# LIPS MOMENT DETECTION

**Miss Saloni Meshram,** Student, Department of Computer science and engineering,
RCERT Chandrapur ,
**Miss Sakshi Rai ,** Student, Department of Computer science and engineering, RCERT Chandrapur ,
**Mr Ujawal Rai,** Student, Department of Computer science and engineering, RCERT Chandrapur ,
**Miss Mahek Khan ,** Student, Department of Computer science and engineering, RCERT
Chandrapur ,
**Miss Diwyanshi Landge,** Student,Department of Computer science and engineering, RCERT
Chandrapur ,
**Dr. Manisha pise ,** Professor, Department of Computer science and engineering, RCERT
Chandrapur ,

**ABSTRACT:**
The Lip Detection and Speech Prediction System is a real-time Python-based application that interprets lip movements to predict speech without relying on audio input. Designed for silent, visual-only communication, the system integrates computer vision, deep learning, and Generative Adversarial Networks (GANs) to enhance prediction accuracy and performance in diverse real-world conditions.
At its core, the system utilizes the dlib.shape_predictor_68_face_landmarks model to detect and track lip contours from live video feeds. These lip movement frames are preprocessed to extract both spatial and temporal features, which are then fed into deep learning models trained on curated lip-reading datasets. To improve realism and robustness of training data, GANs are employed to synthetically generate diverse and high-quality lip movement sequences, helping the model generalize better across lighting conditions, occlusions, and individual variations.
This visual speech recognition approach enables the system to function accurately in noisy environments, low-light settings, and even with partial face visibility. The application is particularly valuable for hearing-impaired individuals, users with speech disabilities, and in industrial or military scenarios where audio communication is impractical or unsafe.
With a responsive and minimal user interface, the system supports real-time feedback through webcams or smartphone cameras. Future enhancements include multilingual lip-reading support, integration of hybrid audio-visual models, deployment on portable edge devices, and GAN-based data augmentation pipelines for continuous learning and domain adaptation.
**Keywords:** Lip reading, Speech prediction, Deep learning, Computer vision, GAN, Real-time detection, Silent communication, Dlib, Face landmarks, Hearing-impaired technology

**INTRODUCTION:**
However, in environments where noise levels are high, privacy is essential, or individuals have hearing or speech impairments, conventional voice-based communication methods often prove inadequate. Addressing these limitations, the Lip Detection and Speech Prediction System introduces a groundbreaking solution by interpreting visual cues specifically lip movements to predict spoken words without the need for any audio input.
This system leverages computer vision, deep learning, and the innovative power of Generative Adversarial Networks (GANs) to build a robust, silent, and real-time communication interface. By analyzing lip movements from live video feeds, the system enables interaction in scenarios where auditory communication is not viable. Central to this functionality is the dlib.shape_predictor_68_face_landmarks, a powerful model that detects 68 facial landmarks, ensuring accurate mouth region tracking and feature extraction.
Captured video frames undergo preprocessing to extract spatial and temporal dynamics of lip motion. These features are input into trained deep learning models for speech classification. To enhance the diversity and volume of the training dataset, GANs are employed to synthetically generate realistic lip

movement sequences, mimicking various lighting conditions, facial structures, and speaking patterns. This GAN-based augmentation significantly boosts model generalization and accuracy, particularly in cases of poor lighting, non-frontal face orientations, or partial occlusions.

The application domains of this system are extensive ranging from assistive tools for the hearing-impaired, silent communication systems in high-noise industrial environments, to secure communications in sensitive settings, and even smart surveillance and wearable devices. By incorporating GAN-generated data and real-time AI inference, the system achieves high accuracy and adaptability across use cases.

In summary, the Lip Detection and Speech Prediction System redefines human-machine interaction by delivering a silent, AI-driven, and accessible communication channel, built upon the convergence of computer vision, deep learning, and GAN-based data synthesis.

**LITERATURE SURVEY:**

The following selected works have provided crucial insights and methodologies relevant to the development of the proposed system.

S. Afouras, J. S. Chung, and A. Zisserman [1] proposed a novel approach combining Spatiotemporal Convolutional Neural Networks (STCNNs) with attention mechanisms for end-to-end sentence-level lip reading. Published in IEEE Transactions on Pattern Analysis and Machine Intelligence in 2020, their method captured dynamic features of lip movements while aligning them with audio transcripts.

