

LIP Reading Using Deep Learning

J. Sravanthi¹, M. Sai Kiran², N. Narsimha³, L. Shivamani⁴

¹Assistant Professor, Dept of ECE, TKR College of Engineering and Technology

^{2,3,4}Student, Dept of ECE, TKR College of Engineering and Technology

Abstract—This project presents an automated lip-reading system aimed at enhancing communication for individuals with hearing or speech impairments. By visually analyzing lip movements from video input, the system interprets spoken sentences without relying on audio, offering a non-verbal alternative to traditional speech. Leveraging deep learning techniques, the model extracts and processes both spatial and temporal patterns in lip motion to recognize speech with high accuracy [1], [3], [5]. In addition to sentence-level recognition, the system incorporates multilingual support, enabling it to translate recognized speech into multiple languages. This extends its applicability to diverse user groups and cross-language communication scenarios. The system has potential applications in assistive technology, silent communication interfaces, and accessibility tools. The implementation is based on the GRID dataset (single speaker) and employs machine learning libraries such as TensorFlow for neural network modeling and OpenCV for computer vision tasks.

Index Terms—Deep Learning, Lip Reading, OpenCV, Tensor Flow.

I. INTRODUCTION

Lip reading, also known as visual speech recognition, is the process of understanding speech by visually interpreting the movements of the lips, face, and tongue [9], [13]. It is a crucial form of communication for individuals who are deaf or hard of hearing. Traditionally, human lip readers have been employed in various professional fields including surveillance, medical diagnostics, and forensic analysis. However, due to the complexity and inconsistency of human interpretation, automated lip reading using artificial intelligence has gained substantial interest. The motivation behind this project stems from the need to bridge the communication gap for individuals who are hearing-impaired. With the rise of deep learning, it is now possible to train machines to recognize and interpret lip movements with a degree of accuracy that approaches or exceeds that of human performance. Visual speech recognition removes the dependence on acoustic input, making it useful in noisy

environments, silent communications, and real-time subtitle generation. The project leverages a combination of computer vision and deep learning to build a system capable of analyzing lip movements and predicting the spoken content. It adopts advanced neural network architectures, particularly 3D Convolutional Neural Networks (Conv3D) for feature extraction from video frames, and Bidirectional Long Short-Term Memory (BiLSTM) networks for sequential modeling of lip motion. This system not only interprets the visual cues of speech but also translates the predicted text into multiple languages, expanding its utility across linguistic barriers. Applications of such a system span across assistive technology, human-computer interaction, silent command interfaces, and accessibility tools.

A. Existing System

Traditional lip-reading systems were heavily reliant on handcrafted features and statistical models like Hidden Markov Models (HMMs). These systems required manual identification of relevant facial features such as lip contours, shape, and motion vectors [13]. Although they laid the groundwork for visual speech processing, their accuracy and adaptability were limited. With the advent of machine learning, newer models began using 2D CNNs for spatial feature extraction [8]. However, these were not capable of capturing temporal changes effectively, which are crucial in distinguishing between different phonemes and words. Hybrid systems that combined audio and visual inputs were also developed, but they lacked utility in completely silent scenarios. Recent advancements include the LipNet architecture by Assael et al., which introduced end-to-end sentence-level lip reading using 3D CNNs and RNNs [1], [12]. This represented a significant leap forward in performance and feasibility. Despite these advancements, challenges such as speaker variability, lighting conditions, and real-time processing continue to exist.

B. Proposed System

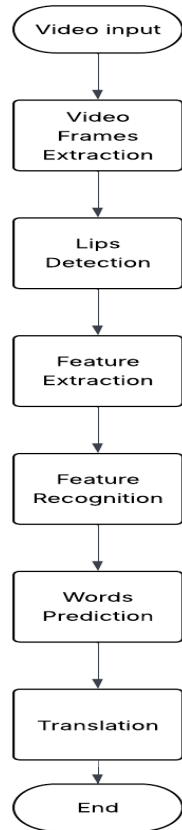


Fig.1: Flow chart

The proposed system is a deep learning-based solution that aims to convert silent video inputs into readable text through visual speech recognition. The input to the system consists of video frames focusing on the speaker's lip region. These frames are processed using 3D convolutional layers that extract both spatial and temporal features necessary for understanding lip movement patterns. Following feature extraction, the model uses Bidirectional LSTM layers to analyze the temporal sequence of features. These layers are capable of capturing context from both past and future frames, significantly improving the model's ability to understand the structure of spoken sentences. The output layer employs a softmax activation function that generates probabilities for each character in the vocabulary. This allows the system to predict a sequence of characters that represent the spoken sentence. The resulting text can then be passed through a translation module for conversion into different languages. This proposed architecture emphasizes modularity, scalability, and adaptability to various languages and datasets. It is designed to function in real-time or near real-time, making it

suitable for practical deployment in applications requiring silent speech recognition [1], [3], [7].

