

# TALKSYNC: AI POWERED REAL TIME SPEECH TO TEXT TRANSLATION FOR VIRTUAL MEETINGS

Niraj Bhausaheb Sabale, Dhyan Ramesh Parasiya, Rushikesh Prabhakar Mhatre, Aniket Tatyaba Sarvade  
*BE.IT, Pillai HOC College Of Engineering and Technology, Navi Mumbai, India*

**Abstract--** Virtual meetings have turn out to be an important mode of verbal exchange in today's digital global, enabling remote collaboration by means of transmitting audio and video records over the net. However, language barriers and accessibility demanding situations can abate powerful verbal exchange. To deal with this, our project, goals to enhance digital communiqué via offering actual-time speech-to-textual content translation and automated assembly summarization. The platform permits contributors to pick their preferred spoken and subtitle languages, ensuring seamless multilingual interaction. Using advanced AI models together with Whisper AI and integrating numerous APIs, the machine transcribes spoken content material in real-time and shops it in a database. Additionally, an wise backend approaches the transcript to generate concise assembly summaries routinely. The mission is constructed the use of Angular for the user interface and deployed on Vercel for efficient virtual meeting management. This answer complements accessibility, collaboration, and productivity in virtual meetings, this mission is call as Talk sync in which realtimespeech to textual content translation is executed.

**Index Terms -** Real-time speech-to-text, AI-powered translation, Speech recognition, Multilingual communication, Natural language processing (NLP), Virtual Meetings

## I. INTRODUCTION

Virtual conferences have become an important part of how people and agencies work together in today's virtual world. With the increasing recognition of remote work, to gain online knowledge and global business operations, these systems allowed the groups to connect immediately, no matter what they were. They offer real -time sounds and video oral exchanges online, making the interaction more efficient and practical [13], [4], but in addition to virtual conferences, in addition to this, language barriers, information discussions in information discussions and surprise of facts include situations such as situations. Checking these problems is prominent to make online oral exchanges more powerful and accessible to any body.

From simple sound conferences to AI-operated collaborative tools [4], [13], Virtual Meeting has developed with a lot of time. Originally, they rely on telephone -based sound conference, which limited the interaction. With the use of excessive speed networks, La video conferencing

structures such as Zoom and Microsoft team added features that include screen sharing, chat and cloud recording, which is more attractive for collaboration. The latest progressive benefit from AI for real-time transcript, computer-controlled translation and mounting summaries, and ensures better access and productivity. As the time goes on, Destiny reforms such as AR and VR conferences will also increase digital interactions, making them more engrossing and increasingly.

Our AI-operated things specialize in two major abilities: real-time speech reading material transcription and automatic meeting summary. Transcript characteristics, operated from Whisper AI, immediately convert phrases to text material [1], [2]. This allows individuals to see live subtitles, break language limits and interact on hand. The system also collects transcription in a database so that they can be reviewed later. Summary characteristics take these ties and use advanced language treatment techniques to remove the main points, leading to a short pre -collection of the meeting [15]. These guides want word technology and guarantee that members can easily review the important discussion with a full copy. By combining these tasks, the platform makes our digital conferences more green, inclusive and effective..

It discusses the structure, implementation and distribution of the Paper Talk wash, and emphasizes its scalability, protection and adaptability for diverse digital meeting environment [8]. Breaking language barriers, promotes synchronization and improves cooperation around the world. The proposed AI-operated speech-to-text translation unit is designed for spontaneous integration and consumer-enjoyable operations at virtual meetings. It has a personal-friendly interface that is integrated with well-known platforms such as Zoom and Microsoft team, and ensures extra layout [9]. Real -time seizures enables uninterrupted speech treatment, even better speech starts their favorite speakers and subtitle languages, allowing seamless multilingual communication [6]. Backnd automatically produces a brief summary, which helps users to review the big discussion points quickly. This technique

increases access and productivity in virtual meetings, especially for international teams.

