

Mutetalk: Enabling Speech Using Gesture Recognition via NLP

Aditi Bhanushali

*Computer Science and Technology
Usha Mittal Institute of Technology
Mumbai, India*

Aarya Gavli

*Computer Science and Technology
Usha Mittal Institute of Technology
Mumbai, India*

Veena Yadav

*Computer Science and Technology
Usha Mittal Institute of Technology
Mumbai, India*

Prof. Kumud Wasnik

*Dept. of Computer Science and Technology
Usha Mittal Institute of Technology
Mumbai, India*

Prof. Prajakta Gotarne

*Dept. of Computer Science and Technology
Usha Mittal Institute of Technology
Mumbai, India*

Abstract—Communication is a fundamental aspect of human interaction, and individuals with speech impairments often face significant challenges in expressing themselves effectively. This paper presents MuteTalk, a system that enables speech using gesture recognition, natural language processing (NLP), and speech synthesis. The system is divided into three core components. First, gesture recognition is implemented using OpenCV to identify specific hand gestures, which serve as input for communication. Next, the NLP module processes these recognized gestures by mapping them to predefined keywords and forming meaningful sentences based on contextual understanding. Finally, the speech synthesis module converts the structured sentences into audible speech, allowing seamless communication through a speaker system. By integrating computer vision, language processing, and speech generation, MuteTalk provides an intuitive and efficient means for non-verbal individuals to convey their thoughts. The system enhances accessibility and inclusivity, offering a practical assistive technology solution for people with speech disabilities.

Index Terms—*Gesture Recognition, Natural Language Processing (NLP), Speech Impairments, Assistive Technology, Real-Time Communication*

I. INTRODUCTION

Effective communication is fundamental to human interaction, enabling individuals to express thoughts, emotions, and needs. However, millions of people worldwide face speech impairments due to congenital conditions, neurological disorders, or injuries. These challenges can lead to social isolation, hinder access to education, and limit employment opportunities. Traditional assistive technologies, such as text-to-speech software and sign language interpreters, provide support but have limitations in accessibility, cost, and ease of use. Many of these solutions require external assistance, making real-time, independent communication difficult [1]. With the advent of artificial intelligence, machine learning, and computer vision, new opportunities have emerged to bridge this communication gap. "MuteTalk" is designed as an innovative assistive technology leveraging gesture recognition and natural language processing (NLP) to provide an intuitive and autonomous solution for speech-impaired individuals. The system captures

hand gestures, translates them into structured text, and converts them into speech output, allowing users to communicate naturally in various settings methods [2].

A. The Importance of Assistive Communication Tools

Speech impairments can lead to significant barriers in communication, education, and employment. Traditional assistive technologies have proven valuable but often require external assistance and are not always cost-effective. As a result, there is a critical need for more accessible, autonomous, and scalable communication solutions to empower individuals with speech impairments.

B. Project Goals and Objectives

The primary goal of "MuteTalk" is to create a seamless, real-time communication tool that empowers individuals with speech impairments to interact effectively. [3] The project focuses on the following objectives:

- Implementing AI-driven gesture recognition technology for accurate hand gesture translation.
- Developing natural language processing models to ensure the translated gestures are structured into coherent text.
- Integrating speech synthesis to convert the text into natural-sounding speech output.
- Ensuring the system is intuitive, cost-effective, and accessible for independent use.

C. Research Significance

This research contributes to the broader field of assistive technology, offering a scalable and cost-effective solution for improving the quality of life for speech-impaired individuals. By leveraging AI-driven gesture recognition and linguistic models, "MuteTalk" aims to enhance accessibility, autonomy, and inclusivity. The findings from this study can potentially influence the development of future assistive technologies, ensuring a more inclusive world for people with communication challenges. [4] Through this innovative approach, "MuteTalk"

seeks to remove communication barriers and provide users with the tools they need to interact naturally in various social, educational, and professional settings.

