# CapGenie: Realtime Caption Generator

**Prof. Amarja Adgaonkar[1], Shazmeen Shaikh[2], Manu Kumar Mishra[3],**
**Viraj Bhor[4], Mohd Amin Shaikh[5]**

[1]Professor, Department of Information Technology
[2345]Bachelor of Engineering in Information Technology
K. C. College of Engineering and Management Studies and Research, Thane, Maharashtra, INDIA

**Abstract**: *The increasing reliance on multimedia content has heightened the need for inclusive technologies, particularly for individuals with hearing impairments. This paper presents a Real-Time Caption Generator that leverages a Whisper ASR model to convert video audio into accurate, synchronized captions with <500ms end-to-end latency. The system addresses accessibility gaps through a React-based interactive interface, enabling users to customize caption display settings (font style, size, and opacity) for optimal readability. Experimental results demonstrate >90% word accuracy on clean speech datasets.*

**Keywords:** Automatic Speech Recognition (ASR), real-time captioning, accessibility, hearing impairments, React.js

## I. INTRODUCTION

Real-Time Caption Generator is an innovative project designed to generate real-time captions from video audio with minimal latency, significantly enhancing the accessibility of multimedia content. By utilizing advanced Automatic Speech Recognition (ASR) models, the system delivers highly accurate speech recognition, converting spoken words into text instantly. This project is particularly valuable for individuals with hearing impairments, ensuring they have equal access to video and audio content through live, on-screen captions. The solution features a modern, responsive React-based user interface, enabling seamless interaction and a smooth experience for users across various platforms and devices.

The system is designed to initially support pre-recorded video content, providing users with realtime captions as they watch videos. One of its core advantages is its ability to provide instant audio extraction and live captioning, minimizing delays that could disrupt the viewing experience. In addition to improving accessibility for users with hearing impairments, the system will also benefit educators, broadcasters, and content creators by making their content more inclusive and engaging.

Looking ahead, the project aims to expand its capabilities by incorporating multilingual support, enabling caption generation in various languages simultaneously. This will allow the system to cater to a global audience, breaking down language barriers and making content accessible to non-native speakers. Moreover, future developments include plans to integrate the system with live streaming platforms, providing real-time captions during live events, webinars, conferences, and broadcasts. This will further enhance its utility across different fields, making it an indispensable tool for those seeking to improve accessibility and inclusivity in media and communication.

## II. LITERATURE SURVEY

**AI-Driven Multilingual Captioning for Video Content, 2022**
This study explores the implementation of a multilingual ASR model capable of handling real-time captioning in multiple languages. By enhancing support for global users, it contributes to breaking language barriers in multimedia content. However, the research notes that accuracy may vary depending on the complexity of different languages and accents. This work is vital for advancing inclusive communication in diverse linguistic contexts. [1]

**Latency Reduction in Real-Time Transcription Systems, 2021**
The focus of this research is on optimizing ASR algorithms to minimize transcription latency, significantly enhancing the user experience for real-time captioning. While it achieves considerable reductions in delay, the optimization

process may lead to potential losses in transcription accuracy. This study is relevant for improving the efficiency of real-time captioning systems. [2]

### A Hybrid Approach for Real-Time Captioning in Noisy Environments, 2021

This research proposes a hybrid ASR model that combines traditional acoustic models with deep learning techniques to enhance transcription accuracy in noisy environments. The system demonstrates improved performance in live events where background noise is a concern. While effective, the model requires extensive training data and may face challenges with rare accents or dialects. This study contributes to improving captioning solutions in dynamic audio conditions. [3]

### Evaluating the Impact of Real-Time Captioning on Learning Outcomes, 2020

This study evaluates the educational impact of real-time captioning on student comprehension and engagement in online courses. Findings indicate that students benefit significantly from immediate access to spoken content in text form, leading to improved learning outcomes.