#### Background of Speech-to-Text Translation:

Speech-to-text (STT) technology has evolved significantly due to artificial intelligence (AI), Natural Language Processing (NLP) and progress in deep learning. Previously, the STT system was a rule-mainly-based and statistical model, which required Big Guide School Education and forced accuracy [16]. Later, the construction of the hidden Markov model (HMM) and the Gaussian Mixer Model (GMM) presented the reputation, but these techniques were still facing challenges that include background noise, accent and fast speech. At the forefront, the lip -based -based -based perfect speech, especially recurrent nerve networks (RNN), triggered long -term short -term memory (LSTM) and transformer -based architecture such as Whisper AI, VAV 2 Vec 2.0 and deep speech. These models benefit from large -scale data sets and self -protected mastery to improve accuracy. In addition, Sky-Mainly-based STT structures enables real-time transcript and translation with minimal delay, making them ideal for digital meetings

Recent reforms in speech-to-reading material translation and automatic summary of mounting have added many AI-operated strategies. End-to-end ASR models such as Whisper AI and VAV 2 VEC 2.0 wants different acoustic, pronunciation and language fashion, and more correct transcription in the coming. An operated machine translation uses pre-trained multilingual fashion, as well as mart and m2m- with a hundred, to convert speech tape into unique languages, allows smooth passenger-language oral exchange. [15] Real-time-speech-to-text streaming is made easier through technology such as stream video API and WebrTC, which allows immediate transcription and translation with minimal putters. The automatic meeting summary extracts (selection of larger sentences) and abstract (rewriting summary) strategies, often driven through GPT elements based and Burt NLP fashion. In addition, the speaker Diary generation allows the difference between multiple sound systems, making the tape more structured and readable. The platform provides great blessings, including scalability for integration with accurate transcription, seamless language translation, AI-generated summaries and current digital meeting platforms. Study at the Project:

This assignment integrates advanced AI technology to ensure the speech recognition of high compatibility, real-time translation and automatic condensation. By using a Whisper AI for speech-to-text-material conversion, a multilingual NLP version for translation, and LLM-based perfectly perfectly text material summary, the gadget can automate the entire meeting documentation process. The

platform offers major blessings including unique transcription, seamless language translation, AI-Birthed Summary and scalability for integration with existing digital mounting structures. By taking advantage of AI, NLP and Cloud-Base. Related work

#### Speech-to-Text (STT) Implementation:

The research thesis is a function in real-time speech-to-speech translation for virtual meetings, which uses the system to get fashion to the system. It detects strategies that include automated speech recognition (ASR), nerve gadget translation (NMT) and textile content-to-eating (TTS) synthesis, enables uninterrupted multilingual oral exchange. Inspection possibly appoints models such as Whisper AI for ASR, transformer-based perfect architecture for NMT and advanced TTS system for natural voice output. ] In addition, it illuminates the conditions for real -time treatment, language shades and scalability requirements.

This research thesis deals with the improvement of a speech-recycling translation forum to facilitate communication among students with different language backgrounds, [2] specializes in English and Yoruba language. The platform uses asp.net with C# for growth, [5] for a responsive design signals CSS, Bootstrap, Jquery and JavaScript, and a SQL Server database for Safe Record Garage. The object-oriented feature (OOM) was hired for the size of the software, which aims to deliver the consumer-enjoying interfaces that allow real-time interaction and bridges the communication hole between the English and Yoruba sound systems.