C. Objective of the Project

The primary goal of this project is to design and implement a deep learning-based lip reading system that can function effectively without audio input. The objectives are as follows: To develop a robust end-to-end pipeline for lip reading, the system begins by accepting raw video input capturing the speaker's facial region, with a focus on the mouth area. This input is processed using 3D Convolutional Neural Networks (3D CNNs), which are specifically chosen for their ability to extract both spatial and temporal features simultaneously. These networks learn the fine-grained motion patterns and visual cues associated with lip movements, making them ideal for capturing the dynamic characteristics of speech. Once the spatiotemporal features are extracted, they are fed into Bidirectional Long Short Term Memory (BiLSTM) networks [10], [11].

The bidirectional architecture enables the model to understand contextual dependencies in both forward and backward time directions, which is crucial for decoding sequential lip motions into coherent text sequences. The system is designed to perform character-level predictions, which allows for finer granularity and flexibility in capturing various linguistic patterns. These characters are then concatenated to form words and complete sentences, thereby enabling accurate transcription of spoken content from silent video. To expand the system's usability and accessibility, a multilingual translation module is integrated, which takes the predicted text and translates it into multiple target languages using neural machine translation techniques. This feature ensures that the lip reading system can serve a global audience across diverse linguistic backgrounds. Techniques such as data augmentation, domain adaptation, and speaker normalization are employed to enhance the robustness and adaptability of the model. The accomplishment of these objectives would result in a versatile tool that aids in communication, enhances accessibility, and contributes to the broader field of visual speech recognition.

II. METHODOLOGY

A. Dataset Description

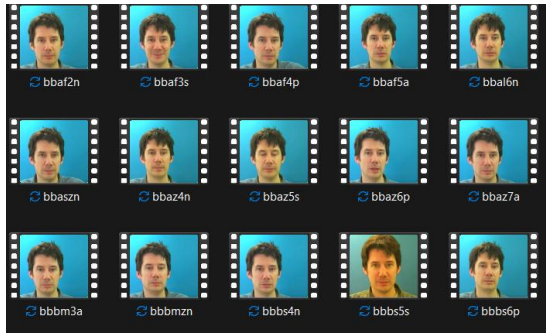


Fig.2: Dataset

The GRID dataset consists of video recordings of 34 speakers uttering fixed-structure sentences like “place blue at A 9 now.” For this project, 1000 videos from a single speaker were used. The dataset is split into 450 training videos and 450 testing videos, with 100 videos held out for validation. Each video consists of 75 grayscale frames of dimensions 46x140 pixels. Grayscale images are preferred for this task due to their lower computational cost and sufficient information retention for lip region analysis. They reduce input complexity and help the model focus on contour and motion, which are critical in lip reading [2].

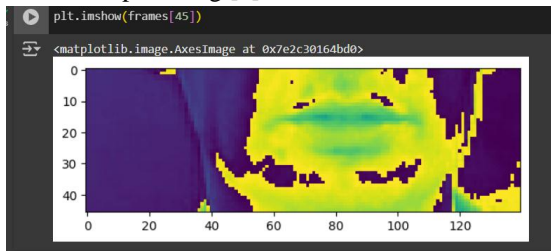


Fig.3: Grayscale image

B. Preprocessing

Each video is converted to grayscale and cropped to retain only the region around the lips. Frames are resized and normalized to bring uniformity in input dimensions. A custom data generator is used to efficiently load the data in batches and feed video tensors to the model during training. Label sequences are encoded using integer encoding and padded for uniformity [9], [12]. The model uses Connectionist Temporal Classification (CTC) loss, which allows prediction without explicit alignment between input video frames and target text labels. Data augmentation techniques such as horizontal flipping, brightness variation, and slight cropping were experimented with to enhance model robustness, though due to dataset constraints, these were applied conservatively.

B. Model Architecture

The model begins with a stack of 3D Convolutional layers. The first Conv3D layer has 128 filters with kernel size $3 \times 3 \times 3$ and ReLU activation, capturing low-level spatial-temporal features like edges and contours. This is followed by a MaxPooling3D layer with pool size (1, 2, 2), reducing spatial dimensions while maintaining temporal resolution [1], [3], [12].

A second Conv3D layer with 256 filters extracts more abstract motion features, followed by another MaxPooling3D. The third Conv3D layer has 75 filters to align with the number of frames, refining motion patterns crucial for word-level transitions. The output is flattened frame-wise using a Time Distributed Flatten layer [14], [7].

