

HearNSign: Bidirectional Communication System for Signers and Non-Signers in Digital Meeting Environments

Bhavadharani V¹, Harithra R², Mercy Catherine D³, Mrs.K.Anandhi⁴

U.G. Student, Department of Information Technology, R.M.D Engineering College, Kavaraipettai, Tamil Nadu, India¹

U.G. Student, Department of Information Technology, R.M.D Engineering College, Kavaraipettai, Tamil Nadu, India²

U.G. Student, Department of Information Technology, R.M.D Engineering College, Kavaraipettai, Tamil Nadu, India³

Assistant Professor, Department of Information Technology, R.M.D Engineering College, Kavaraipettai,
Tamil Nadu, India⁴

ABSTRACT: Virtual meetings pose communication barriers to users who primarily rely on sign language. Current digital platforms have very limited functionality in supporting real time sign language interpretation in the form of an avatar. The assistive technologies available for the visually challenged have many limitations such as lack of accuracy, speed and do not support visually appealing representations required in virtual meetings. This paper thus presents an intelligent communication framework using Artificial Intelligence (AI) techniques for effective two-way communication between the sign language users and non-sign language users in virtual meetings.

The proposed framework has three layers namely Sign Recognition Module, Speech Recognition and Synthesis Module and Avatar Module. The Sign Recognition Module is based on the Transformer based gesture encoder to recognize the Indian Sign Language (ISL) from the video streams of users by analysing the hand movements, facial expressions and body language. The Speech Recognition and Synthesis Module comprises of RNN-T, CTC and DNN models for speech to text transcription of speech signal. The Avatar Module uses the neural avatar synthesis for generation of signs for the given speech or text.

KEYWORDS: Sign Language Recognition, Artificial Intelligence, Virtual Communication, Gesture Recognition, Neural Avatar Synthesis, Accessibility.

I. INTRODUCTION

Sign language is a mode of physical communication that is used by the deaf people. It is different from one country to the other. The sign language gestures and symbols are used in a linguistic way and for each individual sign is a word. Every sign has three aspects, namely: Handshape-Location of the hands-Motion of the hands. In India, American Sign Language (ASL) is used the most.

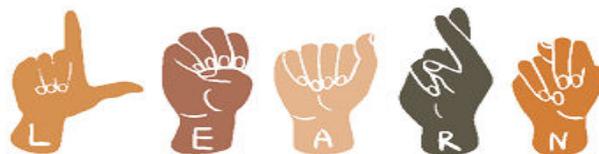


Fig 1.1 Sign Language

Sign language (also called finger spelling) is a language that uses hand and body movements, including facial expressions and lip movements. Sign language uses hand shape, orientation and movement as well as body language to convey meaning. Often a sign will convey the whole idea of a word, rather than the word itself. A common aspect of

sign language is finger spelling. Finger spelling is a system of hand positions to spell out individual letters in a word. A sign language is developed by a deaf community which is made up of Deaf/heard people, including their families, friends and interpreters.

II. PROBLEM STATEMENT

Virtual meeting apps are widely used in schools, workplaces, and social environments to connect people. Currently, these platforms are not accessible to users who depend on signing and this creates a barrier to communication between signers and non-signers. Deaf, mute or hard of hearing people are excluded from interactive participation in virtual meetings due to lack of real-time human or automatic sign language interpretation. At the same time, hearing people have difficulties understanding the sign language of the Deaf, mute or hard of hearing participants in the virtual meeting. The current assistive technologies including speech-to-text applications and gesture recognition applications offer only unidirectional communication, low accuracy, and delays. There is currently no real-time human or automatic sign language to speech or text to natural sign language video interpretation that addresses the needs of Virtual Meetings. An artificial intelligence-based solution is required for facilitating interactive communication in virtual meetings.

III. OBJECTIVE OF THE STUDY

Developing an AI-based communication technology for real-time virtual meeting spaces. Aim is to create a Sign Language (SL) recognition module for real time continuous Indian Sign Language (ISL) recognition from video streams using a Transformer-based gesture encoder, as well as a speech recognition system using state-of-the-art algorithms such as RNN-Transducer, Connectionist Temporal Classification (CTC) and Deep Neural Networks for speech to text transcription in real time. The proposed solution includes a neural avatar synthesis module for creating artificial sign language for given spoken or typed text. The overall objective is to provide real time communication and to enable deaf and mute, hard of hearing and non-signing people to access and participate in virtual meetings more effectively, using a single interface that has SL to text and text to SL functionality.