#### Text Translation Implementation:

By using Angular, Django and Whisper API, you will integrate a translation supplier into the workflow. [12] First, Whisper API may be responsible for changing the speech in the text in real time. This transferred text is then sent to Django Backend, where it is treated for translation. You can use the Google Translation API, Deepl API or a self -confidence translation model to convert the lesson into the preferred language. ] In Django, after receiving transcape text content from Whipper API, the machine will detect supply language and bypass the text through a translation supplier. If the use of Google translates API, Backand API can send the requests to translate the course content correctly. The translated response will then be stored and will be designed for the use of Django channels for oral exchange in real time. On angular fronts, an online contact connection or a HTTP vote mechanism can be used to bring translated text dynamically. [4] UI should be

designed to display captions in many languages with minimal delay. Users should also be able to take their favorite translation language. Further adaptation, which involves constant translation and handling of language detection errors, can improve efficiency. [12] By integrating these additives in the first place, Talkync may be able to offer strong multilingual communication in virtual meetings, allowing a real-time experience of speech-to-surface.

## II. OVERALL SYSTEM DESCRIPTION

### A. Existing System

In digital conferences, specially those involving participants from exclusive areas and languages, conversation boundaries frequently get up due to language differences. These demanding situations can sluggish down selection-making, prevent collaboration, and exclude members who aren't fluent within the dominant language of the assembly. Moreover, existing actual-time transcription answers may also lack accuracy, conflict with more than one audio system or accents, or fail to combine seamlessly with digital assembly platforms, leading to inefficient communique.[1],[2].

Rev combines AI and human transcription services to provide high-quality real-time transcription, along side translation abilities for various languages. Sonix.Ai offers a comparable service with actual-time transcription, translation, and integration with gear like Zoom and Google Meet. Another sturdy contender is Trent, which also affords AI-based totally transcription offerings with actual-time captioning and multi-language assist, making it perfect for global groups[4],[5]. For extra customizable answers, Microsoft Azure Cognitive Services (Speech API) and Google Cloud Speech-to-Text offer effective, scalable gear for real-time speech-to-text conversion with help for a extensive variety of languages and integration flexibility, specifically for users who need more manipulate or have unique enterprise needs[6],[7]. These structures all offer strong answers for seamless, correct transcription and translation in digital meetings.

### Proposed System:

The proposed machine is designed for spontaneous integration in virtual meetings and user -friendly operations. It provides capacity for a user-dried interface that is integrated with popular platforms such as Zoom and Microsoft teams, which get easy to become easy with more layouts. Real -time sound price in real time allows spontaneous speech treatment, even advanced speech popularity appoints the machine's mastery to provide immediate and correct transcription.

Machine language supports translation so that customers can communicate effectively in specific languages. The speaker's identity improves transcript cladding by differences between a type of speaker. Customized performance options allow users to regulate font size, color and format for advanced clarity.

A stable cloud frastructure guarantees safe treatment and storage of mounting ties, privatizing facts. In addition, the automatic prismatic generation removes important factors and decisions, and published illiterates increase accessibility. These features make the machine a green, spontaneous and inclusive solution for digital communication.

The device will regularly stumble on more than one sound system in many languages and dialects and provide real -time texts. For groups around the world, talk washing can help rapid translation between special languages, which can ensure that simple oral exchange despite language restrictions. In addition, it will allow customers to customize translation settings, so that industry examination or conditions can be moved and translated effectively.

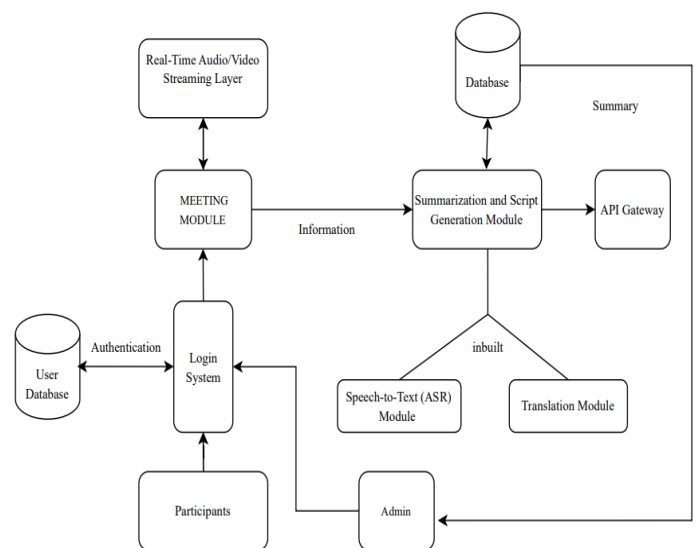


Fig.2. Architecture Diagram

This diagram represents the architecture of our project, Let's break it down:

### Key Components :-

#### 1. Login System (Future Enhancement)

In our system right now, users join meetings via a direct link without authentication. In future updates, login authentication can be implemented to allow only registered users to join.

