers for audiences throughout the world. Featuring an intuitive user interface, strong security features,

and an emphasis on increased interaction, the innovation is a trailblazing approach to promoting inclusivity and global understanding in the consumption of video material.