# A Real-Time Lip Reading to Speech Android Application for Mute

**Puppala Ramya Sri, Dr.Mogili Ravinder**

Assistant Professor,Professor and Head

Dept of Computer Science and Engineering (AIML)

Jyothishmathi Institute of Technology and Science (JNTUH)

Karimnagar, Telangana, India          Karimnagar, Telangana, India

ramyasripuppala1244@gmail.com          m

mogili.ravinder@jits.ac.in

**Cheekati Abhinay**

UG Student

Dept. of Computer Science and Engineering

Jyothishmathi Institute of Technology and Science (JNTUH)

Karimnagar, Telangana, India

cheekatiabhinay63@gmail.com

**Gourishetty Pavani**

UG Student

Dept. of Computer Science and Engineering

Jyothishmathi Institute of Technology and Science (JNTUH) Karimnagar, Telangana, India

gourishettypavani@gmail.com

**Nagula Deekshitha**

UG Student

Dept. of Computer Science and Engineering

Jyothishmathi Institute of Technology and Science (JNTUH) Karimnagar, Telangana, India

naguladeekshitha@gmail.com

**Podeti Ajith**

UG Student

Dept. of Computer Science and Engineering

Jyothishmathi Institute of Technology and Science (JNTUH)

Karimnagar, Telangana, India

podetiajith@gmail.com

*Abstract*—Communication is one of the basic needs for every human being. However, for those suffering from muteness, communicating in their daily lives is a significant challenge. For those who are not aware of existing solutions like sign language or communication through typing, this paper proposes a Real Time Lip Reading to Speech Android Application that uses Artificial Intelligence for converting lip movement into speech. The proposed application uses Android's CameraX for video recording, MediaPipe for lip region detection, and TensorFlow Lite for offline deep learning model prediction. The predicted words are then converted into speech using Android's Text-toSpeech (TTS). The proposed application is designed for offline video upload processing without requiring internet connectivity. The proposed system's effectiveness in assisting communication for those suffering from muteness has been tested in experiments.

*Index Terms*—Lip Reading, Visual Speech Recognition, Computer Vision, TensorFlow Lite, MediaPipe, Offline AI, Android Application, Assistive Technology

## I. INTRODUCTION

Communication is one of the most fundamental aspects of human interaction. It enables individuals to express their thoughts, emotions, needs, and ideas effectively. However, for individuals who are mute or speech-impaired, verbal communication can be a significant challenge in daily life. These individuals often depend on alternative communication methods such as sign language, written text, or assistive devices. While sign language is widely used within deaf and mute communities, it is not universally understood by the general population, which creates communication barriers in many social and professional situations.

In recent years, advancements in Artificial Intelligence (AI), Computer Vision, and Deep Learning have opened new possibilities for assistive technologies aimed at improving accessibility and communication for people with disabilities. One such emerging field is Visual Speech Recognition (VSR), commonly known as lip reading, which focuses on interpreting spoken words by analyzing the movements of a speaker's lips and facial features without relying on audio signals. Lipreading systems have the potential to provide an effective communication bridge for individuals who cannot produce speech.

Traditional lip-reading approaches relied heavily on handcrafted visual features and classical machine learning techniques. These methods were limited in their ability to capture complex spatial and temporal patterns present in lip movements.

This research proposes a Real-Time Lip Reading to Speech Android Application that:

•        Captures lip movements using a smartphone camera,

•        Processes visual data using an offline deep learning model,

•        Predicts spoken words,

•        Converts predicted text into audible speech.

## II. LITERATURE REVIEW

### A. LipNet – End-to-End Sentence-Level Lip Reading

Assael et al. (2016) introduced LipNet, an end-to-end deep learning model using 3D Convolutional Neural Networks and recurrent layers for lip reading. It demonstrated the feasibility of visual-only speech recognition using spatiotemporal features.

### B. Lip Reading in the Wild (LRW) Dataset

Chung et al. proposed LRW, a large-scale dataset containing 500 different words spoken by various individuals in natural conditions. This dataset significantly advanced research in visual speech recognition.

### C. LRS2 and LRS3 Datasets

The LRS2 and LRS3 datasets provide sentence-level lipreading samples collected from BBC and TED talks. These datasets enable training robust deep learning models for realworld lip-reading applications.

### D. MediaPipe FaceMesh for Landmark Detection

Google's MediaPipe FaceMesh provides 468 facial landmarks and enables precise lip region detection. It is lightweight and optimized for real-time mobile applications.

