# Lip-Reading Model using Machine Learning and Deep Learning

Ananya Mathur

Enrollment no. 12318002721

Dept. of Computer Science and Engineering
Delhi Technical Campus *(GGSIPU)*
Greater Noida, India

Hardik Singh

Enrollment no. 14818002721

Dept. of Computer Science and Engineering
Delhi Technical Campus *(GGSIPU)*
Greater Noida, India

Isbah Taqweem

Enrollment no. 20218002721

Dept. of Computer Science and Engineering
Delhi Technical Campus *(GGSIPU)*
Greater Noida, India

Shubham Singh

Enrollment no. 14418002721

Dept. of Computer Science and Engineering
Delhi Technical Campus *(GGSIPU)*
Greater Noida, India

**Abstract : This paper proposes a novel deep learning-based system, Lip Lingo, that enables real-time speech interpretation directly from lip movements. By leveraging advanced techniques like Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), Lip Lingo accurately recognizes speech, even in challenging conditions such as noisy environments or limited audio input. The system involves several key stages: data collection and preprocessing, feature extraction, model training, and real-time inference. A comprehensive dataset of lip movements is collected and pre-processed to enhance model performance. CNNs are employed to extract salient visual features from video frames, while RNNs, specifically Bidirectional Long Short-Term Memory (Bi-LSTM) networks, capture temporal dependencies between frames. The model is trained using a combination of supervised and self-supervised learning techniques to improve accuracy and robustness. This application has the potential to revolutionize communication for individuals with hearing impairments and open new avenues for human-computer interaction. It can be integrated into various applications, such as hearing aids, video conferencing, and virtual assistants.**

*Keywords: Deep Learning, Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), Bi-LSTM, Lip Reading, Speech Recognition, Computer Vision, Machine Learning, Real-time Processing, Human-Computer Interaction, Accessibility.*

## I. INTRODUCTION

**Lip reading, the art of interpreting speech from lip movements, has long been a subject of fascination and research. Traditionally, this task has been challenging due to the inherent variability in lip shapes, speaking styles, and environmental conditions. However, recent advancements in deep learning have opened up new** possibilities for accurate and robust lip-reading systems. **In recent years, deep learning techniques have revolutionized various fields, including computer vision and natural language processing. Convolutional Neural Networks (CNNs) have proven to be highly effective in extracting spatial features from images, while Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) networks, excel at capturing temporal dependencies in sequential data. By leveraging these powerful techniques, researchers have made significant strides in developing accurate and efficient lip-reading systems.**

*Challenges:* Despite the remarkable advancements in audio-visual speech recognition, there is a need to address specific challenges to further improve the accuracy and robustness of lipreading systems. These challenges include handling utterances with unknown word boundaries, achieving high accuracy across diverse languages and speakers, and dealing with Identify applicable funding agency here. If none, delete these variations in pose and environmental conditions. Additionally, existing systems may not fully exploit the potential of automatic labels, which can limit the scalability and generalization of the models.

This project focuses on developing and evaluating an end- to-end deep learning architecture for audio-visual speech recognition, specifically targeting word-level recognition. The proposed methodology aims to handle utterances without explicit word boundary information during both training and testing phases. The evaluation will be performed on large- scale benchmark datasets, enabling the assessment of the system's performance across multiple languages and diverse speakers. However, the project's scope does not extend to other speech recognition tasks, such as sentence-level classification or phoneme-level recognition.