#### 2. Meeting Module:

Manages the virtual meeting and it use vercel platform to connect more than one (peer-to-peer) participant in meeting. Module is connects with the real-time

audio/video streaming layer. It gathers information from the meeting and passes it to the Summarization and Script Generation Module. Stream Video API used in the Meeting Module Handles real-time video/audio streaming in virtual meetings. Ensures smooth communication between participants.

### 3. Summarization and Script Generation Module:

Receives information from the Meeting Module. Uses inbuilt Speech-to-Text (ASR) Module (powered by Whisper AI API) to convert speech to text. Uses the Translation Module for real-time language translation. Whisper AI API help in improving accessibility. Stores the generated summary in the Database and sends it via the API Gateway.

### 4. User & Admin Interaction:

Participants can join meetings and interact with the system. Admins have control over managing meetings and transcripts.

### 5. Databases

User Database: Stores authentication and participant details (for future login system). Main Database: Stores meeting transcripts, summaries, and translations.

## III. METHODOLOGY

It takes a look at the feature used in the improvement of AI-integrated real-time speech reading material translation tools, Talksinks designed for digital conferences. The feature includes several degrees, including statistics collection, preprosares, speech-to-text transcript, unit translation, real-time treatment, evaluation and perfection. In order to produce a strong and scalable speech-to-layer content and translation models, we created different data sets from some resources. Speech recognition information is publicly transformed into acquired from available data sets [5] Librispeech, Mozilla Common Voice, Ted-Liam and Meeting-Settle Corpora. The translation protocol was taken from parallel text Corpora, which included WMT, Opensubtitles and Europarl Corpus, which has been used to learn language fashion. In addition, multimodal data rate containing sound, ties and translation is included to increase relevant accuracy. Similarly, to decorate domain-perfect accuracy, [2] realistic meeting recordings were collected and transmitted manually.

Before the version of education, the collected facts had to undergo pre-compuls. ] Functional extraction was for the use of the flour frequency cepstral coefficient (MFCCS) and spectrograms, while employing the speaker derogation to separate the sound system in derivation calls. [10] Lesson spreklamation concerned about breaking, extermination of the phrase to create some cleanser input information, prevent generalization, lemmatization and punctuation generalization.

For real-time transcript, we have hired the modern deep mastering-mainly-based ASR model. Whisper API was used for perfectly accurate transcription, and utilized its better speech recognition skills. [4] Acoustic modeling did not do so -cooked the audio signal for foam and expression and expression of the use of LSTM and CNN. [5] Language modeling, BURT and GPT-based perfectly used model, transcription was integrated to decorate accuracy and expect relevant appropriate sentences.

In order to resume uninterrupted real -time treatment, several adaptation strategies were performed. The low -line model was used to reduce computational overheads, to include amounts of dystylbert. A hybrid method is used to beautify the general performance (local estimate) and shooter [16]. Streaming API integration was able to facilitate oral exchanges in real time using the WeBRTC, GRPC and Socket-based complete streaming.

The evaluation of the gadget is evaluated by the use of industry standard viewing measurements. Word fault rate (Wer) is used to assess the ASR transcription accuracy, [16] even blue and roser ranging were used to evaluate the high quality of the translations. Late measurement secured real -time treatment inside Milceconds for purposes in digital meetings. In addition, real -time recognition system with customers on live conferences was done through A/B check and commentary series to increase accuracy and accountability.