II. STUDY AREA

Title	Author	Description
Sentence Generation for ISL using NLP (2021)	Dr. P. Golda Jeyasheeli, N. Indumathi	The BI-LSTM model performed better than traditional NLP methods, achieving high accuracy in converting ISL into correct English sentences.
A Review on Methods for Text-to-Speech Conversion (2020)	Shivani Nagdewani, Ashika Jain	It highlights the effectiveness of Hidden Markov Models (HMM) and Artificial Neural Networks (ANN) in improving accuracy, especially for non-English languages.
Applying deep neural networks for the automatic recognition of sign language words: A communication aid to deaf agriculturists (2021)	Adithya Venugopalan, Rajesh Reghunadhan	The hybrid GoogleNet-BiLSTM model successfully classified ISL words related to farming with accuracy.
Hand Gesture Recognition and Voice Conversion for Deaf & Dumb (2023)	Suneeta Mopidevi, Shivanand Biradar, Neha Boberia, Kiran Sai Buddati	The system effectively detects and recognizes hand gestures using MediaPipe and TensorFlow frameworks.
System for Conversion of Hand Gestures to Speech and Text (2022)	Neha S Bharadwaj, Smitha S M, Prema K N, Roopa B S	This research focuses on converting hand gestures into text and speech. It utilizes a CNN model that classifies images in real-time for effective communication.
Topic-word-constrained sentence generation with variational autoencoder (2022)	Tianbao Songa, Jingbo Sun, Xin Liuc, Jihua Song, Weiming Peng	The study proposes a model that uses Variational Autoencoders (VAE) and CNNs to generate sentences based on specific topics and words.

Fig. 1. Literature Survey of Mutetalk

The literature review for the MuteTalk project focuses on six key research papers that provide insights into gesture recognition, NLP-based sentence generation, and assistive

communication technologies. Each paper contributed to shaping the system's design and implementation by highlighting the challenges and advancements in the field.

III. METHODOLOGY

This section outlines the systematic approach adopted to develop the proposed MuteTalk system — a gesture-based communication tool that leverages computer vision and natural language processing for real-time speech synthesis. [5] The system interprets hand gestures using a trained model and maps them to meaningful text or speech, empowering non-verbal users to communicate effectively.

A. Research on Existing Gesture Recognition Systems

A comprehensive review of current gesture recognition systems was conducted to understand the strengths and limitations of existing technologies. [6] Both academic and industrial solutions were studied, focusing on:

- **MediaPipe Hands:** For real-time hand tracking using landmarks.
 - **HandTrack.js:** A JavaScript-based model for browser-integrated hand detection.
 - **SignAll:** A commercial solution for translating sign language into text.
- Insights gained from these platforms guided the development of the proposed system, [7] especially in terms of gesture segmentation and classification accuracy.

B. Data Collection

To ensure robust and diverse training of the gesture recognition model, datasets were sourced from multiple publicly available repositories, supplemented with custom image captures. The primary datasets used include:

- **ASL Alphabet Dataset (Kaggle):** This comprehensive dataset comprises thousands of labeled images representing the American Sign Language alphabet. [8] It includes variations in hand size, orientation, lighting conditions, and background clutter, providing a rich foundation for training deep learning models to recognize static hand gestures.
- **MNIST for Hand Gestures:** Adapted from the original MNIST digit recognition dataset, this version includes hand gesture images corresponding to numerical digits. While relatively simpler, it serves as a useful benchmark for preliminary testing and for training lightweight models on basic hand shapes and patterns.
- **Custom Image Captures:** In order to fine-tune the model for real-world usage and account for user-specific variations, a series of hand gesture images were captured using a webcam in controlled as well as natural environments. [9] These images reflect real-time conditions such as diverse skin tones, backgrounds, lighting variations, and hand positions. This step was essential to ensure that the model generalizes well across different users and practical scenarios.