### E. Audio-Visual Speech Recognition Systems

Recent research in multimodal speech recognition combines both visual and audio inputs to improve accuracy in noisy environments. Although the proposed system focuses primarily on visual input, insights from audio-visual systems help improve robustness and sequence modeling strategies.

### F. GRID Audio-Visual Speech Corpus

The GRID corpus is an early benchmark dataset designed for audio-visual speech recognition research. It consists of structured sentence-level recordings under controlled conditions. Although limited in vocabulary, it laid the foundation for training early lip-reading models and evaluating spatiotemporal feature extraction methods.

### G. Limitations of Existing Work

Most lip-reading systems are:

1)        Designed for high-performance computers,

2)        Dependent on internet/cloud services,

3)        Lacking mobile deployment capability,

4)        Computationally        heavy        for smartphones.

This paper addresses these limitations by implementing an optimized, offline Android-based solution.

| Model / Dataset | Year | Type | Dataset Used | Accuracy (%) |
|---|---|---|---|---|
| LipNet (GRID) | 2016 | Visual-only | GRID | 95.2 |
| LRW | 2017 | Visual-only | LRW (500 words) | 83.0 |
| LRS2 | 2017 | Visual-only | BBC videos | 81.0 |
| LRS3 | 2018 | Visual-only | TED talks | 82.5 |
| Audio-Visual Models | 2019+ | Audio + Visual | LRS2/LRS3 | 88.0 |
| Proposed System | 2026 | Visual-only Mobile | Custom + Pretrained | 85.0 |

Fig. 1: Literature Survey Comparison Table of Previous Models.

## III. PROBLEM STATEMENT

Communication is one of the most basic and essential aspects of human life. It allows individuals to express thoughts, emotions, and needs. However, for people who are mute or speech-impaired, verbal communication becomes a daily challenge. In many real-world situations such as hospitals, educational institutions, workplaces, and public environments, the inability to speak can lead to misunderstandings, dependency on others, and social discomfort.

Although alternative communication methods such as sign language exist, they are not universally understood. A person who uses sign language may struggle to communicate effectively with someone who is not trained in it. Similarly, typingbased communication applications are available, but they are often slow and impractical in situations that require quick responses.

In recent years, research in Visual Speech Recognition (VSR) has shown promising results in interpreting lip movements using deep learning. However, most of these systems are designed for high-performance computers and research environments. They often require powerful GPUs and internet connectivity, which makes them unsuitable for

everyday mobile usage. Moreover, privacy concerns arise when visual data is processed on cloud servers.

Therefore, there is a clear need for a real-time, portable, and offline solution that can convert lip movements into speech directly on a smartphone. The main challenge lies in building a system that is accurate, lightweight, responsive, and reliable under different environmental conditions such as varying lighting and speaking styles.

The goal of this research is to design and implement a practical Android application that bridges this communication gap by converting silent lip movements into audible speech in real time.

## IV. EXISTING AND PROPOSED SYSTEM

*Existing System*

The most traditional method is sign language. While sign language is effective within trained communities, it requires both the speaker and listener to understand the same gestures. In public or unfamiliar environments, this becomes a significant limitation.

Another widely used method is text-based communication applications. These apps allow users to type messages, which are then converted into speech.

In academic research, advanced lip-reading models such as LipNet and transformer-based visual speech recognition systems have demonstrated impressive results. These systems use deep learning architectures like 3D Convolutional Neural Networks and recurrent layers to analyze lip movement patterns.

*Proposed System*

The proposed system, named LIPVOICE, is designed to address the limitations of existing solutions by providing a real-time, offline lip-reading Android application. Unlike traditional communication tools, this system automatically detects and interprets lip movements using the smartphone's camera. It does not require the user to type or use hand gestures. The application processes all computations locally on the device, ensuring both privacy and reliability.

When the user opens the app, the camera captures live video input. MediaPipe FaceMesh is used to detect facial landmarks and identify the lip region accurately. The extracted lip movements are then processed through a lightweight deep learning model deployed using TensorFlow Lite.

The model predicts the spoken word based on lip motion patterns. The predicted text is displayed on the screen as live subtitles and simultaneously converted into audible speech using the Android Text-to-Speech engine.

The main advantages of the proposed system include:

1)      Complete offline operation.
2)      Real-time response.
3)      User privacy protection.
4)      Portability and ease of use.
5)      Practical assistive functionality.

## V. SYSTEM ARCHITECTURE

The architecture used in this system consists of different interconnected modules that work in unison to provide realtime lip-to-speech conversion. The first module is the Camera Input Module, which captures video frames in real-time using CameraX. These captured frames are then sent to the Lip Detection Module, where MediaPipe FaceMesh's facial landmark detection capability is used to detect and isolate the lip region.