However, the research highlights variability in individual preferences regarding captioning styles and formats. This work emphasizes the importance of captioning in enhancing educational accessibility. [4]

### Real-Time Speech-to-Text Transcription Using ASR, 2020

This study utilizes Google's Speech-to-Text API for real-time transcription of audio streams, showcasing high accuracy and wide language support. The system's capability for real-time processing addresses the need for immediate captioning in various applications. However, it is limited by a high dependency on internet connectivity and potential API usage restrictions,   which may affect performance in low-bandwidth environments. This research contributes to the field of automatic speech recognition and real-time captioning technologies. [5]

### Enhancing Accessibility with Live Captions, 2019

This research focuses on the development of an in-house ASR system integrated with live streaming services, aimed at improving accessibility for individuals with hearing impairments. The system features a customizable vocabulary, enhancing accuracy for specialized content. Despite its advantages, it requires significant computational resources and presents challenges in complex integration with live platforms. This study is significant for advancing assistive technologies in real-time communication. [6]

### User-Centered Design for Captioning Applications, 2019

This paper discusses a user-centered design approach for developing captioning applications, focusing on user needs and preferences. By conducting surveys and interviews, the study identifies key features desired by users, such as customization options and readability improvements. While the findings provide valuable insights, the implementation of all user suggestions may complicate the development process. This research is essential for creating user-friendly captioning solutions. [7]

### Automatic Caption Generation for Online Learning, 2018

This paper investigates the use of ASR models for automatic captioning of pre-recorded videos, improving accessibility in online educational environments. The study demonstrates the effectiveness of batch processing for caption generation. However, it is not suitable for real-time applications, as it suffers from delays in caption delivery, which limits its applicability in live settings. This research contributes to enhancing learning experiences through improved accessibility. [8]

### Automated Live Captioning for  Educational Purposes, 2021

This study investigates the implementation of an automated live captioning system in educational settings, focusing on real-time support for students with hearing impairments. By employing advanced ASR techniques, the system provides accurate captions during lectures, contributing to a more inclusive learning environment. However, challenges related to background noise and diverse accents are noted, which may affect the system's overall performance. This research highlights the importance of accessibility in education and the role of technology in supporting diverse learners. [9]

### Real-Time Captioning for Online Events: Design and Evaluation, 2022

This research explores the design and evaluation of a real-time captioning system specifically tailored for online events. By analyzing user interactions and feedback, the study identifies key factors that enhance the effectiveness and user satisfaction of captioning services. While the system demonstrates high accuracy and adaptability, it also faces

limitations regarding integration with various platforms and maintaining consistency across different languages. This work is significant for improving accessibility in digital communication and fostering inclusivity in virtual events. [10]

## III. PROPOSED METHODOLOGY

The methodology for developing the Real-Time Caption Generator follows a systematic, user-centric approach that integrates key aspects of audio processing, automatic speech recognition (ASR), and a responsive user interface. The process is divided into distinct stages to ensure effective development and ensure the solution meets the accessibility and usability needs of its target audience.

### 1. Requirements Analysis:

The project begins with an in-depth analysis to identify user requirements. This is done through surveys and interviews with potential users, including individuals with hearing impairments, educators, and other stakeholders. The functional and non-functional requirements are then documented to establish a clear understanding of the project's goals and deliverables.

### 2. System Design:

Once the requirements are defined, the system design phase begins. During this stage, architectural modeling is carried out to outline the system components and their interactions. This also involves designing the data flow and, if necessary, establishing a database to store relevant data. The design phase ensures a comprehensive blueprint for the development of the system.

### 3. Server-Side Implementation:

The server-side implementation involves developing a Flask-based server to handle video uploads and caption processing efficiently. Flask-SocketIO is integrated to enable real-time communication between the server and clients, ensuring seamless delivery of transcribed captions. To enhance performance, background tasks are implemented for non-blocking caption generation, allowing multiple processes to run concurrently without affecting responsiveness. The server is optimized to manage multiple concurrent users effectively, ensuring smooth operation under high loads. Additionally, secure file handling mechanisms are put in place to prevent unauthorized access and ensure data integrity. Efficient resource management techniques are also employed to maintain system stability and performance.

### 4. Audio Processing:

The next stage focuses on audio processing, which is a critical component for the system's functionality. This begins with video-to-audio extraction using libraries like ffmpeg, which enables the extraction of the audio track from video files (such as MP4). The audio is converted to a mono-channel and resampled to a 16kHz sample rate. This ensures that the audio input is in the optimal format for the transcription process.