The very last talk sink is primarily distributed as a machine based on API and integrated into different platforms. Cloud Purpose is hosted on AWS, Google Cloud and Microsoft Azure for excessive availability and scalability. The system, together with the Zoom, Microsoft team and virtual meeting structures, together with Google are with Google. Utility's advance status developed the use of angular and secured the experience of a responsible and interactive person. In addition, passport platform support is ensured by developing a standalone tool using flute (mobile) and electrons (laptops) for maximum access.

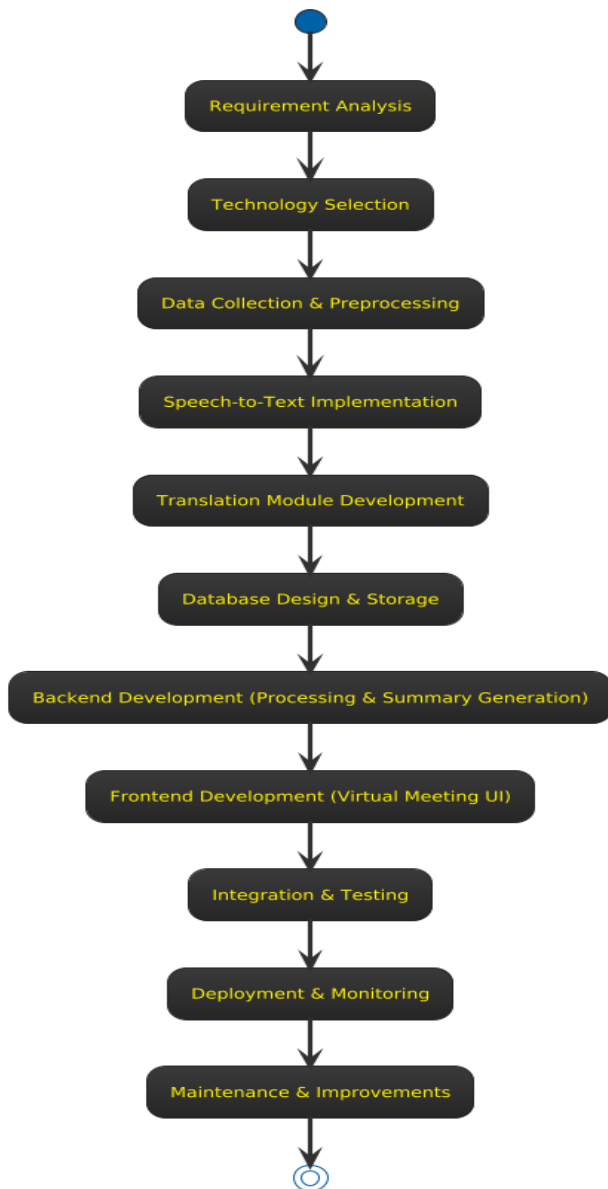


Fig.3 Development Stages.

#### IV. RESULT

##### Meeting Platform-

We have created a meeting platform using Versel, so that many participants from all over the world can attend a meeting. The user interface (UI) is developed using Angular, and ensures a smooth and interactive experience. For video conferences, we integrate stream video APIs, enabling high quality truth communication between participants.

##### Real time speech recognition -

For real -time speech recognition, we use Whisper AI API, which enables accurate speech text translation. This system

not only converts the speech language into a real -time text, but also stores lesson tape in a database for future reference. Automatic meeting summary, where the system calls the entire meeting, helps participants to review the discussion points quickly. Multilingual speech recognition.