IV. LITERATURE SURVEY

In this section, we discuss existing research and technology that has been explored to facilitate Sign Language (SL) interpretation and communication. The technology available for SL recognition and real-time two-way communication in Virtual Meetings (VMs) is very limited. Currently, all accessible technologies can deal with isolated tasks such as gesture recognition or speech recognition. Therefore, none of them has the ability to act as a complete communication aid. Some of the existing research and technologies explored so far are SL recognition, Computer Vision based gesture recognition, speech recognition and assistive communication systems. All these are briefly discussed in the following subsections.

In paper [1], S. Gupta and R. Sharma proposed a deep learning-based real-time sign language recognition system for human-computer interaction. The authors implemented a vision-based model that detects and interprets hand gestures using deep learning techniques to enable communication between sign language users and computers. The system achieved improved recognition accuracy and faster gesture detection compared with traditional machine learning methods. However, the proposed system mainly focuses on gesture recognition and does not fully address real-time bidirectional communication between speech and sign language users.

In paper [2], L. Zhou, W. Li, and H. Li proposed a vision-based continuous sign language recognition system using transformer networks. The study utilized transformer-based deep learning architecture to recognize continuous sign language sequences from video input. The approach improved contextual understanding and sequence modeling compared to conventional CNN-based methods. Although the system demonstrated better performance in recognizing continuous gestures, it mainly focuses on recognition and does not integrate speech-to-sign or sign-to-speech communication features.

In paper [3], N. C. Camgoz et al. introduced Sign Language Transformers, a framework designed for joint end-to-end sign language recognition and translation. The authors used transformer architectures to directly translate sign language

videos into spoken language text. Their approach improved translation accuracy and reduced dependency on intermediate processing stages. However, the model requires large datasets and high computational resources, making real-time deployment challenging in practical communication systems.

In paper [4], H. Rastgoo, K. Kiani, and S. Escalera presented a comprehensive review of deep learning techniques used in sign language recognition systems. The paper analyzed different computer vision approaches, datasets, and deep learning models used in gesture recognition and sign language interpretation. It also discussed current challenges and future research directions in this field. However, the study mainly provides a survey of existing techniques rather than proposing a real-time implementation system.

V. SYSTEM ARCHITECTURE

The proposed system allows the hearing people and the sign language users to communicate real time. The system is consisting of three layers input layer, processing layer and output layer. The first layer contains the camera that transfers the images to the second layer. The second layer translates the images transferred from the first layer into sounds that are transferred to the third layer. The third layer converts the transferred sounds to speech output that allows real time communication between the hearing people and the sign language users. In addition the system uses cloud hosting technology.

In the input layer the user provides the data by means of live video in sign language, by speech or by writing a message. The camera recognizes the sign language of the user, the microphone records the speech of the user and the data is processed by noise reduction in order to be ready for the following stage.

