

Video Audio Replacement Proof of Concept (PoC)

Introduction

This document provides a detailed explanation of the Video Audio Replacement Proof of Concept (PoC) application. This innovative solution leverages advanced AI technologies to enhance video content by replacing suboptimal audio with an improved, AI-generated version.

Key Features

1. **Video Upload:** Users can upload video files in various formats (mp4, mov, avi).
2. **Audio Transcription:** Utilizes Google's Speech-to-Text API for accurate transcription with word-level timestamps.
3. **Text Correction:** Employs Azure's GPT-4o model to correct grammatical errors and remove filler words.
4. **Text-to-Speech Conversion:** Uses Google's Text-to-Speech API with the Journey voice model for high-quality audio generation.
5. **Intelligent Audio Synchronization:** Aligns new audio with the original video using advanced timing algorithms.
6. **Audio Quality Analysis:** Provides metrics on both original and processed audio for quality comparison.
7. **User-Friendly Interface:** Built with Streamlit for an intuitive, web-based user experience.

How It Works

1. Video Upload and Processing

- Users upload a video file through the Streamlit interface.
- The application extracts the audio from the uploaded video.

2. Audio Transcription

- The extracted audio is sent to Google's Speech-to-Text API.
- The API returns a detailed transcription with word-level timestamps.
- The original transcription is displayed to the user.

3. Text Correction

- The transcription is sent to Azure's GPT-4o model.
- The model corrects grammatical errors and removes filler words.
- The corrected text is displayed to the user for comparison.

4. Text-to-Speech Conversion

- The corrected text is processed by Google's Text-to-Speech API.
- The API generates high-quality audio using the Journey voice model.

5. Audio Synchronization and Replacement

- The new audio is aligned with the original video using the word timestamps.
- Advanced algorithms ensure proper synchronization, including:
 - Splitting audio on silence for precise word alignment.
 - Adjusting chunk durations to match original timing.
 - Adding fade-in and fade-out effects for smooth transitions.

6. Audio Quality Analysis

- Both the original and new audio are analyzed for quality metrics:
 - Signal-to-Noise Ratio (SNR)
 - Spectral Centroid (brightness)
 - Spectral Bandwidth
- These metrics are displayed for user comparison.

7. Final Output

- The processed video with replaced audio is rendered and displayed in the Streamlit interface.
- Users can play the video directly in the browser to verify the results.

Technical Components

1. **Streamlit**: Powers the user interface and handles file uploads.
2. **MoviePy**: Used for video and audio file manipulation.
3. **Google Cloud Speech-to-Text API**: Provides accurate transcription with timing information.
4. **Azure OpenAI API (GPT-4o model)**: Performs advanced text correction.
5. **Google Cloud Text-to-Speech API**: Generates high-quality synthetic speech.
6. **PyDub**: Enables precise audio manipulation and silence detection.
7. **Librosa**: Used for advanced audio analysis and quality metrics calculation.

How to Use

1. Setup:

- Ensure all required libraries are installed (`streamlit` , `moviepy` , `google-cloud-speech` , `google-cloud-texttospeech` , `openai` , `pydub` , `librosa`).
- Set up Google Cloud credentials for Speech-to-Text and Text-to-Speech APIs.
- Configure Azure OpenAI API key and endpoint.

2. Running the Application:

- Execute the script to start the Streamlit server.
- Access the application through a web browser.

3. Processing a Video:

- Upload a video file using the provided interface.
- Click the "Process Video" button to start the enhancement process.
- Wait for the processing to complete (duration depends on video length).

4. Reviewing Results:

- Compare the original and corrected transcriptions.
- Analyze the audio quality metrics for both versions.

- Play the processed video to verify audio-visual synchronization and overall quality.

Potential Applications

1. **Content Creation:** Enhance poorly recorded videos or audio for professional-quality output.
2. **Education:** Improve lecture recordings or educational content for better clarity.
3. **Podcasting:** Clean up interview recordings by removing filler words and improving audio quality.
4. **Video Localization:** Facilitate easier dubbing of videos into different languages.
5. **Accessibility:** Generate clearer audio for hearing-impaired viewers.

Limitations and Considerations

- Processing time may be significant for longer videos.
- Perfect synchronization may be challenging for videos with rapid speech or significant background noise.
- API usage costs should be considered for large-scale applications.
- Privacy considerations when handling potentially sensitive video content.

Future Enhancements

1. Implement more sophisticated audio alignment algorithms (e.g., dynamic time warping).
2. Add user controls for adjusting speech rate, pitch, and voice selection.
3. Introduce a progress bar for processing longer videos.
4. Enhance error handling for API failures or unexpected inputs.
5. Optimize for processing longer videos or handling batch processing.

Conclusion

This Video Audio Replacement PoC demonstrates the powerful capabilities of combining multiple AI technologies to enhance video content. By automating

the process of transcription, correction, and audio replacement, it offers a unique solution for improving video quality with minimal manual intervention.