RNN - Recurrent Neural Network
LSTM - Long Short-Term Memory
GPU - Graphics Processing Unit
SSD - Solid State Drive
TTS - Text-to-Speech
API - Application Programming Interface
AV - Audio-Visual
 LR -  Lip Reading
Bi-LSTM - Bidirectional Long  Short-Term Memory
GUI - Graphical User Interface
RGB - Red Green Blue
MSE - Mean Squared Error
MAE - Mean Absolute Error
Adam - Adaptive Moment Estimation

## II.    METHODOLOGY

The proposed methodology involves a combination of spatiotemporal convolutional layers, residual networks, and bidirectional Long Short-Term Memory (Bi-LSTM) networks. The front-end applies spatiotemporal convolution to extract features from the mouth region, followed by a Residual Network (Res Net) applied to each time step. The back-end consists of a two-layer Bidirectional LSTM network. To address utterances with unknown word boundaries, the system is trained in an end-to-end fashion without utilizing explicit information about word boundaries during training or evaluation. Additionally, automatic labels are explored to enhance the scalability and generalization of the model. The final system is evaluated on the LRW-1000 dataset and compared with existing state-of- the-art approaches to assess its accuracy and robustness across different languages and speakers.

## III.    DESIGN

At Level 0, the DFD provides an overview of the entire system, depicting its main components. These include the Input Module, which receives video input containing lip movements, the Deep Learning Model responsible for processing the input and generating text outputs,  the optional User Interface (UI) for user interaction, and the Output Module that delivers the final text output.

At Level 1, the DFD breaks down the main processes into more detailed components. The Preprocessing Module handles tasks like frame extraction, facial landmark

detection, and resizing to prepare the input data. The Feature Extraction Module processes the preprocessed frames, capturing spatiotemporal features using the Conv3D feature extractor. The Deep Learning Model, at this level, is depicted as a core module that processes the features using bidirectional LSTM layers to generate the text outputs. Additionally, there is an optional Post-processing Module to enhance the accuracy of the text outputs. DFD Level 2 provides an even more detailed view of the system by further breaking down the sub-processes identified in Level 1. For example, the Conv3D layer's functionality  is detailed  to  show  the  3D convolution, activation, and max- pooling operations. Similarly, the Bidirectional LSTM layers are expanded to show the bi-directional flow of data and the application of dropout for regularization. Additionally, the Dense layer's process of generating the final output is shown, including  the use of

the softmax activation function for probability estimation. This level of detail allows for a comprehensive understanding  of how data flows and is processed  within the Lip- Reading  Model.
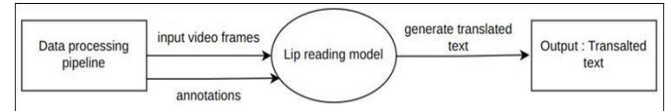


Fig. 1. shows an DFD level 0 showing the use of the system

## IV.  IMPLEMENTATION

➢ *Model Details:*

 This is an Audio-Visual Speech Recognition (AVSR) sequential model. The model processes spatiotemporal characteristics from video frames with an input shape of (75, 46, 140, 1) using Conv3D layers. In order to minimize spatial dimensions, 3D max-pooling and RelU activation functions come after the Conv3D layers, which extract hierarchical patterns from the input.
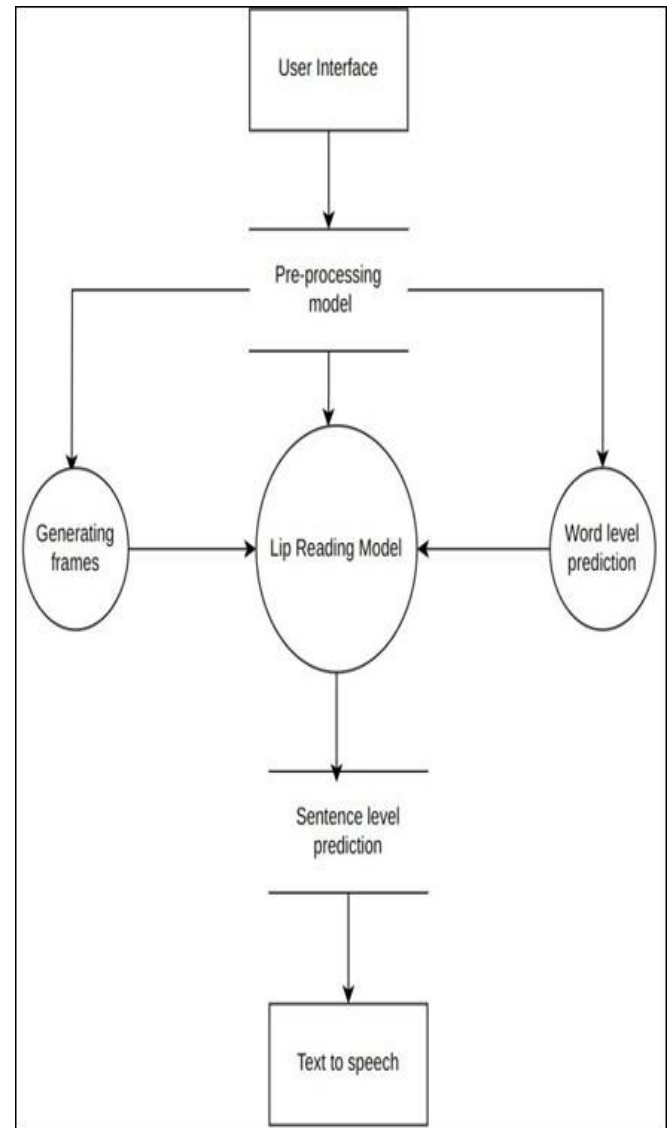