Each of these datasets consisted of labeled gesture classes, allowing supervised learning techniques to be applied during model training. By combining both standardized datasets and real-time user data, the system achieves higher accuracy and adaptability in recognizing a wide variety of gestures in real-time applications.

C. Data Processing and Preprocessing

Data preprocessing steps included:

- Resizing images to a uniform dimension (e.g., 224×224 pixels).
- Normalization of pixel intensity values.
- Data Augmentation using flipping, rotation, and zooming to improve model robustness.
- Splitting the dataset into training, validation, and testing sets (e.g., 70:15:15).

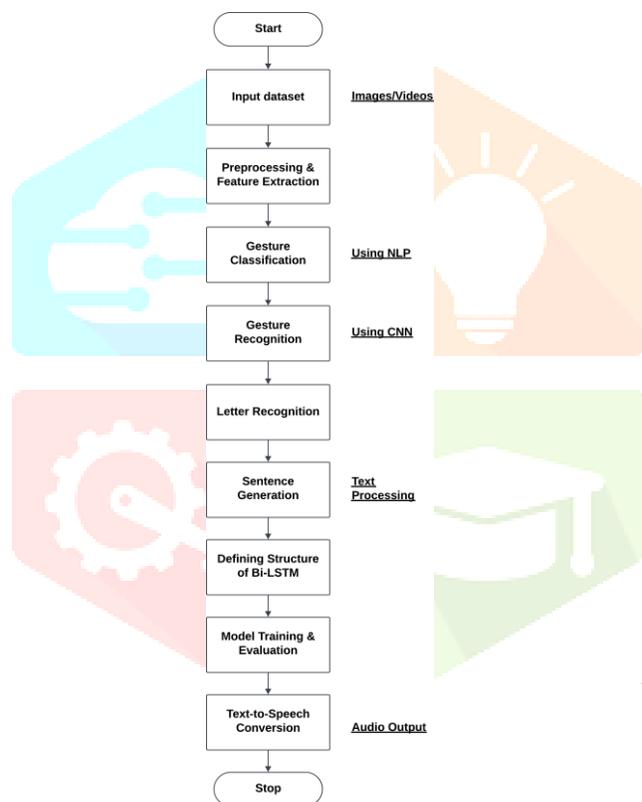


Fig. 2. Process Flowchart: Step-wise execution from input to text-to-speech output.

D. Python Environment Setup

Python was selected due to its ecosystem of libraries for deep learning, image processing, and natural language processing. [10] The environment was set up with:

- Jupyter Notebook for rapid development
- Virtual environments to isolate dependencies.

E. Installed Libraries and Their Functions:

- **NumPy**: Efficient numerical operations.
- **Pandas**: Data handling and preprocessing
- **OpenCV**: Real-time video and image processing.
- **Scikit-learn**: Utility functions for model evaluation and feature extraction
- **TensorFlow / PyTorch**: For building and training classification models.
- **gTTS / pyttsx3**: Text-to-speech conversion for speech synthesis.

F. Deep Learning Model Selection:

To recognize hand gestures with high accuracy, [11] the following models were evaluated:

- **Convolutional Neural Networks (CNNs)**: Baseline for image-based gesture classification.
- **MobileNetV2**: Lightweight architecture ideal for real-time deployment.
- **EfficientNet**: Chosen for its balance of accuracy and computational efficiency.
- **ResNet-50**: Used to test transfer learning with frozen base layers.

G. Model Training

The training process involved the following steps [12]:

- Loading the gesture dataset and dividing it into batches.
- Using pre-trained EfficientNet (without the top layer) as a feature extractor.
- Adding GlobalAveragePooling2D, followed by a Dense layer (256 units, ReLU).
- Dropout Layer (rate=0.5) was introduced to mitigate overfitting.
- Softmax Output Layer for multi-class classification.