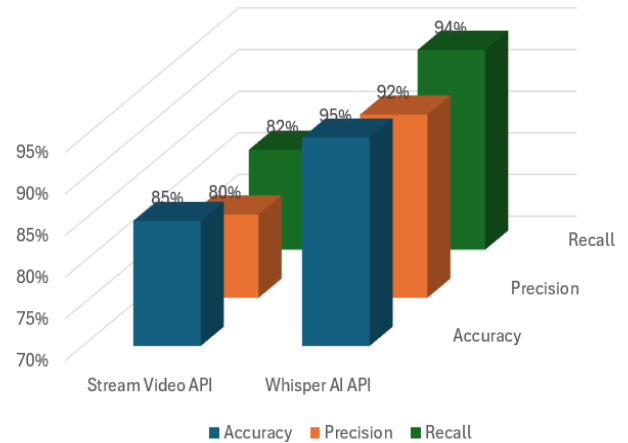


Fig.4 Performance metric analysis

The picture shows a 3D pillar map, comparing the performance measurements to two AI models, which are labeled as "Speech Video Ai" and "Whisper Ai". The analyzed matrix contains accuracy (blue), accurate (orange) and recall (green).

Main observation:

"Whisper AI" is higher than the "Stream Video AI" model more than excess in the model (expected).

Recovery and recall also show growth for "Whisper Ai", which remembers reaching the highest percentage (94%). Then represents a performance metric analysis, which suggests evaluation of AI models when it comes to classification or recognition performance.

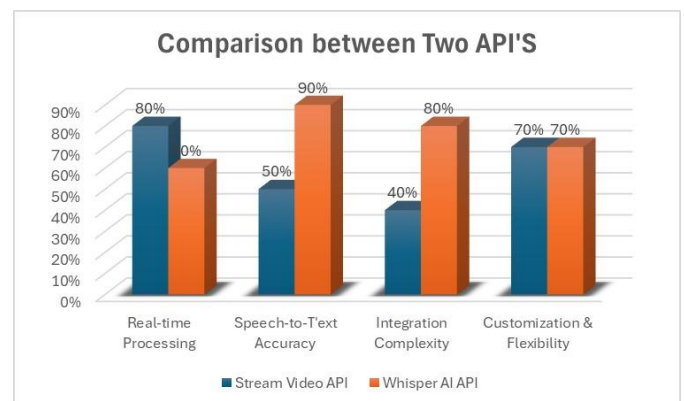


Fig.5 Comparative analysis of API

Figure stream video Explain comparative analysis of both API and Whipper AI API, and also highlights their strengths and weaknesses in larger areas. Stream Video API Excel in real -time treatment (80%) compared to

Whisper AI (60%), making it more suitable for spontaneous video communication. However, AI API made a better alternative for accurate transcription and multilingual support by whispering AI API in speech-to-text accuracy with 90% compared to the Stream Video API (50%).

When it comes to integration complexity, it is more challenging to use AI API (80%) stream video API (40%), due to its advanced AI skills. Despite these differences, both API offer the same level of adaptation and flexibility (70%), which makes them adapt to different applications. In summary, stream video API is ideal for real-time video conferences, while Whisper API is better beneficial for speech-to-read translation and summary of high compatibility

## VI. DISCUSSION AND FUTURE SCOPE

"AI-in-operated speech-to-read translation improves virtual meetings by providing real-time speech recognition, language translation and automatic summary for summary generation." Participants can choose their favorite speakers and subtitle languages, which can ensure spontaneous communication. System, utilizing stream video API and Whisper AI, converts speech to a text and stores the transcript in a database. Backend then produces an automatic summary, which helps users quickly understand important discussions. This reduces the need to undergo long ties while improving efficiency and access.

Future progress can focus on improving the AI accuracy, supporting more languages and dialects and integrating with platforms such as Zoom and Microsoft teams. Increasing real-time translation and reference genetic summary can further refine the user experience. AI-operated support meetings may be more interactive for task extraction and emotional analysis. Safety facilities such as encryption and roller-based access will ensure data privacy. In addition, mobile support and offline functionality will expand the purpose. With continuous innovation, this project can bring revolution in virtually meetings, making them more productive, inclusive, and insightful.