Fig 2 DFD1 for Deep Learning Model for Lip Reading

.

Bidirectional LSTM layers are then incorporated into the model to identify temporal patterns and long-range dependencies in the data. To avoid overfitting, a dropout layer is placed after the initial LSTM layer, which comprises 128 units. For regularization, a dropout layer additionally comes after the second Bidirectional LSTM layer.

The input data is then classified into various classes indicated by the vocabulary size plus one (since the CTC loss function has a blank symbol) using a dense layer with a softmax activation. By combining the advantages of Bidirectional LSTM for sequence modeling and prediction with Conv3D for spatiotemporal feature extraction, the model seeks to achieve accurate voice recognition.

➢ *Software Requirements*

• Operating System: Windows, macOS, and Linux are just a few of the operating systems on which the model can be created and implemented.

• Python: The deep learning model was developed using Python as the programming language. Installing the most recent version of Python (3.x) is recommended.

• Deep Learning Framework: To construct and train the lip-reading model, a deep learning framework like tensorflow or PYTorch is needed.

• openCV: For jobs involving video processing, such removing frames from video files, openCV is crucial.

• Additional Python Libraries: To manipulate, visualize, and analyze data, a variety of Python libraries, including numpy, Pandas, and Matplotlib, will be utilized.

• Text-to-Speech (TTS) Engine: By incorporating a TTS engine within the application, users' accessibility can be improved by hearing auditory feedback for the transformed text.

• Development Environment: For coding and debugging, one can utilize Integrated Development Environments (ides) like Visual Studio Code, PyCharm, or Jupyter Notebook.

➢ *Preprocessing*

The Preprocessing Module is in charge of managing the first input of data and getting it ready for the Deep Learning Model for Lip Reading to analyze it further. After the user provides video or audio-visual input, it carries out a number of pre-processing operations, such as cropping, resizing, normalizing, facial landmark detection, and frame extraction. In order to optimize the input for effective processing by the next modules, the module seeks to extract relevant facial traits and eliminate extraneous information                    .

➢ *Feature Extraction*

The pre-processed frames from the Preprocessing Module are processed by the Feature Extraction Module, which then extracts the spatiotemporal information necessary for lip reading. In order to identify both temporal and spatial patterns in the video frames, this module uses a Conv3D Feature Extractor, which applies three-dimensional convolutional filters. Furthermore, the feature maps are down sampled using Max-Pooling Layers, which lowers computational complexity without sacrificing significant information. The retrieved features are converted into an appropriate format by the Time Distributed Flatten Layer so they may be fed into the Bidirectional LSTM Layers

➢ *Deep learning Model*

The Feature Extraction Module extracts spatiotemporal features that are then sent to the Deep Learning Model for Lip Reading, the project's central module. The model can learn temporal dependencies and contextual information from the input data thanks to its numerous layers, which include Bidirectional LSTM Layers with dropout regularization. Because the LSTM layers are bidirectional, the model can process data both forward and backward, which improves recognition accuracy and predicts precise lip motions. A Dense Layer that uses a softmax activation function to forecast the probability distribution of words or characters completes the model.