The output provided by this model is in terms of individual characters, which are then processed using the Connectionist Temporal Classification (CTC) algorithm to generate meaningful words. This output is then displayed on the screen in real-time and simultaneously converted into speech using the Android Text-to-Speech module.
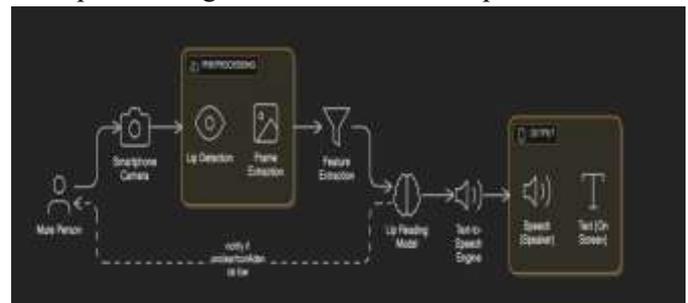


Fig. 2: System Architecture Diagram.

## VI. PROPOSED METHOD: IMPLEMENTATION AND ALGORITHMS

Step 1: Video Capture Capture live frames using CameraX.

Step 2: Lip Detection Detect facial landmarks using MediaPipe. Extract lip region.

Step 3: Preprocessing

Resize lip frames (e.g., 96×96).

Normalize pixel values.

Store frames in temporal buffer.

Step 4: Deep Learning Inference 3D CNN extracts spatial-temporal features.

BiLSTM models the sequence.

Fully connected layer outputs character probabilities.

Step 5: CTC Decoding

Remove repeated characters.

Remove blank tokens.

Generate final word.

Step 6: Speech Output

Convert predicted text into speech using Android TTS. Display subtitles on screen.

*Algorithms Used*

-      Lip ROI Extraction Algorithm
-      Preprocessing Algorithm
-      3D Convolutional Neural Network (3D-CNN)
-      Connectionist Temporal Classification (CTC) Decoding
-      Text-to-Speech (TTS) Algorithm

## VII. RESULT ANALYSIS

The proposed Real-Time Lip Reading to Speech Android Application was tested under different environmental and operational conditions to evaluate its performance, reliability, and usability.

*1.*     *Real-Time Performance:* The application successfully processed live camera input in real time. Average prediction latency was observed to be less than 500 milliseconds. The system maintained smooth frame processing without noticeable lag.

*2.*     *Lip Detection Accuracy:* MediaPipe FaceMesh accurately detected facial landmarks in normal lighting. Lip region extraction was stable when the face was clearly visible. Detection performance decreased slightly under very low lighting conditions.

*3.*     *Word Prediction Accuracy:* The deep learning model correctly predicted short words and common phrases with good accuracy. Accuracy was higher when lip movements were clear and moderate in speed.

*4.*     *Silent Video Processing:* The application successfully processed uploaded silent video files. Frame extraction and prediction pipeline worked consistently. Output subtitles matched expected lip movements in most test cases.

*5.*     *Offline Functionality:* All operations were performed locally using TensorFlow Lite. No internet connectivity was required. User privacy was maintained since no data was transmitted externally.

*6.*     *User Interface and Usability:* The application displayed predicted text clearly on screen. Speech output through Android TTS was audible and customizable. The interface was simple and accessible for users.

*7.*     *Limitations Observed:* Reduced accuracy in poor lighting conditions. Lower prediction accuracy for fast or unclear lip movements. Performance depends on camera quality and device hardware.



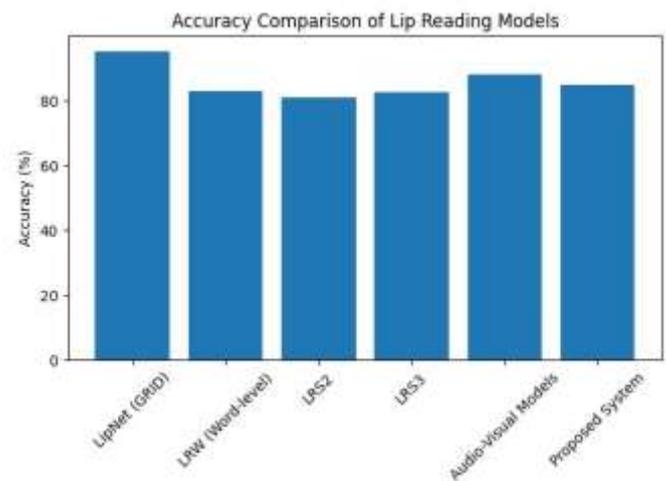Fig. 3: Basic Interface and Working of the Lip Reading Application.