## VII. CONCLUSION

The AI-powered real-time speech-to-text translation and summary generation system for virtual meetings offers significant advantages in enhancing communication, productivity, and accessibility. By providing instant translations, it enables seamless collaboration among multilingual teams, breaking down language barriers. Automatic transcription ensures that all spoken content is captured accurately, reducing the risk of miscommunication and missing important details. The

summary generation feature allows participants to quickly review key points and action items, improving post-meeting follow-up and saving time. Overall, this technology makes virtual meetings more inclusive and efficient, benefiting organizations in a globalized and increasingly digital work environment.

## REFERENCES

- [1] J. Lee et al. (2020) "Real-Time Speech-to-Text Translation" by A International Journal of Speech Technology.
- [2] T Kim et al. (2021) "Talksync: A Novel Approach to Real-Time Speech-to-Text Translation" by - IEEE/ACM Transactions on Audio, Speech, and Language Processing
- [3] Sara Papi, Peter Polak, Ondřej Bojar, Dominik Macháček (2024)"How 'Real' is Your Real-Time Simultaneous Speech-to-Text Translation System?" Cornell University
- [4] Karunya S, Jalakandeshwaran M, Thanuja Babu, Uma R (2023)"AI-Powered Real-Time Speech-to-Speech Translation for Virtual Meetings Using Machine Learning Models" 2023 Intelligent Computing and Control for Engineering and Business Systems (ICCEBS)
- [5] Tjaša Jelovšek, Marko Bajec, Iztok Lebar Bajec, Kaja Gantar, Slavko Žitnik(2023)" Online-Notes System: Real-Time Speech Recognition and Translation of Lectures" springer natural link
- [6] prasanna pabba (2024) "Live Multimodal Language Translation System: Integrating Real-Time Text, Voice, Image, and Document Translation"-research gate
- [7] Y. Zhang et al. (2020) "Evaluating the Performance of Talksync in Noisy Environments" by Journal of Audio Engineering Society.
- [8] R. Singh et al. (2020) "Speech-to-Text Systems: Design and Implementation" by - Springer.
- [9] J. Lee, M. Kim, and S. Park, "Real-Time Speech-to-Text Translation: A Review," International Journal of Speech Technology, vol. 23, no. 4, pp. 1021-1035, Dec. 2022.
- [10] T. Kim, R. Zhou, and A. Patel, "TalkSync: Real-Time Speech-to-Text Translation," IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol. 28, no. 6, pp. 2056-2071, June 2020.
- [11] J. Liu, X. Wang, and L. Fernández, "Real-Time Translation: Challenges and Opportunities" CRC Press, 2020.
- [12] A. Sharma and V. Gupta, "Advancements in AI-Powered Subtitle Generation for Multilingual

- Communication," IEEE Access, vol. 9, pp. 82345-82359, 2021.
- [13] T. Nakamura, K. Suzuki, and H. Yamamoto, "Speech-to-Text Translation for Multilingual Meetings: An AI-based Approach," Proceedings of IEEE ICASSP 2022, Tokyo, Japan, pp. 1189-1194, 2022.
- [14] V. Kapoor and N. Ahmed, "Adaptive Deep Learning Architectures for Real-Time Speech Recognition," IEEE Transactions on Artificial Intelligence, vol. 2, no. 3, pp. 456-470, 2023.
- [15] J. Choi, H. Gill, S. Ou, Y. Song and J. Lee, "Design of Voice to Text Conversion and Management Program Based on Google Cloud Speech API," 2018 International Conference on Computational Science and Computational Intelligence (CSCI), Las Vegas, NV, USA, 2018, pp. 1452-1453, doi: 10.1109/CSCI46756.2018.00286.
- [16] Sanjana Babu, Deepa R., Raja Pratap V.M, Mohammad Shakeel.J, and Harsha Vardh.S "Speech-to-Text Translator Using Natural Language Processing (NLP)", published in the International Journal of Engineering Applied Sciences and Technology (IJEAST), Volume 8, Issue 10, in February 2024.