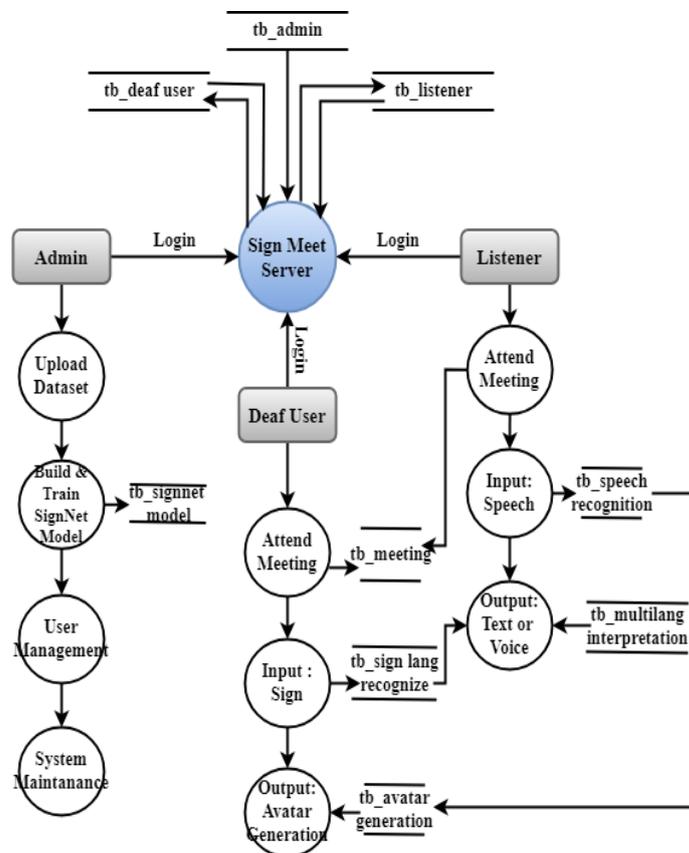


Fig 1.2 Data Flow Diagram

The processing layer comprises a number of AI units that instantly recognise input data. SignNet learns how to extract sign language gestures from videos and the Transformer-Based Gesture Encoder translates gesture sequences into text. The Speech Recognition Module recognises speech using various models such as RNN-Transducer (RNN-T), Connectionist Temporal Classification (CTC) and Deep Neural Networks (DNN) to output readable text. This output is then passed to the Neural Avatar Synthesis module that produces a lifelike video of the sign language being conveyed.

Sign Language To Speech recognition in Output layer – The sign language gestures performed by hearing user are converted into speech or text output for hearing. Conversely, the input speech or text of hearing user is converted into ASL (American Sign Language) animation Avatars for deaf user WebRTC (Web Real-Time Communication) is the underlying technology for real time communication between the hearing and deaf users. This module manages the whole backend admin panel functionality for user authentication and admin panel for the updations of model.

VI. OVERVIEW OF THE SYSTEM

The proposed system is a tool for an Artificial Intelligence (AI) communication application known as the AI Communication Platform (AICP). The AICP will enable deaf, mute and hearing individuals to communicate with each other in real time in a virtual environment. The system will be empowered with several features such as gesture recognition, speech to text, neural avatar and multi-language translation. The frontend of the proposed system is developed using the combinations of HTML, CSS, JavaScript, and WebRTC for real time video streaming and sign language avatars. The frontend communicates with the proposed system's backend using APIs. The proposed system's backend is built using Python and Flask framework for handling API routes, AI models and database connectivity.

The user details, gesture databases and session data are all stored in the MySQL or PostgreSQL databases. SignNet is capable of detecting continuous sign language movements in a real time video. The Transformer-Based Gesture Encoder leverages the temporal information within the gesture sequences for the accurate translation of gestures to speech. In addition, the speech recognition module of the proposed system utilizes several advanced speech recognition models such as Recurrent Neural Network Transducer (RNN-T), Connectionist Temporal Classification (CTC) and Deep Neural Networks (DNN). By utilizing the Neural Avatar Synthesis, the system can convert text and speech into realistic signs that can be easily understood by the deaf community. Hence, the various functionalities are interconnected by utilizing the Flask APIs and WebRTC for real-time communication, that too in a secured way. A comprehensive set of unit tests, integration tests and user acceptability tests have been carried out prior to deployment on a local or cloud server to validate the system for its effectiveness, performance and user acceptability.

V. MODULE DESCRIPTION

The proposed system consists of the following functional modules:

Sign Meet Web App:

The Sign Meet Web App is a common platform for the deaf and hearing to communicate each other in real time. This module has integrated all the AI and communication components into a common platform utilizing the backend made in Python, with web application built with Flask, and database in MySQL, and using Bootstrap for designing front end. It consists of live video streaming, real time gesture recognition, speech to text, and the avatar that translates and generates sign language. This module has the functionality of session management, routing, authentication and logging. This allows the deaf and hearing users to communicate naturally using video, text and animated sign language without the use of intermediaries or specialized hardware. The internet program enables the deaf and hearing users to communicate naturally using video, text and animated sign language, without the use of intermediaries or specialized hardware.