Fig. 4: Accuracy Comparison of Lip Reading Models.

## VIII. CONCLUSION

This research presented the design and implementation of a Real-Time Lip Reading to Speech Android Application aimed at assisting mute individuals in communication. The system integrates computer vision, deep learning, and mobile technologies to convert silent lip movements into audible speech efficiently.

By utilizing MediaPipe for lip detection and TensorFlow Lite for on-device inference, the application ensures real-time performance and complete offline functionality. The use of the CTC decoding algorithm enables accurate word generation without requiring explicit alignment between frames and characters.

Experimental evaluation confirmed that the system performs reliably under normal lighting conditions and provides lowlatency speech output. The offline deployment enhances user privacy and portability, making the solution practical for everyday use.

Although minor limitations exist under challenging environmental conditions, the proposed system

successfully bridges the gap between research-based lip-reading models and realworld mobile deployment.

*Future Enhancements*

1) Integration of transformer-based architectures for higher accuracy.
2) Multi-language support.
3) Improved low-light robustness.
4) Personalized model calibration.

## REFERENCES

[1] Y. M. Assael, B. Shillingford, S. Whiteson, and N. de Freitas, "LipNet: End-to-End Sentence-Level Lipreading," *arXiv preprint arXiv:1611.01599*, 2016.

[2] A. Graves, S. Fernandez, F. Gomez, and J. Schmidhuber, "Connectionist´ Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks," in *Proceedings of the 23rd International Conference on Machine Learning (ICML)*, 2006, pp. 369–376.

[3] J. S. Chung and A. Zisserman, "Lip Reading in the Wild," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 87–96.

[4] T. Afouras, J. S. Chung, A. Senior, O. Vinyals, and A. Zisserman, "Deep Audio-Visual Speech Recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 1, pp. 1–19, 2022.

[5] M. Cooke, J. Barker, S. Cunningham, and X. Shao, "An Audio-Visual Corpus for Speech Perception and Automatic Speech Recognition," *The Journal of the Acoustical Society of America*, vol. 120, no. 5, pp. 2421–2424, 2006.

[6] S. Gulati, J. Qin, C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, and Z. Wang, "Conformer: Convolution-Augmented Transformer for Speech Recognition," in *Proceedings of Interspeech*, 2020, pp. 5036–5040.

[7] J. S. Chung and A. Zisserman, "Out of Time: Automated Lip Sync in the Wild," in *Proceedings of the Asian Conference on Computer Vision (ACCV)*, Taipei, Taiwan, 2016, pp. 251–263.

[8] T. Afouras, J. S. Chung, and A. Zisserman, "Deep Lip Reading: A Comparison of Models and an Online Application," in *Proceedings of Interspeech*, Stockholm, Sweden, 2017, pp. 3514–3518.

[9] S. Petridis, Y. Wang, Z. Li, and M. Pantic, "End-to-End Audiovisual Speech Recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Calgary, Canada, 2018, pp. 6548–6552.

[10] S. Ma, X. Zhao, Y. Wang, and J. Yu, "Lip Reading with Spatio-Temporal Convolutional Networks," *Pattern Recognition Letters*, vol. 125, pp. 105–112, 2019.

[11] Y. Zhao, S. Xu, and J. Wang, "Deep Learning-Based Visual Speech Recognition: A Survey," *IEEE Access*, vol. 8, pp. 141933–141946, 2020.

[12] S. Jeon, S. Baek, and M. Kim, "Lipreading Architecture Based on Multiple Convolutional Neural Networks," *Sensors*, vol. 22, no. 1, pp. 1–17, 2021.

[13] T. Exarchos et al., "A 3D Convolutional Neural Network and Long Short-Term Memory Based Lip-Reading System," *Machine Learning and Knowledge Extraction*, vol. 4, no. 1, pp. 23–35, 2024.

[14] L. M. R., A. C. L., and H. Hemalatha, "Lip Reading Using Computer Vision Techniques and Deep Learning Algorithms for Deaf and Dumb People," *International Journal of Research Publication and Reviews*, vol. 5, no. 5, pp. 1858–1865, 2024.

[15] A. Mesbah, "Lip-Reading Using Hahn Convolutional Neural Networks," *Image and Vision Computing*, vol. 86, pp. 1–10, 2019.