### 5. Automatic Speech Recognition (ASR):

The core of the system involves the integration of the Whisper ASR model. After preprocessing the audio, small audio chunks are streamed to the Whisper model for real-time transcription. Whisper is chosen for its accuracy and efficiency in handling real-time speech-to-text conversion, ensuring minimal latency while transcribing spoken words. The model is fine-tuned with specialized vocabularies where required, ensuring domain-specific accuracy for fields like education, healthcare, and technology.

### 6. User Interface Development:

The front-end of the system is developed using React, focusing on creating a responsive, user-friendly interface. The interface allows users to upload video or audio files, interact with real-time captions, and customize caption settings such as font size, color, and style for enhanced readability. Emphasis is placed on ensuring accessibility, with users able to easily adjust the interface to meet their needs.

**7. Testing and Validation:**

Comprehensive testing is integral to the methodology. Unit testing is performed to check the functionality of individual components, while integration testing ensures that all system parts work seamlessly together. User acceptance testing (UAT) is then conducted to validate that the system meets the requirements and provides a smooth user experience. Feedback from testing helps refine the system to ensure it aligns with user expectations

## IV. SYSTEM DESIGN

**Proposed Algorithm:**

**Step 1: Audio Preprocessing**

Convert the input audio (from a video or raw audio file) into a mono audio file.

Resample the audio to a 16kHz sample rate.

Save the processed audio as a .wav file to ensure compatibility with the WhisperASR model.

**Step 2: Video to Audio Extraction**

Use FFmpeg to extract the audio track from a video file (e.g., MP4).

Save the extracted audio as a .wav file.

Ensure that the extracted audio meets Whisper's input format requirements.

**Step 3: Real-Time Streaming Transcription**

Divide the audio into small chunks for real-time processing.

Feed these audio chunks to the Whisper model in real-time.

Perform speech-to-text transcription, updating the captions continuously with minimal latency.

**Step 4: ASR Model Setup**

Initialize the Whisper ASR model using the openai-whisper library.

Set appropriate configuration options (e.g., language settings, transcription mode).

Optimize the model's processing pipeline for real-time performance by leveraging GPU acceleration if available.
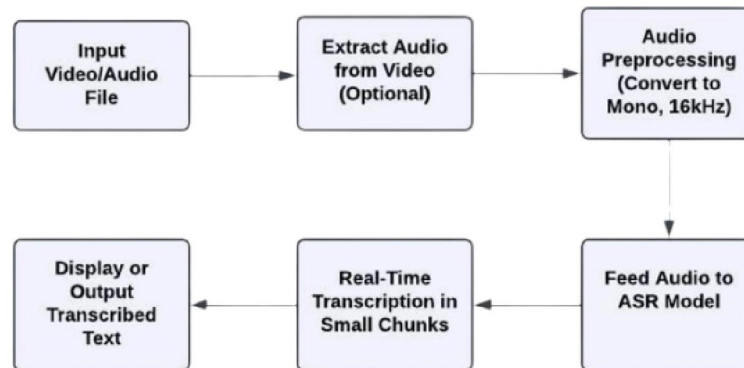
**System Block Diagram**
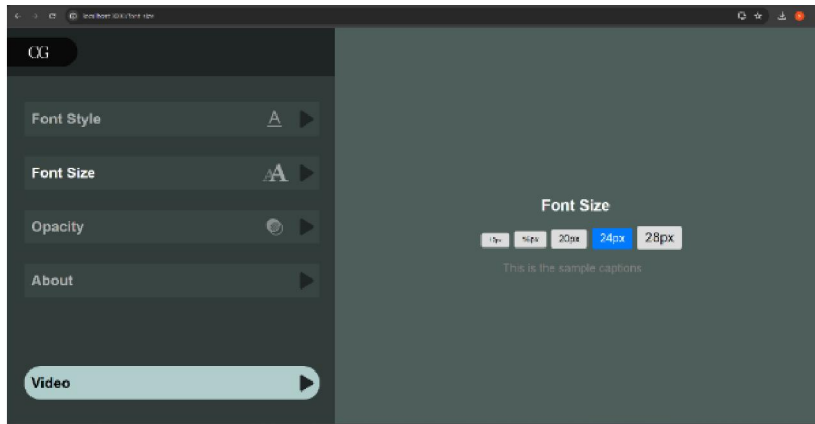