End User Interface:

Our system has an End User Interface (EUI) comprising of three functional modules - Administrator (Admin), Deaf User (DU) and Non-Deaf User (NDU). Admin functions include user administration, dataset upload, deploying machine learning models, system surveillance, reporting faults/bugs as well as enforcing security measures for the whole system. The DU interface caters for sign language interpretation via webcam-based gesture input. Deaf users can facilitate communication using their hand gestures that are being simultaneously interpreted into the written form on the computer screen in real time through machine learning techniques and animated chat bubbles which act as avatars.

Similarly, in the NDU interface, hearing users can key-in voice information or text, which is subsequently converted to sign language through the animation feature and the written text output into the computer screen.

SignNet Model: Build and Train:

The SignNet is a deep learning model used for motion-based sign language recognition from images and video streams. It uses the continuous video stream from the webcam to extract the hand movements, facial expressions and body language. The frames extracted from the video are processed by a series of techniques such as resizing, conversion to grayscale, removal of noise, thresholding and background subtraction. Deep learning techniques such as Convolutional Neural Networks (CNNs) with activation function ReLU and pooling are used to extract gesture and motion features from the collected frames. The collected features are classified using a fully connected layer and softmax layer. The SignNet model is trained on sign language datasets for classification and is used for real time gesture recognition.

Sign Language Recognition:

Sign Language Recognition module: It is able to recognize the movements present in the live video feed. Our device captures videos from webcam at a constant pace, showing the hand movements and facial expressions along with upper body movements, in real time. The important part is that it retains the transitions in between the gestures as well as the non-manual markers like eyebrows and head tilts that are quite crucial in passing on the ISL information. Frames are processed using Transformer-Based Gesture Encoder built upon the SignNet. It utilizes the self-attention to analyze the spatial features as well as the temporal movement patterns between frames and outputs a classification predicting the presence of a gesture along with the corresponding isolated word/phrase in Indian Sign Language

Multilanguage Interpretation:

This module is an accessibility tool designed to identify and translate common hand gestures of sign language into natural spoken languages using APIs. It can transliterate and translate the output text to languages such as English, Hindi and Spanish to facilitate communication and make it easier to understand among people speaking different languages.

Speech Recognition and Synthesis:

This module enables hearing people to hold conversations via speech input. Speech is captured via the microphone and recognised via the RNN-Transducer (RNN-T), Connectionist Temporal Classification (CTC) algorithms and Deep Neural Networks (DNNs) to translate speech to text accurately, thereby facilitating reliable real-time communication.

Avatar Generation:

This module utilizes neural avatar rendering. After speech-to-text or sign recognition technology translates speech into text, this module turns that text into the language of the Avataring system, producing natural hand movements and facial expressions associated with proper sign language interpretation. Deaf people are able to visually read information and hearing people are able to carry on natural and unimpeded conversations with Deaf people using this technology in real time.

Visual Communication:

This module integrates all the outputs generated by the system, such as gestures, translations, written language and body language of the avatar. Therefore, it provides a common platform that in the case of gesture recognition it gives the transcription of the gesture in written language for hearing users, and in the case of translations it provides written language and the body language of the avatar for both hearing and deaf users to have a two-way real time bidirectional interaction that synchronizes text and hand signs for both hearing and deaf users.

VI. METHODOLOGY AND ALGORITHMS

The Transformer-Based Gesture Encoder algorithm is a machine learning approach that takes sequential gesture data (e.g., sensor signals, video keypoints, or image patches) as input and learns to represent the data with a contextualized embedding, i.e., encoding, that is then classified using a fully connected layer (soft-max) similar to traditional machine learning approaches but leverages the self-attention mechanism from the well-known Transformer model.

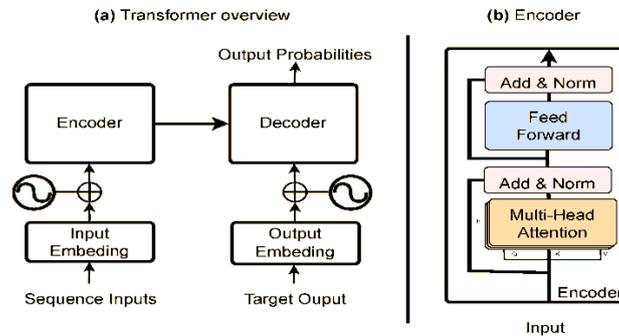


Fig 1.3 Transformer overview

A. RECURRENT NEURAL NETWORK TRANSDUCER (RNN-T)

Recurrent Neural Network Transducer (RNN-T) is an end-to-end, real-time streaming speech recognition model that can transcribe audio sequences to text sequences in real-time without needing any pre-processing on the audio data. RNN-T model utilizes the encoder, predictor (decoder) and the joint network to model the auditory and language interdependencies.

