

AI-Powered Silent Speech to Text Transformation

Ms. Pooja K, Mr. Pradeep T, Ms. Shanmugapriyaa V R

Mr. Madhan K S P [Mentor]

ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING

SRI SHAKTHI INSTITUTE OF ENGINEERING AND TECHNOLOGY , COIMBATORE.

Abstract - This project introduces a silent speech-to-text system driven by artificial intelligence that uses a web-based interface to translate lip movements into readable text. By combining two crucial preprocessing modules—Frame Blur Removal and Lip ROI Stabilization—the system improves the performance of an already-existing lip-reading model. By identifying and fixing motion-blurred frames, the blur removal module enhances clarity and makes lip contours easier to see. Using facial landmark tracking, the lip ROI stabilization module reduces positional variations brought on by user movement by maintaining a uniform and centered mouth region throughout frames. These improvements give the model cleaner, more reliable input, which increases prediction accuracy. By taking frames from a webcam, processing them on the backend, and displaying the predicted text on the frontend interface, the suggested system runs in real time. According to experimental findings, the additional modules greatly improve the model's resilience to real-world scenarios like head movement, lighting fluctuations, and camera instability. In addition to supporting future development toward assistive technologies and real-time speechless interaction, the solution shows a workable strategy for silent communication systems.

Key Words: Lip Reading, Silent Speech Recognition, Computer Vision, Deep Learning, Dataset Preprocessing.

1.INTRODUCTION

The goal of the developing field of artificial intelligence known as "silent speech recognition" is to comprehend speech without the use of audio. Rather, it predicts the spoken text by analyzing a person's lip movements. This

method is helpful in circumstances where sound cannot be accurately recorded, such as in noisy settings, private conversations, or people with speech or hearing impairments. However, real-world problems like blurry frames, unstable lip regions, and head movements frequently cause lip-reading systems to lose accuracy. Our project offers an AI-powered silent speech-to-text system with two significant improvements to address these issues: a Lip ROI Stabilization module to preserve a consistent mouth region across frames and a Frame Blur Removal module to enhance visual clarity. When the model is used via a real-time web interface, these enhancements enable it to generate predictions that are more accurate.

2. Body of Paper

A. Overview of the System

The suggested system is a silent speech-to-text platform driven by AI that uses a web interface to translate lip movements into legible text. The user's webcam records video frames, which the system then processes through a number of enhancement and prediction stages in real time. In order to ensure more accurate text output during real-world use, the system's primary objective is to enhance the clarity and stability of lip-region frames prior to their delivery to the lip-reading model.

There are three main parts to the system:

Web interface or frontend

The user launches the website and grants access to the camera. Video frames are continuously captured by the frontend and sent to the server.

Engine for Preprocessing

To enhance video quality and preserve a consistent lip region across frames, the incoming frames are processed through two new modules: Frame Blur Removal and Lip ROI Stabilization.

Model Backend

Following preprocessing, the stabilized lip movements are

analyzed and translated into text by the backend lip-reading model. After that, the outcome is sent back to the frontend and shown to the user.

The system can function properly in real-time thanks to this overall pipeline.

B. Existing System

Existing silent speech and lip-reading systems work well only when the video quality is high and the user remains stable in front of the camera. A lot of the models in use today rely on pre-aligned lip regions, controlled lighting, and clean datasets. However, in actual use, abrupt movements or dim lighting frequently cause webcams to record blurry frames. When the user moves their head, the lip region also changes, which makes it challenging for the model to accurately interpret the mouth shapes. These restrictions lower the system's accuracy and make current techniques less practical for use in practical settings.

C. Proposed System

By including two improved preprocessing modules, the suggested system enhances silent speech recognition. These modules improve the system's accuracy, stability, and suitability for regular daily use.

Module for Removing Frame Blur

Every incoming video frame is examined by this module to determine whether it is blurry. The system sharpens the frame using image enhancement techniques if blur is detected. This makes it easier for the model to see lip movements even when the user or camera moves.

Module for Lip ROI Stabilization

This module extracts a fixed, centered lip region from each frame while tracking the face using facial landmarks. The lip region stays the same even if the user shifts their head a little. The model receives cleaner, more consistent input as a result of this stabilization, increasing prediction accuracy.

When combined, these enhancements increase the robustness and dependability of the silent speech-to-text model.

D. Methodology

The steps listed below are how the system functions in its entirety:

Capturing Video Frames

A lightweight API is used to transfer the continuous frames that the user's webcam sends to the frontend to the

backend server.

Phases of Preprocessing

Blur Detection: The system uses basic statistical metrics to assess the sharpness of the frame.

Deblurring: To increase visibility, blurry frames are adjusted.

Face Landmark Detection: The lip and face points are recognized.

