

LipReaderAI: Advanced Multi-Speaker Lipreading for Enhanced Communication Accessibility

Rohan Nain, Shivam

Department of Physics & Astronomy,

The University of Tennessee

March 14, 2024

Abstract

LipReaderAI presents a novel approach in assistive technologies, designed to overcome communication barriers for individuals with hearing impairments. Our model leverages the synergy between spatiotemporal convolutions and recurrent neural network architectures to accurately transcribe spoken language from visual input across multiple speakers. To eliminate the need for manual segmentation, we employ Connectionist Temporal Classification (CTC) loss, which facilitates direct, end-to-end training from unprocessed video input to textual transcription. This project tests the robustness of the model by evaluating its accuracy on a custom dataset featuring recordings of the author, ensuring it can reliably translate visual speech. LipReaderAI not only aims to improve accessibility for those with hearing impairments but also extends the use of lipreading into loud environments where traditional speech recognition systems that rely solely on audio are ineffective.

Contents

1	Introduction	3
1.1	Background:	3
1.2	Motivation: Related Work	3
2	Data Set	4
2.1	Data Description	4
2.1.1	Learning Objectives	5
2.1.2	Data set Characteristics	5
2.1.3	Source and Methodology	5
3	Model and Implementation	6
3.1	Model Architecture	6
3.2	Loss Function	7
3.3	Training	8
4	Model Evaluation	9
4.1	Bench marking Methods	10
4.2	Optimization Techniques	10
4.3	Testing and Validation	10
5	Results	11
5.1	Performance Metrics	11
6	Discussion	13
6.1	Interpretation of the Results	13
6.2	Limitations and Challenges	13
6.3	Future Directions	13
7	Conclusion	14
8	Appendix	14
8.1	Additional Resources	14
	References	15

1 Introduction

1.1 Background:

Speech recognition technology has evolved significantly and has come a long way, tapping into a myriad of applications in various fields extending from accessibility, security, and human-computer interactions. Traditional speech recognition systems heavily rely on the audio input which gets highly impacted in a noisy environment and complete silence is required for complete accuracy. However, where the audio signal falters, visual speech recognition (such as LipReaderAI) offers a very promising substitute by interpreting speech based only on visual cues.

Lip-reading technology has traditionally been very challenging because of the subtle variations in the lip movements and, very high precision is required for capturing these movements. The rise of the advanced machine learning techniques, particularly deep learning, has opened new doors for improving the accuracy and efficiency of the lip reading systems. These systems can now leverage vast amounts of data to learn the intricate lip movements and associate them with the corresponding speech.

1.2 Motivation: Related Work

The primary motivation for LipReaderAI project is to enhance the usability of the speech recognition systems to noisy environments particularly for the people with the hearing impairments. As communication can be severely hampered by hearing difficulties, while sign language and hearing aids can help, a computational model that can read lips automatically adds another level of accessibility and involvement.

Several works laid the foundation of applying the deep learning to lipreading, notably the LipNet model developed by researchers at the University of Oxford (Assael, Shillingford, Whiteson, & De Freitas, 2016). This LipNet paper has been the fundamental motivation for this research, as this study showed promising results by using spatiotemporal convolutions to capture the dynamics of the lip movements. Later on, to improve the temporal resolutions of recorded lip movements, more intricate designs involving 3D convolutions and recurrent neural networks have been put forth. However, the majority of models already in use are designed for single-speaker scenarios and are frequently tested and trained in controlled environments.

This project aims to reproduce the results of previous work on lip-reading by training a deep neural network (DNN) model on data from 10 specific speakers, rather than attempting

to handle multiple arbitrary speakers as stated in the original work. The use of Connectionist Temporal Classification (CTC) loss still allows for training the model end-to-end without the need for segmenting the video into predefined units, which remains a significant advancement over traditional methods requiring manual segmentation.

Although reproducing state-of-the-art findings on a targeted 10-speaker dataset is more constrained than handling totally arbitrary speakers in real-time, it nevertheless pushes the limits of what is currently possible with deep learning for lip-reading. This focused approach also offers a model for gradual advancement and insights that can direct future research efforts to improve the robustness of visual speech recognition across many circumstances and speakers.

2 Data Set

The dataset used in this project is based on the GRID audiovisual sentence corpus, as utilized by Assael et al. in their groundbreaking research on the LipNet model (Assael et al., 2016).