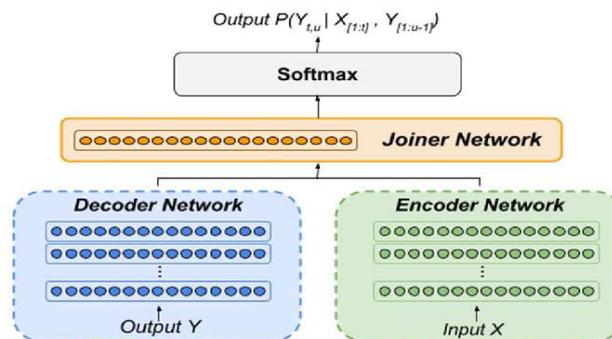


Fig 1.4 RNN-T

The RNN-Transducer (RNN-T) model is composed of three sub-networks, namely an acoustic encoder, a prediction network and a combined network for processing acoustic input and predicting a sequence of speech tokens. RNN-T models are popular choices for real-time speech recognition, due to their ability to recognize speech in real-time, without having to spend time waiting for the entire input, which in many cases can be highly variable. This model is capable of handling silence in speech due to the alignment of blank labels. As a result, RNN-T speech recognition models are increasingly used in speech recognition applications such as voice assistants and mobile dictation.

B. CTC

Connectionist Temporal Classification (CTC) is a way to train deep learning models to recognize sequences like speech or handwriting. By bypassing the need to understand the alignment between input and output sequences, the model can learn arbitrary relationships. CTC uses blank labels and a dynamic alignment between the input and output sequences, and therefore, the transcription model does not need to be manually trained on a large dataset, which is very useful in voice recognition where the training sets can be enormous.

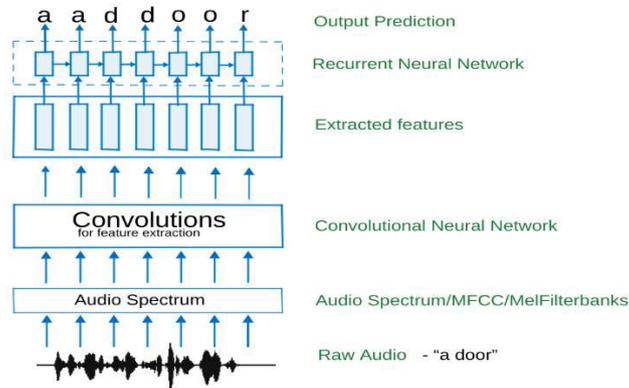


Fig 1.5 Connectionist Temporal Classification

Connectionist Temporal Classification (CTC) is a algorithm for dealing with the sequence learning problem with unequal input and output lengths. For an input sequence $X = [X_1, X_2, \dots, X_m]$ and an output sequence $Y = [Y_1, Y_2, \dots, Y_n]$, the CTC algorithm will automatically search for the optimal alignment during training. Thus the neural network can deal with variable input/output lengths directly and can be commonly used in speech recognition and hand writing recognition.

C. DEEP NEURAL NETWORK

A Deep Neural Network (DNN) is a type of artificial neural network model that consists of several hidden layers between the input and output layers, inspired by the human brain. The deep term in DNN refers to the network ability to learn several hierarchical and abstract levels of representations from a large amount of data.

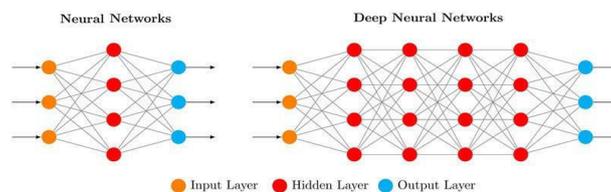


Fig 1.6 Deep Neural Network

Deep Neural Networks (DNNs) have an input layer that receives input data, multiple hidden layers that performs sophisticated processing of data using complex functions on connected weights and a output layer that gives the final prediction for a given input data. By adjusting its weights and biases to fit a given loss function, the model is said to have learned how to utilise the information contained within the data to identify and model highly complex patterns and relationships.