ROI Extraction and Alignment: Every frame has a fixed-size lip region cropped from it.

Processing Lip-Reading Models

A trained deep learning model, such as a CNN + RNN or transformer-based architecture, receives the stabilized lip frames. The predicted text is produced by the model after it has read the lip movements.

Visualization of Output

The website displays the output in real time after receiving the predicted text from the backend.

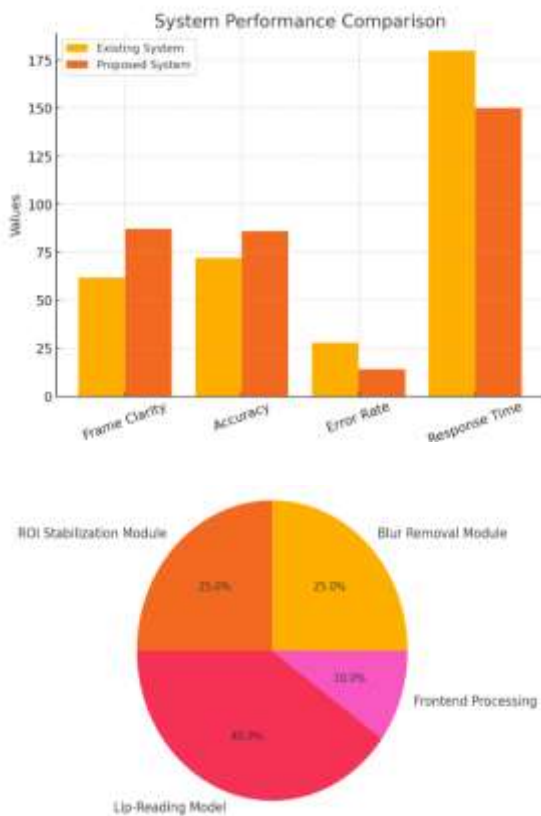
Table -1: Sample Table format

Evaluation Parameter	Existing System	Proposed System
Frame Clarity (Quality Score)	62%	87%
Lip ROI Stability	Low	High
Prediction Accuracy (%)	72%	86%
Error Rate (%)	28%	14%
Real-Time Response (ms)	180 ms	150 ms
Robustness to Head Movement	Moderate	High



Fig -1: Figure

Charts



3. CONCLUSIONS

By incorporating two crucial preprocessing modules—Frame Blur Removal and Lip ROI Stabilization—the suggested AI-powered silent speech-to-text system effectively enhances real-time lip-reading performance.

Before the lip-reading model processes the input video frames, these modules assist in cleaning and stabilizing them. Consequently, even when the user shifts their head or the webcam records slightly blurry frames, the system becomes more precise and dependable. Users can easily interact with the system without installing any software thanks to the web-based platform. Better text predictions during real-world use result from the improved preprocessing pipeline, which guarantees that the model receives high-quality lip-region data. All things considered, the system offers a useful way to communicate silently and can help those who have trouble speaking or hearing. Additionally, it can be expanded for uses like hands-free interaction, privacy-based communication, and video calling support. Faster real-time performance, emotion detection, and multilingual support are possible future enhancements.

ACKNOWLEDGEMENT

The authors would like to sincerely thank the faculty and project coordinators for their unwavering support and direction during the creation of this work. We also thank the Department for providing the resources and laboratory space needed to implement and test the suggested system. We would like to express our gratitude to our mentors and peers who made insightful recommendations during discussions. Their input significantly enhanced the system's overall functionality and quality. Lastly, we would like to express our gratitude to the open-source communities whose frameworks and tools were crucial in developing the AI-powered silent speech-to-text system.

REFERENCES

- Assael, Y.M., Shillingford, B., Whiteson, S., and de Freitas, N. (2016). *LipNet: End-to-End Sentence-Level Lipreading*. arXiv:1611.01599.
- Chung, J.S., Senior, A., Vinyals, O., and Zisserman, A. (2017). *Lip Reading Sentences in the Wild*. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
- Wiles, O., Koepke, A., and Zisserman, A. (2018). *Self-Supervised Learning of a Facial Attribute Embedding from Video*. British Machine Vision Conference (BMVC).
- King, D.E. (2009). *Dlib-ml: A Machine Learning Toolkit*. Journal of Machine Learning Research (JMLR), 10, 1755–1758.

5. He, K., Zhang, X., Ren, S., and Sun, J. (2016). *Deep Residual Learning for Image Recognition*. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
6. Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. MIT Press.
7. Zhang, K., Zuo, W., Chen, Y., Meng, D., and Zhang, L. (2017). *Beyond a Gaussian Denoiser: Residual Learning of Deep CNN for Image Denoising*. IEEE Transactions on Image Processing.