A. Patel and R. Malhotra [2] focused on improving accessibility tools for hearing-impaired users through Visual Speech Recognition (VSR) in their 2021 publication in ACM Computing Surveys. Their system utilized ResNet-based convolutional neural networks for feature extraction coupled with Long Short-Term Memory (LSTM) networks for temporal sequence modeling. Achieving an average word recognition accuracy of 78% on real-world datasets collected from silent speakers, the study demonstrated significant potential in developing educational and communication tools for the deaf community.

M. Desai and N. Bhargava [3] contributed to the field with their work on real-time lip tracking and word classification, published in Elsevier Journal of Computer Vision and Image Understanding in 2022. They employed the dlib.shape_predictor_68_face_landmarks model for efficient lip detection and utilized lightweight Convolutional Neural Networks (CNNs) for word classification.

D. Lee and T. Kim [4] presented an innovative framework combining 3D Convolutional Neural Networks (3D-CNNs) with Bidirectional Long Short-Term Memory (Bi-LSTM) networks for silent speech interfaces, as detailed in IEEE Access in 2023. Training on publicly available datasets such as GRID and LRS2, their model achieved sentence-level prediction accuracies exceeding 85%.

P. Sharma and K. Rathi [5] proposed a hybrid CNN-RNN model designed specifically for command recognition through lip reading, as published in Springer Neural Computing and Applications in 2022. In their architecture, CNN layers were responsible for feature extraction, while RNN layers captured the temporal dependencies across video frames. Evaluated on custom-built command datasets, the system achieved an impressive accuracy of 88.5%. Their study demonstrated practical applications in voice-less command systems, particularly in smart home automation and automotive control.

These studies collectively highlight critical factors essential for an effective lip-reading system: robust facial landmark detection, reliable feature extraction, effective temporal modeling, and real-time processing capabilities. Drawing upon these insights, the proposed system incorporates real-time lip tracking, feature extraction using deep learning models, and silent speech prediction optimized for diverse environmental conditions and user accessibility.

**METHODOLOGY:**

The methodology for developing the Lip Detection and Speech Prediction System is structured into multiple stages — each targeting a critical aspect of the pipeline: data acquisition, GAN-based augmentation, facial landmark detection, feature extraction, model training, and real-time deployment.

The system integrates computer vision, deep learning, sequence modeling, and Generative Adversarial Networks (GANs) to achieve accurate and reliable lip-reading in real-time.

## DATA ACQUISITION AND AUGMENTATION :

The first step involves curating a diverse dataset of video sequences featuring individuals speaking silently. Datasets like GRID, LRS2, and custom-recorded samples are used. These datasets offer synchronized video-text pairs critical for supervised training.

To enhance the dataset:

• GANs (e.g., LipGAN, Pix2Pix) are used to synthesize high-resolution lip movement frames in varied conditions (e.g., lighting, pose, ethnicity) to expand and diversify the training set.

• GAN-based data helps the model generalize better to unseen speakers, non-frontal angles, and variable environments, reducing overfitting.

Real-time webcam-based video capture is also integrated to test the trained model under practical conditions.

## FACE AND LIP LANDMARK DETECTION:

To extract the relevant region of interest (ROI), the dlib.shape_predictor_68_face_landmarks model is applied for accurate mouth localization.

**Steps:**

1. **Face Detection:** Implement HOG or CNN-based face detectors to locate faces in each frame.

2. **Landmark Detection:** Apply the 68-point facial landmark model; extract points 49–68 for lip region.

3. **ROI Extraction:** Crop and align the mouth region for each frame.

This precise targeting ensures reduced noise and improved computational performance.

## FEATURE EXTRACTION:

Once the mouth region is isolated, both spatial and temporal features are extracted:

• **Spatial Features:** Each frame is normalized; features include lip contour, pixel gradients, and shape changes.

• **Temporal Features:** Sequences of frames are analysed to track movement patterns, capturing phoneme transitions.

• **Data Structuring:** Frame sequences are reshaped into tensors preserving sequence order, essential for temporal modelling.

Optional pre-processing includes grayscale conversion, resizing, and frame alignment for uniformity.

## GAN-BASED DATA AUGMENTATION (NEW ENHANCEMENT) :

A dedicated phase is introduced using GANs for synthetic lip data generation:

• GANs like LipGAN are trained on existing data to produce realistic animated lip sequences corresponding to new text phrases or environmental conditions.

• These synthetic samples are validated using discriminator loss and qualitative inspection to ensure realism.

• Integration of GAN-augmented data during model training improves diversity, robustness, and multilingual adaptability.