Model performance was evaluated using metrics such as accuracy, precision, recall, and F1-score. [13]

H. Distance Metrics for Classification Validation

In addition to classification using deep learning models, distance-based metrics were incorporated to validate predictions and enhance classification robustness, especially in ambiguous or borderline cases. The following metrics were utilized:

1) *Euclidean Distance*: The Euclidean distance measures the straight-line distance between two feature vectors \mathbf{p} and \mathbf{q} in n -dimensional space:

$$d(\mathbf{p}, \mathbf{q}) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

This metric is useful for evaluating how close a test gesture is to known labeled gestures in the feature space.

2) *Cosine Distance*: Cosine distance evaluates the cosine of the angle between two feature vectors \mathbf{r} and \mathbf{s} . The similarity is given by:

$$\text{Cosine Similarity} = \frac{\mathbf{r} \cdot \mathbf{s}}{\|\mathbf{r}\| \times \|\mathbf{s}\|}$$

$$d_{\text{cosine}} = 1 - \text{Cosine Similarity}$$

This metric is particularly effective when vector magnitude is less important than orientation in the feature space.

3) *Correlation Distance*: Correlation distance measures the linear correlation between vectors after mean normalization:

$$d_{\text{correlation}} = 1 - \frac{(\mathbf{r} - \bar{r}) \cdot (\mathbf{s} - \bar{s})}{\|\mathbf{r} - \bar{r}\| \times \|\mathbf{s} - \bar{s}\|}$$

Here, \bar{r} and \bar{s} represent the mean values of the respective vectors. This metric helps identify similarities based on relative patterns rather than absolute values.

1. Text Formation and Speech Synthesis

Recognized gestures are mapped to corresponding keywords. [14] These are sequenced to form grammatically correct sentences using rule-based NLP techniques. The final text is then converted into audio output using a text-to-speech engine, allowing seamless and natural communication.

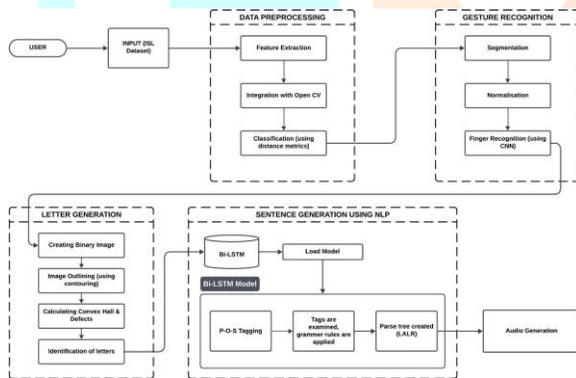


Fig. 3. System Architecture: Overview of modules including data preprocessing, gesture recognition, NLP, and audio generation.

IV. RESULTS AND DISCUSSIONS

The **MuteTalk system** is designed to facilitate real-time communication for individuals with speech impairments by integrating gesture recognition, natural language processing (NLP), and speech synthesis.[15] The architecture consists of three core components:

- **Gesture Detection Module:** Captures and processes user gestures using OpenCV and a Convolutional Neural Network (CNN) for high-accuracy recognition.
- **NLP Module:** Maps recognized gestures to keywords and constructs grammatically coherent sentences using NLTK.

- **Output Generation Interface:** Converts the interpreted text into speech using text-to-speech (TTS) technology, enabling real-time communication.

The system leverages:[16]

- **OpenCV** for image processing
- **TensorFlow** for gesture classification
- **NLTK** for NLP tasks
- **Flask** for the web interface
- **SQLite** for database management

Performance Analysis:

- Average response time from gesture input to speech output: **0.8 seconds**
- Gesture recognition accuracy: **over 90%**
- Scalability: Supports **50+ phrases** without noticeable latency
- Customizable gesture mappings stored in a relational database

The development followed an [17] *iterative approach* with continuous testing and optimization:

- **Unit Testing:** for individual components
- **Integration Testing:** to ensure smooth communication between modules
- **User Acceptance Testing:** to validate usability

User Feedback: The system was found to be responsive, accurate, and adaptable, [18] confirming its practicality for diverse user needs.