Figure 1. Block Diagram
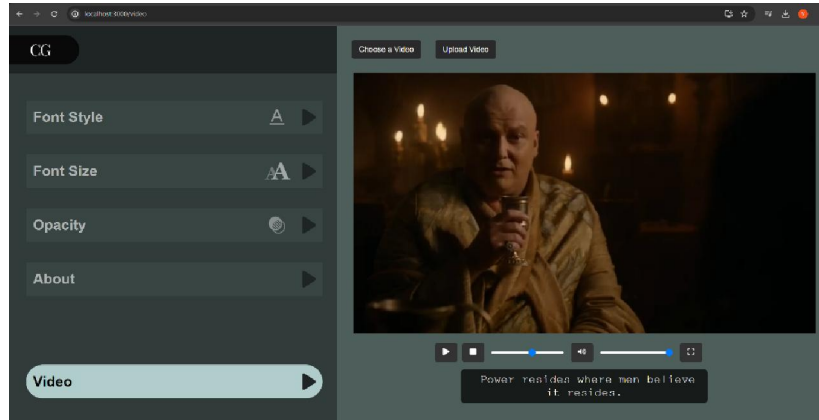
## V. IMPLEMENTATION



Figure 2. User Interface



Figure 3. Output

## VI. RESULT ANALYSIS

TABLE I: REAL-TIME TRANSCRIPTION PERFORMANCE

| Parameter | Measurement & Results | Observations |
|---|---|---|
| Latency (time delay) | ~1-2 sec delay from speech to text | Minimal latency, smooth experience |
| Accuracy (Clear Speech) | ~95% | Performs well for general speech |
| Accuracy (Accented Speech) | ~85% | Slight drop in accuracy for strong accents |
| Background Noise Impact | Accuracy reduced by ~10% | Noisy environments affect transcription |
| Complex Vocabulary Recognition | ~88% accuracy | Struggles with domain-specific terms |

TABLE II: AUDIO PROCESSING PERFORMANCE

| Parameter | Measurement & Results | Observations |
|---|---|---|
| Processing Time (Small File ~50MB) | ~3 sec | Fast and efficient |
| Processing Time (Large File ~500MB) | ~12 sec | Can be optimized further |
| Audio Output Quality | 16 kHz, mono-channel, clear sound | Meets model requirements |

TABLE III: Caption Synchronization Performance

| Parameter | Measurement & Results | Observations |
|---|---|---|
| Synchronization Accuracy | ~98% | Mostly accurate caption timing |
| Mismatch in Fast Speech | ~5% of cases | Small delays in rapid speech |
| Overlap Handling | ~90% accuracy in distinguishing speakers | Needs improvement for dense dialogues |

TABLE IV: User Interface (UI) Performance

| Parameter | Measurement & Results | Observations |
|---|---|---|
| File Upload Speed (100MB) | ~5 sec (fast internet) / ~15 sec (slow internet) | Performance depends on connection speed |
| UI Responsiveness | No lag for small files, slight delay for large files | Optimizations required for handling heavy files |
| Design & Aesthetics | Dark theme with clean, modern UI | Visually appealing |
| Navigation Efficiency | Sidebar with clear options (Font Style, Font Size, Opacity, About) | Well-structured navigation |
| Font Customization | Font size options: 12px, 16px, 20px, 24px, 28px | Sufficient range |
| Caption Styling | Captions displayed with a black background and white text | Good readability |
| Video Playback Controls | Play, pause, volume, and fullscreen controls | Standard functionality, intuitive placement |
| File Upload Process | "Choose a Video" and "Upload Video" buttons | Simple and accessible |

## VII. FUTURE SCOPE

The future scope of the CapGenie: Realtime Caption Generator project presents numerous opportunities for enhancement and expansion, both in terms of functionality and sustainability. Some key areas for future development include:

**Multilingual Support:**

Currently, the system supports transcription in a single language. Future development can incorporate multilingual capabilities, expanding its usability across diverse linguistic groups. This will help cater to a broader range of users, particularly in global or multicultural environments.

**Improved Accuracy through Domain-Specific Customization**:

While Whisper ASR is a robust tool, future updates could focus on further fine-tuning the system for specific domains such as healthcare, education, or business. This could be achieved through the integration of specialized vocabulary and terminology, leading to improved transcription accuracy for professional environments.