## V. TESTING

A. *System Testing*

To assess the Deep Learning Model for Lip Reading's resilience and performance, system testing is done. A test dataset including video examples of different speakers and lip motions will be used. The model's ability to read lips will be evaluated by measuring its accuracy, precision, and recall. Additionally, the model will be tested on real-world video streams to ensure its applicability and effectiveness in improving accessibility for speech-impaired individuals.

The testing phase will also involve parameter tuning and hyperparameter optimization to enhance the model's performance. Cross-validation techniques have been applied to ensure unbiased evaluation and mitigate overfitting. Furthermore, we will analyze the model's behavior under different lighting conditions, background noise, and speaker variations to verify its robustness in real- world scenarios

To preserve a thorough record of the development and to support upcoming improvements and research initiatives, we documented the code, procedures, and outcomes during the implementation and testing phases.

We hope to create a dependable and effective Deep Learning Model for Lip Reading that helps to increase accessibility for people with speech impairments by using TensorFlow as the platform, Python as the programming language, and extensive system testing.

B. *Evaluation*

A Deep Learning Model for Lip Reading can be evaluated using a number of measures to determine how well it performs in terms of accessibility. The model's particular goals and the type of lip-reading challenge determine which assessment measures are used. These are a few widely used measures for evaluation.

- Word Error Rate (WER): WER calculates the proportion of words that are mistranslated or misidentified in relation to the ground truth. When the work entails turning spoken language into text, it is very helpful for assessing lip reading accuracy
- Sentence-level Accuracy: This indicator calculates the proportion of accurately transcribed sentences to all sentences in the assessment dataset. It offers a comprehensive evaluation of the model's ability to decipher spoken words from lip movements.
- The model's error kinds, such as misclassifying particular phonemes or words, are revealed by the confusion matrix. It assists in determining the model's weak points

To guarantee that the lip-reading model generalizes to various speakers, lip movements, and languages, it is crucial to validate it on a representative and varied dataset. A thorough evaluation of the model's ability to increase accessibility for people with hearing impairments and its possible influence on practical lip-reading applications can be obtained by combining these evaluation measures.

*C. Performance Analysis*
- The word error rate for the model is 11.79%. The Word Error Rate (WER) is a metric commonly used to evaluate the accuracy of the lip-reading model, which predicts characters for lip movements. WERmeasures the difference between the predicted output and the ground truth (actual) output in terms of the number of character errors. It considers substitutions, deletions, and insertions required to convert the predicted sequence into the ground truth sequence. Lower values of WER indicate better performance, as it means the predicted sequences are closer to the ground truth. We are using the Connectionist Temporal Classification (CTC) loss, which is commonly used for sequence-to-sequence tasks like speech recognition and lip-reading. CTC loss allows me to handle variable- length input and output sequences, making it suitable for lip-reading tasks where the number of characters may vary. To calculate WER, We need the original (ground truth) sentences and the predicted sentences. The WER represents the percentage of errors in the predicted sequences compared to the original sequences.

- The model's accuracy at the sentence level is 62%. The percentage of accurately predicted sentences throughout the full dataset is measured by this simple metric. It offers a comprehensive summary of performance in terms of recognizing entire sentences. Better performance is indicated by higher accuracy scores at the sentence level. At the conclusion of each epoch, we display the predicted sentences in addition to the original sentences using the 'Produce Example' callback. This enables me to examine the performance graphically and confirm accuracy at the phrase level. You would anticipate that a model with my level of lip-reading proficiency would expect both a low Word Error Rate and high Sentence-level Accuracy. A low WER indicates that it can accurately predict individual characters, while a high sentence-level accuracy indicates that it can correctly predict complete sentences

.