2.1 Data Description

This dataset comprises video files in MPEG format, accompanied by alignment files that serve as labels. These labels map the spoken content to specific video timestamps to facilitate precise synchronization between the visual and textual data. An example of the alignment data is as follows:

```
0 23750 sil
23750 29500 bin
29500 34000 blue
34000 35500 at
35500 41000 f
41000 47250 two
47250 53000 now
53000 74500 sil
```

In this format, 'sil' denotes silence, indicating periods where no speech occurs. Such detailed labeling is crucial for training accurate lipreading models that can distinguish spoken words from non-verbal periods.



Figure 1: Example frames from videos of different speakers, showing the diversity in lip shapes and movements.

2.1.1 Learning Objectives

The primary objective is to develop LipReaderAI, a lipreading model that relies solely on visual inputs to capture lip movements and generate corresponding textual speech outputs. This model aims to surpass traditional speech-to-text models by functioning effectively even in the absence of audio input, which is particularly useful in noisy environments or in situations where sound is unavailable or obscured.

2.1.2 Data set Characteristics

The GRID corpus is extensive, containing video recordings from 33 different speakers. For the initial phase of our project, we focus on a subset of this data, selecting 1001 videos from each of five distinct speakers. This approach provides a balanced and diverse training set of 10010 videos in total, each paired with their corresponding textual alignments. This diversity is critical for ensuring that LipReaderAI can generalize well across different speakers and articulations.

2.1.3 Source and Methodology

The dataset was sourced from the publicly available GRID audiovisual sentence corpus, accessible at <https://spandh.dcs.shef.ac.uk/gridcorpus/>. The preparation of this dataset involved extensive manipulation and preprocessing using the OpenCV library for video processing and the TensorFlow library for data structuring and neural network training. These tools were instrumental in converting raw video files into a format suitable for training our deep learning model, ensuring that each video frame was properly aligned with its corresponding label for

effective learning.

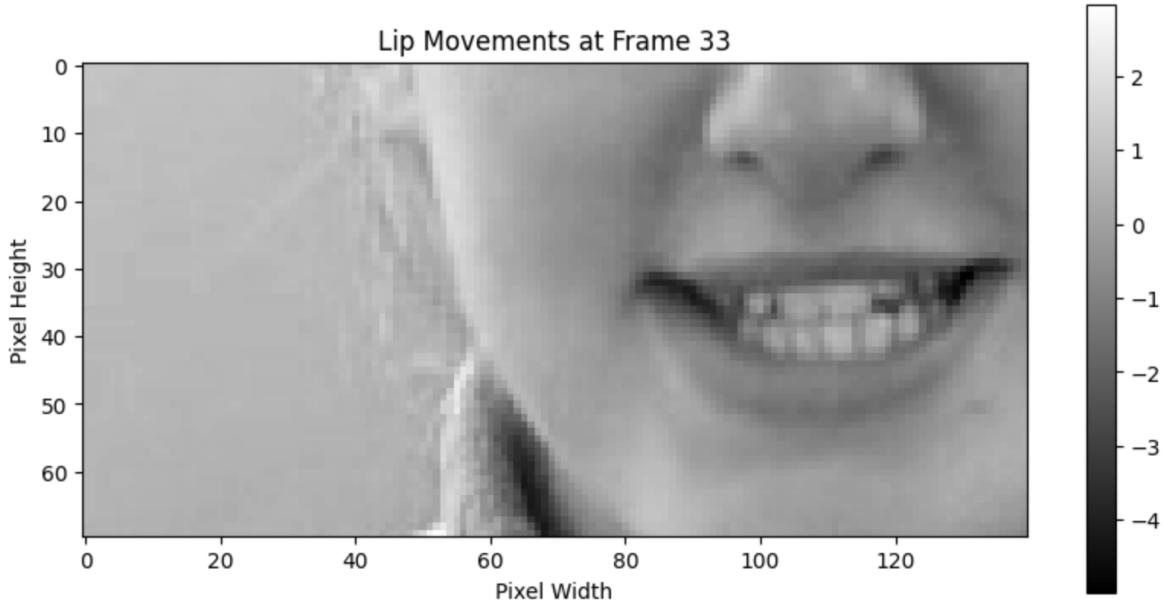


Figure 2: Lip movements at frame 33.

3 Model and Implementation

In this section, we describe the DNN model architecture of the LipReaderAI model, the specialized loss function employed, and the training procedures implemented to optimize model performance.