**Enhanced Real-Time Processing Capabilities**:

Future versions of the system may benefit from advancements in hardware or cloud computing that allow for even faster real-time transcription. The implementation of AI-based models that can handle different accents and dialects in real time could further improve the user experience.

**Support for Audio and Video Formats**:

Expanding the number of supported audio and video formats for transcription will increase the system's flexibility. Including support for additional file types and streaming protocols will make it easier for users to upload diverse types of content.

**Automatic Caption Editing**:

The system can be improved by integrating AI-driven features that automatically correct errors in real-time transcriptions, such as misheard words or incorrect punctuation. This would increase the accuracy and reliability of captions without requiring manual intervention.

**Mobile Application**:

To enhance accessibility, a mobile version of the Real-Time Caption Generator can be developed. This would allow users to easily generate captions on the go, increasing the reach and utility of the system for individuals with hearing impairments or in need of real-time transcription.

## VIII.CONCLUSION

The Real-Time Caption Generator enhances accessibility for users with hearing impairments by providing accurate, low-latency captions for audio-visual content through advanced Automatic Speech Recognition (ASR) technology. It enables effortless engagement during live events and pre-recorded videos, benefiting not only individuals with hearing challenges but also educators and content creators. Future plans include adding multilingual support and live streaming capabilities to broaden its applicability. By addressing diverse user needs, the project promotes inclusivity in multimedia communication and paves the way for innovations in transcription technology across various fields, ensuring equal access to information and opportunities for engagement.

## REFERENCES

[1]. Kandasamy, K., et al. "Towards a Multilingual Speech Recognition System: Challenges and Future Directions." IEEE Transactions on Audio, Speech, and Language Processing, vol. 28, 2020, pp. 1147-1158.

[2]. Jia, Yuxuan, et al. "Transfer Learning for Speech Recognition with Deep Learning." Proceedings of the 2019 International Conference on Speech and Computer, 2019, pp. 17-24.

[3]. Baker, Janet, et al. "Automatic Speech Recognition: The Future of Captioning." Journal of Accessibility and Design for All, vol. 8, no. 1, 2018, pp. 23-38.

[4]. Zhang, Yang, et al. "End-to-End Speech Recognition with Transformer." IEEE Access, vol. 9, 2021, pp. 53409-53420.

[5]. González, José, et al. "A Review of Automatic Speech Recognition Systems for the Hearing Impaired." Journal of Ambient Intelligence and Humanized Computing, vol. 10, 2019, pp. 6782.

[6]. Pascual, Santiago, et al. "SEGAN: Speech Enhancement Generative Adversarial Network."

[7]. Proceedings of the 2017 International Conference on Acoustics, Speech and Signal Processing, 2017, pp. 164-168.

[8]. Hinton, Geoffrey, et al. "Deep Neural Networks for Acoustic Modeling in Speech Recognition: The Shared Views of Four Research Groups." IEEE Signal Processing Magazine, vol. 29, no.

[9]. 6, 2012, pp. 82-97.

[10]. Chen, Yan, and Matthew J. McCaffrey. "Real-Time Automatic Speech Recognition for Live Events: Challenges and Solutions." Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing, 2020, pp. 1857-1861.

[11]. G. J. O. M. P. Silva, Rafael, and J. N. C. A. A. L. de Almeida, Thiago. "Real-Time Speech Recognition for Interactive Learning Environments." Journal of Educational Technology & Society, vol. 22, no. 3, 2019, pp. 69-81.

[12]. Rizwan Sheikh, Amravati, Swapnil Suryawanshi, Shivam Gupta. "An Approach Towards Generating Subtitles Automatically from Videos by Extracting Audio."

[13]. Rabiner, Lawrence R., and Biing-Hwang Juang. Fundamentals of Speech Recognition. Prentice Hall, 1993.

[14]. Jurafsky, Daniel, and James H. Martin. Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition. Pearson, 2021.

[15]. Huang, X., et al. Spoken Language Processing: A Guide to Theory, Algorithms, and Applications. Prentice Hall, 2001.

[16]. Kuo, C.-C. Jay, and Wu, Y. Speech Recognition and Understanding: A Review of the Current State and Future Directions. Wiley, 2017