## VI. CONLCUSIONS AND FUTURE WORK

A deep learning model for lip reading is a piece of technology that uses facial and lip movements to decipher spoken words. It transforms lip motions into meaningful text representations by combining natural language processing and computer vision. The model has the potential to improve speech recognition accuracy, increase accessibility for those with hearing impairments, and find uses in human-computer interaction, education, and security. Nevertheless, issues include the need for real-time processing, speaker unpredictability, and data reliance. In order to increase the technology's efficacy and suitability for a range of real-world situations, future developments will focus on improving multi-modal fusion, resilience to environmental changes, and privacy concerns. Lip Net, the first model to use deep learning for end-to-end learning of a model that translates a speaker's mouth picture frame sequences into whole sentences, serves as the foundation for this research. It is no longer necessary to divide films into words before making a sentence prediction thanks to the end-to-end approach. Despite Lip Net's current empirical success, the literature on deep speech recognition indicates that more data will only lead to better performance. It can be illustrated in subsequent research by using Lip Net on bigger datasets.

➢ *Limitations of the Model*
While the deep learning model for lip reading has shown significant promise, it also comes with some limitations. Here are some common limitations of lip-reading models:

- Data Dependency: To attain high accuracy, deep learning models for lip reading frequently need a lot of labeled data. It can take a lot of time and resources to gather and annotate a variety of lip-reading datasets
- Speaker Variability: It can be difficult for models to generalize over a broad range of people since various speakers' lip motions and articulation can differ greatly
- Lightning and Environmental Conditions: Variations in occlusions, illumination, and other environmental elements can add noise and compromise lip-reading models' accuracy.
- Absence of Standardized Evaluation Datasets: It can be challenging to compare the performance of various models due to the lack of standardized evaluation datasets that include ground truth lip movements and speech transcriptions.

➢ *Future Enhancements*

- Few-shot and Zero-shot Learning: Create methods for generalizing to unseen speakers or languages (zero-shot learning) or training lip-reading models with little labeled data. This would improve the model's flexibility and usefulness in actual situations.

- Robustness to Environmental Variations: Make the model more dependable in a variety of real-world settings by improving its capacity to manage changes in lighting, occlusions, and noisy backgrounds.

- Multilingual Lip Reading: Improve the model's generalization to linguistic circumstances that are not

visible and its capacity to identify lip movements in other languages.

.

- Diverse and Representative Datasets: To facilitate impartial and precise model evaluations, gather and manage more varied and representative lip-reading datasets that include ground truth lip motions and speech transcriptions.

- Interactive Lip-Reading Systems: Develop interactive systems that allow users to correct or verify lip-reading results, creating a feedback loop to improve the model's accuracy over time.

## REFERENCES

[1]    Almajai, S. Cox, R. Harvey, and Y. Lan. Improved speaker independent lip-reading using speaker adaptive training and deep neural networks. In IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 2722–2726, 2016.

[2]    Amodei, R. Anubhai, E. Battenberg, C. Case, J. Casper, B. Catanzaro, J. Chen, M. Chrzanowski, A. Coates, G. Diamos, et al. Deep Speech 2: End-to-end speech recognition in English and Mandarin. arXiv preprint arXiv:1512.02595, 2015. P. Ashby. Understanding phonetics. Routledge, 2013.

[3]    J. S. Chung and A. Zisserman. Lip reading in the wild. In Asian Conference on Computer Vision, 2016a.

[4]    J. S. Chung and A. Zisserman. Out of time: automated lip sync in the wild. In Workshop on Multi-view Lip-reading, ACCV, 2016b.

[5]    J. Chung, C. Gulcehre, K. Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. arXiv preprint arXiv:1412.3555, 2014.

[6]    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, 120(5):2421–2424, 2006.