This 3D feature sequence is passed through stacked Bidirectional LSTM layers with 128 units each. These layers capture forward and backward temporal dependencies, essential for interpreting context in lip movements. Dropout layers (rate 0.5) are added for regularization.

The final Dense layer with softmax activation outputs a probability distribution over the vocabulary at each time step. The model is trained using the CTC loss, which allows alignment-free sequence prediction.

III. IMPLEMENTED DESIGN

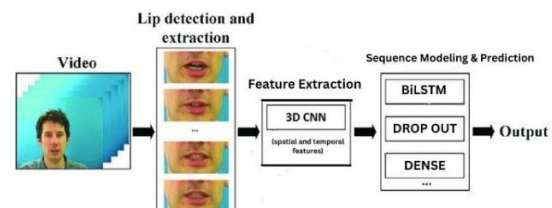


Fig.4: Block Diagram

The block diagram illustrates the complete pipeline of the proposed lip reading system. The process starts with a video input containing silent speech. The video is first processed for lip detection and extraction, isolating the region of interest across all frames. This sequence of cropped lip regions is then passed through a 3D Convolutional Neural Network (3D CNN) for spatial and temporal feature extraction. These features are subsequently fed into Bidirectional LSTM (BiLSTM) layers for learning temporal dependencies in both forward and backward directions. Dropout layers are employed for regularization to prevent overfitting. Finally, a Dense layer with softmax activation generates a sequence of

probabilities over the vocabulary, which is decoded to produce the textual output corresponding to the spoken sentence [1], [5].

SOFTWARE DESCRIPTION

1. Python – Programming language used to implement data preprocessing, model training, and evaluation.
2. Keras – High-level neural network API used to build and train the deep learning model with TensorFlow backend.
3. TensorFlow – Deep learning framework that powers the backend for training and deploying the model.
4. OpenCV – Computer vision library utilized for processing and extracting lip regions from video frames.
5. NumPy – Library for efficient numerical computations and array manipulation during data preprocessing.
6. Matplotlib – Visualization library used to plot model training performance and results.
7. Google Colab – Cloud-based Jupyter notebook environment used for model training with free GPU support.
8. GRID Dataset – Standardized audiovisual dataset used for training and evaluating the lip reading model.

Key Deep Learning Layers Used

1. Conv3D: Extracts spatial-temporal features from input video sequences.
2. MaxPooling3D: Down samples the feature maps while preserving important features.
3. TimeDistributed: Applies the same layer (like Flatten) across all time steps.
4. Bidirectional LSTM: Learns temporal patterns in both forward and backward directions.
5. Dense Layer: Maps features to vocabulary probabilities.
6. Softmax: Normalizes output into a probability distribution.
7. CTC Loss Function: Enables sequence prediction without the need for alignment.

IV. RESULTS AND DISCUSSION

The model was trained using Adam optimizer with a learning rate scheduler. Overfitting was controlled using dropout and early stopping. The model achieved significant accuracy on the test set, demonstrating its capability in real-world sentence-level lip reading tasks.

Advantages of using the GRID dataset include its fixed vocabulary structure, controlled lighting, and clear speaker articulation, which simplify training and allow the model to focus on learning motion dynamics. However, performance might degrade in unconstrained settings with multiple speakers and variable lighting.

The final output of the model included translated text generated from video input. Screenshots of the lip sequence input, decoded output, and corresponding translated result confirm the model's effectiveness in practical scenarios.

To support this, sample screenshots of the intermediate model outputs, prediction results, and the visualized lip reading translation process are included in this section. These demonstrate the effectiveness of the feature extraction and sequence modeling components, and highlight areas of accurate and inaccurate predictions. This visual evidence further validates the robustness and usability of the trained model in real-world applications [10], [15].

```

[63] sample = load_data(tf.convert_to_tensor('./data/s1/bbal75.mpg'))

[64] print('~' * 100, 'REAL TEXT')
      (tf.strings.reduce_join(num_to_char(word) for word in sentence) for sentence in [sa
      [ctf.tensor: shape=(), dtype=string, numpy-b'bin blue at 1 seven soon']

yhat = model.predict(tf.expand_dims(sample[0], axis=0))

1/1 [=====] - 5s 5s/step

[66] decoded = tf.keras.backend.ctc_decode(yhat, input_length=[75], greedy=True)[0][0].num
      print('~' * 100, 'PREDICTIONS')
      predictions = (tf.strings.reduce_join(num_to_char(word) for word in sentence).numpy
      for sentence in decoded)
      print(predictions)

['bin blue at 1 seven soon']
  