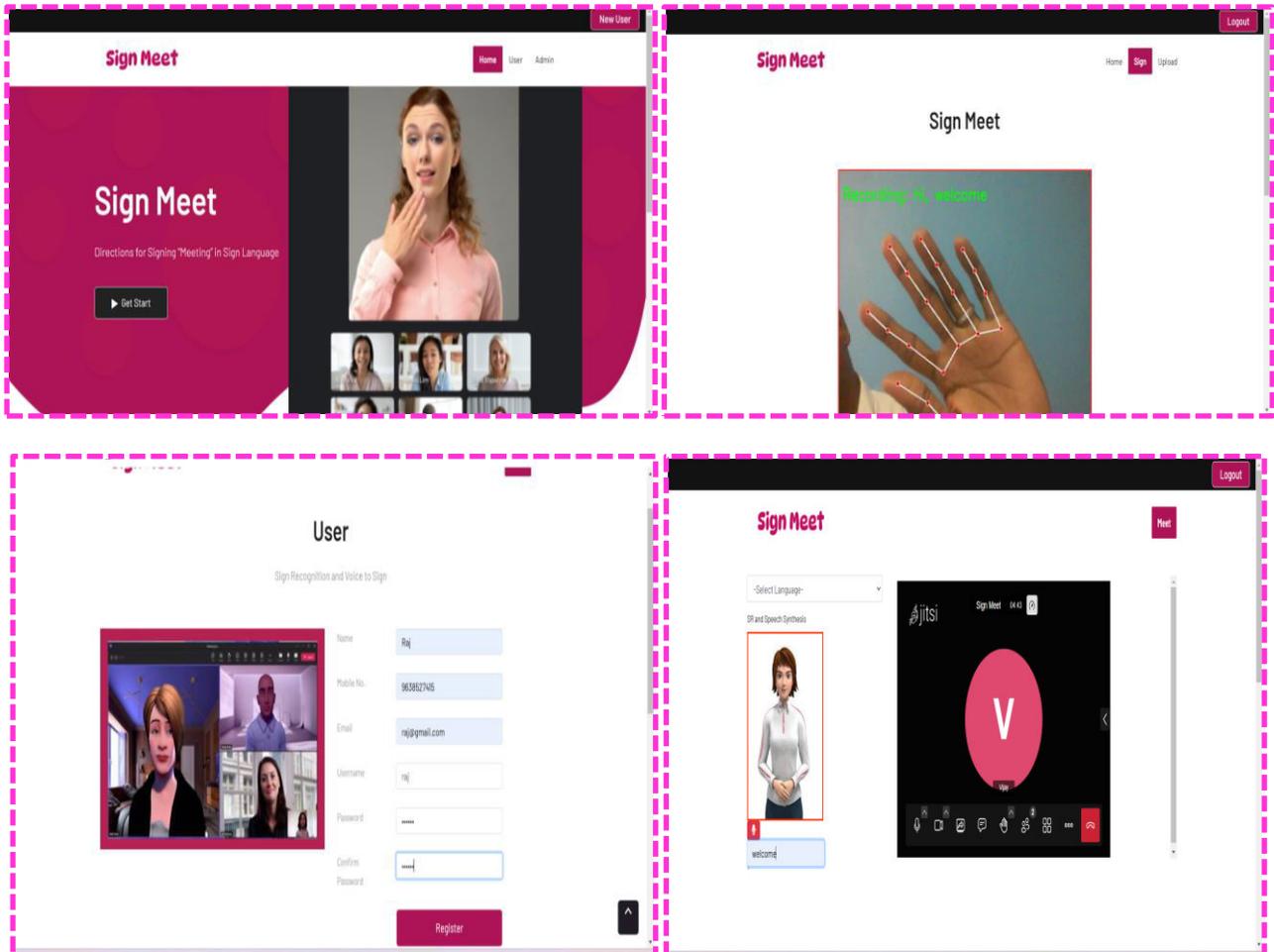
VII. RESULTS AND DISCUSSION

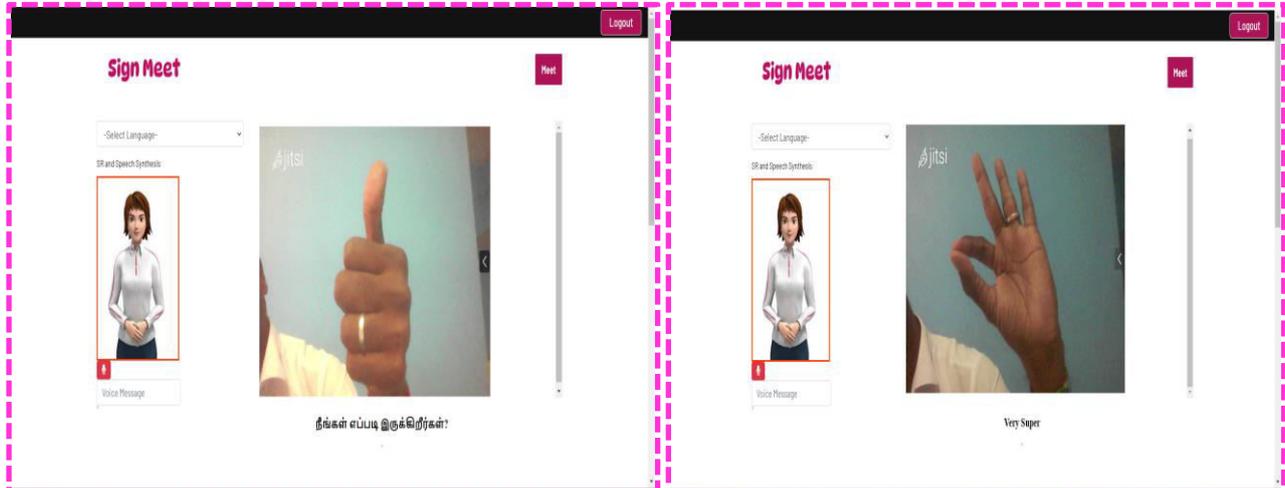
The proposed system is tested to validate the real time sign language recognition and speech generation. The SignNet is modelled and trained on the LSL dataset which includes the labeled motions of Indian Sign Language (ISL). The proposed system is tested using live video from camera. The experimental results validate that the SignNet model can achieve real time sign language recognition for dynamic sign language motions in various lighting conditions and camera distances. In addition, the Transformer-Based Gesture Encoder is employed to further enhance the accuracy of the model by learning temporal patterns in the gesture sequences to understand smooth transitions between signs.

Our Voice Recognition System was tested in several environments using a microphone input. A hybrid model using RNN-Transducer (RNN-T), Connectionist Temporal Classification (CTC) and Deep Neural Networks (DNN) achieved

a reasonable accuracy of speech-to-text in somewhat noisy speech. The output is then translated to animated sign language via the Neural Avatar Synthesis (NAS) module, enabling deaf people to visually capture spoken language.

The experimental results demonstrate that the proposed system supports bidirectional communication between deaf and hearing users via natural sign language. With the aid of the WebRTC-based real-time video streaming technology, the high-accuracy pre-trained models, the proposed system ensures low latency communication. Therefore, the proposed system has the potential of increasing communication accessibility in online platforms using an efficient and effective sign language communication technology.





A. COMPARISON WITH EXISTING SYSTEMS

The proposed solution overcomes some of the challenges and shortcomings of the current communication tools used by hearing and deaf people. The current assistive technology including cochlear implants and hearing aids do not facilitate the sign language communication. Instead, they focus on restoring the hearing ability of the deaf. Text messaging and written notes are considered as a form of manual alternative and require more efforts and may also break the flow of communication. Sign language interpreters do provide a more accurate interpretation of the sign language, but they are often not available. Moreover, it is difficult for them to interpret an impromptu or an online conference that involves a large number of participants.

Current Automatic Speech Recognition (ASR) systems deployed in video conferencing platforms can only display the spoken words as text and do not represent sign language visually. Most machine learning-based gesture recognition approaches such as SVM, k-NN, and HMM depend on pre-defined static gestures and cannot capture the dynamic nature of sign language.