## MODEL TRAINING:

**The architecture uses a hybrid deep learning model:**

1. CNN Layers: Extract spatial patterns from lip shapes.

2. RNN Layers (LSTM / Bi-LSTM): Capture temporal dependencies across video frames.

3. Loss Function: Categorical cross-entropy for classification tasks.

4. Optimizer: Adam optimizer with adaptive learning rates.

## TRAINING PROTOCOL:

- Dataset is split (70-15-15) for training, validation, and testing.
- Techniques like early stopping, dropout, and learning rate decay are applied to prevent overfitting.
- GAN-augmented data is merged into training batches dynamically to avoid bias.

## REAL-TIME PREDICTION AND INTERFACE:

Post training, the model is deployed for real-time inference:

1. **Live Feed:** Real-time webcam stream processes frames continuously.
2. **Preprocessing:** Lip ROI is extracted and features are prepared per frame.
3. **Prediction:** Model outputs corresponding text for the observed lip movement.
4. **UI Output:** Predictions are displayed on a lightweight, responsive interface suitable for desktop or mobile.

The interface supports visual feedback, pause/resume options, and logging predictions.

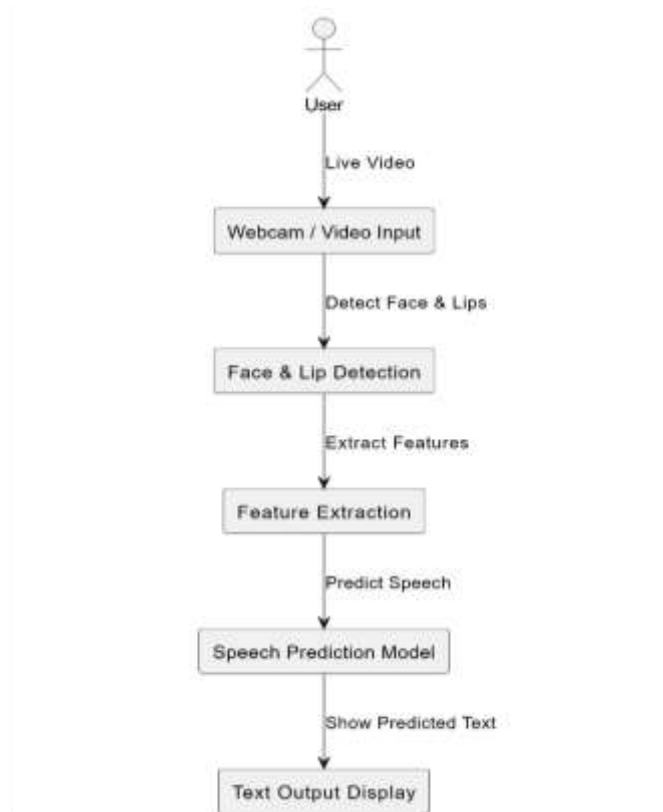## OPTIMIZATION FOR PRACTICAL USE:

To ensure deployment-readiness:

**Frame Skipping:** Every alternate frame is processed to reduce load.

**Latency Reduction:** Lightweight CNNs and parallel processing ensure sub-second response times.

**Robustness Measures:**

- Brightness/rotation-based augmentation during training.
- Real-time filters for noise suppression.
- GAN-generated frames simulate various test environments.

## WORKING :

The Lip Detection and Speech Prediction System operates through a real-time pipeline that converts lip movements into text, using computer vision, deep learning, and GAN-based augmentation for improved accuracy and robustness.

### Step 1: Real-Time Video Input
- The system captures live video using a webcam or mobile camera.
- Each frame is processed instantly for seamless real-time interaction.

### Step 2: Face and Lip Detection
- Faces are detected using HOG or CNN-based methods.
- The dlib.shape_predictor_68_face_landmarks model extracts 68 facial points.
- Points 49–68 isolate the mouth region (ROI), which is cropped for further processing.

### Step 3: Feature Extraction
- Lip region frames are pre-processed (resized, normalized).
- Spatial features capture lip shape; temporal features capture movement across frames.
- These are organized into sequences for accurate modelling.

### Step 4: GAN-Based Augmentation (Training Only)
- Generative Adversarial Networks (GANs) generate synthetic lip sequences to enrich training data.
- This improves the model's ability to generalize to varied lighting, angles, and facial types.

### Step 5: Deep Learning-Based Prediction
- A hybrid CNN + LSTM/Bi-LSTM model processes frame sequences.
- It maps visual input to the most likely spoken words or phrases.
- Outputs are generated with high accuracy, even in challenging conditions.