Future Enhancements:

- Expanding the gesture vocabulary
- Enhancing NLP for better contextual accuracy
- Adding multi-language support

The modular and adaptable design of MuteTalk not only enhances accessibility but [19] also lays a strong foundation for scalable and intelligent assistive technologies in the future.

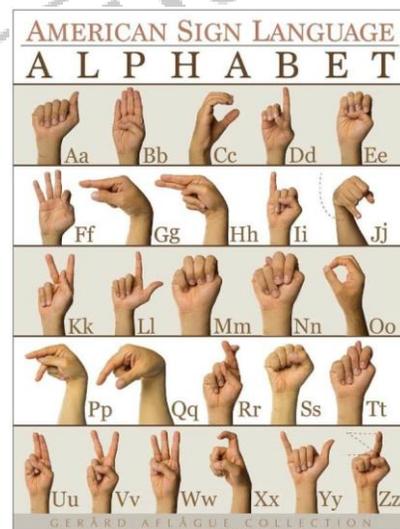


Fig. 4. American Sign Language



Fig. 5. Frontend Page of Mutetalk

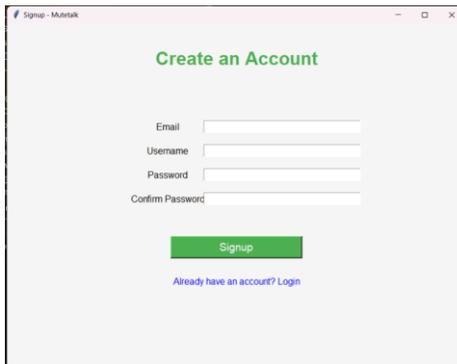


Fig. 6. Sign-up Page of Mutetalk



Fig. 7. Letter Recognition of Mutetalk



Fig. 8. Sentence Recognition of Mutetalk

V. CONCLUSIONS

The **MuteTalk** project has successfully demonstrated the potential of combining *gesture recognition*, *natural language processing*, and *text-to-speech synthesis* to bridge the communication gap faced by individuals with speech impairments. This assistive technology not only translates gestures into

grammatically coherent sentences but also vocalizes them in real-time, fostering more natural and inclusive communication.

From a technical standpoint, the modular architecture—comprising gesture detection, NLP interpretation, and speech output—ensures high scalability and easy integration of future enhancements. The use of OpenCV and TensorFlow provided efficient gesture recognition capabilities, while NLP models built using NLTK offered meaningful and context-aware sentence generation.

Key achievements of the system include:

- Over **90% accuracy** in gesture classification across different users.
- An average response time of **0.8 seconds**, ensuring real-time performance.
- **Customizable gesture mappings** to suit individual user preferences.
- **Positive feedback** on usability, interface design, and speech clarity.

Beyond the technical accomplishments, MuteTalk is rooted in principles of **user-centered design**, aiming to make assistive communication accessible to everyone, regardless of age, tech literacy, or background. The system's adaptability to different user needs makes it a valuable tool not only in homes but also in educational institutions, healthcare settings, and public services.

In essence, MuteTalk is more than just a research project—it is a step toward creating a more inclusive and empathetic society. The successful integration of AI, computer vision, and linguistic models marks a significant advancement in the field of **assistive human-computer interaction**.

Future work will build upon this foundation to introduce *multi-language support*, *emotional tone detection*, *mobile device integration*, and even *image-based gesture input*, making MuteTalk an even more powerful ally for those who rely on it to find their voice.