In contrast, the proposed AI-driven system combines sign language recognition, speech-to-text conversion, and neural avatar generation on a single platform, allowing for real-time bidirectional communication between sign-language users and hearing participants in virtual meeting settings.

VIII. FUTURE SCOPE

The proposed method can be advanced with state-of-the-art technology in order to improve accessibility, performance and user experience. Some of the advancements planned include: - Increasing the language support to include more sign languages and dialects in order to reach out to a wider linguistic variety of users - Improving the customisation of the avatar in order to produce a more natural representation of sign language using advanced 3D animation techniques, such as for example the use of more natural gestures, facial expressions and lip syncing.

The system can be accessed from mobile devices, such as Smartphones and Tablets, to facilitate real-time communication. Meeting notes and transcripts can be produced using AI, in order to facilitate easy accessibility. Creating a mobile application will enhance user experience, increasing accessibility to all users, hence the proposed system will enable access to all stakeholders, enhancing communication within educational institutions.

IX. CONCLUSION

Current solutions fail to meet the diverse needs of deaf, mute, hard-of-hearing and non-signing individuals in virtual meetings. Existing solutions lack real-time or universal communication methods between signers and hearing individuals. Our proposed solution utilizes deep learning technologies to enable efficient real-time communication among deaf, mute, hard-of-hearing and non-signing users in virtual meetings, while at the same time enabling signers

and hearing individuals to communicate seamlessly with one another. This proposed solution is implemented by including sign language detection and speech recognition technologies with a multilingual model and a neural avatar-based visual communication technique. The proposed system uses a Transformer-based gesture encoder that can accurately recognize and classify continuous signs from any video stream. Real-time speech-to-text transcription is enabled by utilizing a mix of RNN-Transducer, CTC and Deep Neural Networks, which enables accurate transcription of spoken language from voice signals and produces real-time text output.

This neural avatar synthesis module will convert input text or speech to ASL(American Sign Language) using AVATOUR - AI-based technology. Our system is powered by a comprehensive suite of software libraries, built using widely known open-source languages and tools such as: Python Flask, WebRTC, MySQL, which enables quick interaction, in a low latency manner. Plans for feature improvements include supporting for multiple languages, more accurate tracking and animating of avatar movements, to enable participants who use ASL to communicate remotely in conferences. This may also be integrated with other multimedia conferencing solutions to increase multimodal accessibility for a more diverse communication set.

REFERENCES

- [1] S. Gupta and R. Sharma, "Deep learning based real-time sign language recognition system for human-computer interaction," *International Journal of Intelligent Systems*, vol. 38, no. 5, pp. 2750–2765, 2024.
- [2] L. Zhou, W. Li, and H. Li, "Vision-based continuous sign language recognition using transformer networks," *IEEE Access*, vol. 12, pp. 15210–15222, 2024.
- [3] N. C. Camgoz, O. Koller, S. Hadfield, and R. Bowden, "Sign language transformers: Joint end-to-end sign language recognition and translation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 1, pp. 105–120, 2023.
- [4] H. Rastgoo, K. Kiani, and S. Escalera, "Sign language recognition: A deep learning review," *IEEE Access*, vol. 11, pp. 13524–13545, 2023.
- [5] A. Moryossef, I. Tsochantaridis, and R. Reichart, "Data augmentation for sign language translation using neural networks," *Computer Vision and Image Understanding*, vol. 231, pp. 103663, 2023.
- [6] T. Reddy Gadekallu, G. Srivastava, and M. Liyanage, "Hand gesture recognition based on a Harris hawks optimized convolution neural network," *Computers & Electrical Engineering*, vol. 100, pp. 1–12, 2022.
- [7] G. Halvardsson, J. Peterson, C. Soto-Valero, and B. Baudry, "Interpretation of Swedish sign language using convolutional neural networks and transfer learning," *SN Computer Science*, vol. 2, no. 3, pp. 1–15, 2021.
- [8] P. Likhar, N. K. Bhagat, and R. G. N., "Deep learning methods for Indian sign language recognition," in *Proc. IEEE 10th Int. Conf. Consumer Electronics (ICCE-Berlin)*, Berlin, Germany, 2020, pp. 1–6.



INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 9940 572 462  6381 907 438  ijirset@gmail.com



www.ijirset.com

Scan to save the contact details