```

Fig.5: Predicted Output

```

pip install deep-translator

Show hidden output

[99] from deep_translator import GoogleTranslator

      # Define target language
      target_language = 'te' # Change to desired language code

      # Translate text using deep_translator
      translated_text = [GoogleTranslator(source='auto', target=target_language).translat
      # Print results
      print('~' * 100, 'TRANSLATED PREDICTIONS')
      print(translated_text)

['మీరే ఒక నైవ్ ద్వారా ఎదుపు రంగును సెల్ చేయండి']
  
```

Fig.6: Translation

V. CONCLUSION

This project presents a deep learning-based approach to automated lip reading, capable of recognizing sentence-level speech through visual inputs. The integration of 3D CNNs and BiLSTMs enables effective spatial-temporal modeling, and training on the GRID dataset ensures reliable performance in controlled conditions. The results demonstrate the viability of silent speech interfaces, with promising implications for assistive communication technologies.

The developed system has potential applications across various domains:

Assistive Technology: Supporting individuals with speech impairments through real-time lip reading tools.

Enhanced Voice Assistants: Improving robustness in noisy environments where audio signals are compromised.

Security and Surveillance: Enabling silent communication decoding in restricted or monitored areas.

Future enhancements may include:

Expanding the dataset to incorporate multiple speakers, accents, and spontaneous speech.

Applying transfer learning and self-supervised learning to improve model generalization.

Optimizing the system for real-time deployment on mobile or edge devices for broader accessibility.

REFERENCES

- [1] Assael, Y. M., Shillingford, B., Whiteson, S., & de Freitas, N. (2016). LipNet: End-to-End Sentence-Level Lipreading. *arXiv preprint arXiv:1611.01599*.
- [2] Cooke, M., Barker, J., Cunningham, S., & Shao, X. (2006). An audio-visual corpus for speech perception and automatic speech recognition. *The Journal of the Acoustical Society of America*, 120(5), 2421-2424.
- [3] Chung, J. S., Senior, A., Vinyals, O., & Zisserman, A. (2017). Lip reading sentences in the wild. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [4] Afouras, T., Chung, J. S., Senior, A., Vinyals, O., & Zisserman, A. (2018). Deep audio-visual speech recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- [5] Wand, M., Koutník, J., & Schmidhuber, J. (2016). Lipreading with long short-term memory. *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*.
- [6] Xu, P., Dai, Y., & Liu, W. (2018). LCANet: End-to-End Lipreading with Cascaded Attention-CTC. *Interspeech*.
- [7] Stafylakis, T., & Tzimiropoulos, G. (2017). Combining residual networks with LSTMs for lipreading. *Interspeech*.
- [8] Chung, J. S., & Zisserman, A. (2016). Lip reading in the wild. *Asian Conference on Computer Vision (ACCV)*.
- [9] Zhao, G., & Pietikäinen, M. (2007). Dynamic texture recognition using local binary patterns with an application to facial expressions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(6), 915-928.
- [10] Graves, A., Fernández, S., Gomez, F., & Schmidhuber, J. (2006). Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks. *ICML*.
- [11] Kingma, D. P., & Ba, J. (2015). Adam: A Method for Stochastic Optimization. *International Conference on Learning Representations (ICLR)*.
- [12] Simonyan, K., & Zisserman, A. (2014). Two-stream convolutional networks for action recognition in videos. *NeurIPS*.
- [13] Saitoh, T., & Okuno, H. G. (2000). Smart lip-reading system for speech recognition. *IEEE International Conference on Systems, Man, and Cybernetics*.
- [14] Chung, J. S., & Zisserman, A. (2017). Out of time: Automated lip sync in the wild. *Workshop on Multimodal Understanding of Social, Affective and Subjective Attributes (CVPR)*.
- [15] Petridis, S., Stafylakis, T., Ma, P., Cai, J., & Pantic, M. (2018). End-to-end multi-view lipreading. *British Machine Vision Conference (BMVC)*.
- [16] Noda, K., Yamaguchi, Y., Nakadai, K., Okuno, H. G., & Ogata, T. (2015). Audio-visual speech recognition using deep learning. *Applied Intelligence*, 42(4), 722-737.
- [17] Huang, J., & Kingsbury, B. (2013). Audio-visual deep learning for noise robust speech recognition. *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*.
- [18] Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep Learning. *MIT Press*.

- [19] Chollet, F. (2015). Keras.
<https://github.com/fchollet/keras>
- [20] Paszke, A., Gross, S., Massa, F., et al. (2019).
PyTorch: An imperative style, high-
performance deep learning library. *NeurIPS*.