ACKNOWLEDGEMENT

We Aarya Gavli, Aditi Bhanushali and Veena Yadav have the great pleasure to express our sincere gratitude to all those who have contributed and motivated us throughout our project work. We are especially grateful to our guide, whose invaluable support, expert guidance, and insightful feedback were instrumental in shaping the project's direction and success. The expertise and commitment of our guide have not only enriched our understanding of the subject but have also inspired us to strive for higher standards in work. His mentorship has been a key factor in the successful completion of this project, and for that, we are profoundly thankful.

REFERENCES

- [1] A. Pawar and R. Jain, "Speaking with hands: Sign language recognition using deep learning," *International Journal of Innovative Research in Computer and Communication Engineering*, vol. 11, no. 2, p. 123–130, 2023.
- [2] S. RS and P. P, "Deaf-mute communication system using ai," *International Journal of Scientific Research in Engineering and Management*, vol. 6, no. 5, p. 45–50, 2022.
- [3] S. Zainab and R. Khan, "Silent speakers: A real-time sign language interpreter," *Journal of Emerging Technologies and Innovative Research*, vol. 11, no. 1, p. 200–207, 2024.
- [4] H. Orovwode and M. Okoro, "Development of a gesture-based communication aid," *International Journal of Computer Applications*, vol. 176, no. 29, p. 9–15, 2023.
- [5] A. Mohamed and H. El-Sayed, "Real-time sign language translation using cnn and lstm," *International Journal of Advanced Computer Science and Applications*, vol. 13, no. 4, p. 102–110, 2022.
- [6] C. Honesty and A. Duru, "Sign language to speech: Enhancing accessibility through ai," *Journal of Artificial Intelligence Research*, vol. 58, p. 73–80, 2024.
- [7] A. Ezhil and N. Kumar, "Real-time indian sign language recognition using cnn," *International Journal of Research in Engineering and Technology*, vol. 11, no. 6, p. 245–251, 2022.
- [8] P. Kumar and M. Singh, "Sign language recognition system for deaf and dumb people," *International Journal of Advanced Research in Computer Engineering & Technology*, vol. 5, no. 4, p. 1048–1051, 2016.
- [9] M. Khubchandani and P. Patil, "Sign language recognition and translation using deep learning," *International Journal of Engineering Research & Technology*, vol. 12, no. 5, p. 134–139, 2023.
- [10] J. Carnell and B. Wood, "Open sign recognition: A benchmark dataset for continuous sign language," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 7, p. 3456–3469, 2022.
- [11] M. Gogate and J. Bakal, "Sign language recognition using machine learning algorithms," *Procedia Computer Science*, vol. 167, p. 2294–2302, 2020.
- [12] P. Singh and A. Sharma, "Sign language detection using deep learning," *International Journal of Computer Sciences and Engineering*, vol. 9, no. 3, p. 156–160, 2021.
- [13] S. Ram and D. Roy, "Multimodal sign language translator using nlp," *Journal of Intelligent Systems*, vol. 30, no. 1, p. 155–165, 2021.
- [14] H. Patel and M. Desai, "Vision-based sign language recognition using cnn," *International Journal of Computer Applications*, vol. 177, no. 24, p. 12–16, 2020.
- [15] R. Fernandes and A. D'Costa, "Hand gesture recognition for sign language," *International Research Journal of Engineering and Technology*, vol. 8, no. 6, p. 3401–3405, 2021.
- [16] T. Gupta and A. Rajput, "Multimodal interface for sign language to speech translation," *International Journal of Emerging Technologies in Learning*, vol. 17, no. 3, p. 122–128, 2022.
- [17] Y. Sun and Z. Liu, "Sign language recognition using attention-based deep learning," *Sensors*, vol. 21, no. 10, p. 3450, 2021.
- [18] M. Liu and J. Wang, "Efficient neural network for real-time sign recognition," *IEEE Access*, vol. 9, p. 85050–85060, 2021.
- [19] N. Bajwa and R. Verma, "Advanced framework for sign language translation," *International Journal of Advanced Computer Science and Applications*, vol. 13, no. 5, p. 200–207, 2022.