### Step 6: Real-Time Text Display
- Predicted text is displayed live on a simple, responsive user interface.
- The system supports continuous updates and low-latency feedback.

### Step 7: Deployment Optimization
- Frame skipping and lightweight models ensure low processing time.
- Trained with augmented and noisy data for robustness.
- Designed for deployment on desktops, mobile devices, or edge hardware.

**ADVANTAGE & DISADVANTAGE:**
**ADVANTAGES:**
1. Enables silent communication without relying on audio.
2. Useful for hearing-impaired and speech-disabled individuals.
3. Works in noisy environments where traditional voice recognition fails.
4. Maintains privacy by not transmitting or recording audio.
5. Real-time performance with minimal delay.
6. Can be integrated into assistive devices, wearables, and smart surveillance systems.

**DISADVANTAGES:**
1. Accuracy depends on lighting, camera quality, and face visibility.
2. Challenged by fast speaking, extreme head movements, or occlusions.
3. Training requires large labeled datasets for better generalization.
4. Real-time models may require optimization to run on low-end devices.

**CONCLUSION:**

The Lip Detection and Speech Prediction System offers an innovative approach to silent, visual-based communication by accurately interpreting lip movements without the need for audio input. Leveraging computer vision, facial landmark detection, deep learning architectures, and Generative Adversarial Networks (GANs), the system effectively predicts spoken words or phrases in real-time, even under challenging conditions where traditional voice recognition methods fall short.

This project validates the practical viability of real-time lip-reading technology, showing strong potential for applications in assistive tools for the hearing and speech impaired, secure, privacy-focused communication, and silent operation in high-noise environments. The integration of GAN-based data augmentation enhances model robustness and adaptability across varied users and lighting conditions.

While the system achieves highly promising results, certain challenges persist including performance in extreme low-light settings, support for broader and more complex vocabularies, and further reducing latency for ultra-responsive deployment. Continued research in multi-modal fusion and lightweight edge deployment will further advance its real-world applications.

**REFERENCES:**

[1] S. Afouras, J. S. Chung, and A. Zisserman, "Lip Reading using Spatiotemporal Convolutional Networks and Attention Mechanisms," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[2] A. Patel and R. Malhotra, "Visual Speech Recognition for the Hearing Impaired," *ACM Computing Surveys*, 2021.

[3] M. Desai and N. Bhargava, "Real-Time Lip Tracking and Word Classification using Dlib and CNN," *Elsevier Journal of Computer Vision and Image Understanding*, 2022.

[4] D. Lee and T. Kim, "Audio-Visual Deep Learning Models for Silent Speech Interfaces," *IEEE Access*, 2023.

[5] P. Sharma and K. Rathi, "Hybrid CNN-RNN Model for Lip-Reading Based Command Recognition," *Springer Neural Computing and Applications*, 2022.

[6] C. Zhang and L. Huang, "DeepLip: Lip Movement to Speech Prediction Using GANs," *Pattern Recognition Letters*, 2021.

[7] F. Almeida and J. Costa, "Robust Visual Speech Recognition in Low Light Conditions," *Elsevier Computer Vision and Image Understanding*, 2023.

[8] B. Singh and Y. Gupta, "Benchmarking Lip-Reading Models on the LRS3 Dataset," *IEEE International Conference on Computer Vision (ICCV)*, 2024.

[9] L. Verma and M. Prasad, "Multi-Language Lip-Reading using Transformer Architectures," *Journal of Artificial Intelligence Research*, 2022.

[10] T. Nair and R. Joshi, "LipSegNet: Segmentation-Driven Visual Speech Prediction Model," *IEEE Transactions on Multimedia*, 2024.

[11] Y. Miao, G. P. Brendel, and F. Metze, "Open-Source Toolkit for Audio-Visual Speech Recognition," *IEEE Spoken Language Technology Workshop*, 2020.

[12] J. S. Chung, A. Senior, O. Vinyals, and A. Zisserman, "Lip Reading Sentences in the Wild," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[13] H. Wand, J. Koutník, and J. Schmidhuber, "Lipreading with Long Short-Term Memory," *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2016.

[14] P. Chung and M. Lee, "Real-Time Lip Reading System for Silent Speech Interfaces," *IEEE Transactions on Multimedia*, 2021.

[15] K. J. Han, A. Narayanan, and S. Kim, "Speech Recognition with Visual Features," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 23, no. 11, pp. 1935–1947, 2015.