3.1 Model Architecture

The LipReaderAI model employs a combination of 3D convolutional neural networks (CNNs) and bidirectional Long Short-Term Memory (LSTM) networks to process video input and produce accurate lipreading results. The architecture is designed to capture both spatial and temporal information effectively.

- **3D Convolutional Layers:** The initial part of the network consists of three 3D convolutional layers. These layers are responsible for extracting spatiotemporal features from the video frames. Each layer is followed by a ReLU activation function to introduce non-linearity, helping the network learn more complex patterns in the data.
 1. The first layer applies 128 filters of size $(3, 3, 3)$, with padding to maintain the spatial dimensions, on the input video sequence shaped as $(75, 70, 140, 1)$, where 75

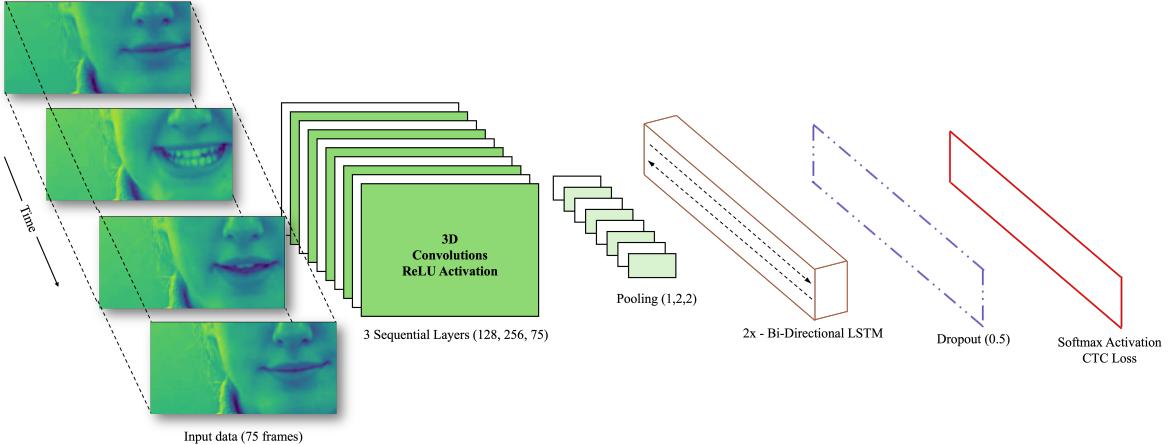


Figure 3: Diagram of the LipReaderAI model architecture.

represents the time dimension (or frames), and 70x140 is the spatial resolution per frame.

2. The second and third layers increase the number of filters to 256 and then to 75, respectively, with each also using a (3, 3, 3) kernel size. These layers continue to refine the features extracted, focusing on higher-order relationships within the data.

Each convolutional layer is followed by a max pooling layer with a (1, 2, 2) pool size to reduce dimensionality and increase the field of view of the filters.

- **Bidirectional LSTM Layers:** After the convolutional layers, the model uses two bidirectional LSTM layers, each with 128 units. These layers are designed to capture temporal dependencies and dynamics from the sequence of features extracted by the convolutional layers. The bidirectional approach processes the data in both forward and reverse temporal directions, enhancing the context available to the network.
- **Output Layer:** The final part of the network is a dense layer with a softmax activation function, designed to output the probabilities of each character in the lipreading vocabulary at each time step.

3.2 Loss Function

The model uses the Connectionist Temporal Classification (CTC) loss to train the network. CTC is particularly well-suited for tasks like speech and lipreading, where the alignment between the input sequences and the target text is not explicitly known. The CTC loss function enables the model to align the sequence of predicted probabilities with the actual labels efficiently,

optimizing the training process by automatically finding the most likely alignment between inputs and labels.

This loss function is directly copied from the CTCLoss defined for the Automatic Speech Recognition using CTC by Authors: Mohamed Reda Bouadjenek and Ngoc Dung Huynh

```
def CTCLoss(y_true, y_pred):
    batch_len = tf.cast(tf.shape(y_true)[0], dtype="int64")
    input_length = tf.cast(tf.shape(y_pred)[1], dtype="int64")
    label_length = tf.cast(tf.shape(y_true)[1], dtype="int64")

    input_length = input_length * tf.ones(shape=(batch_len, 1), dtype="int64")
    label_length = label_length * tf.ones(shape=(batch_len, 1), dtype="int64")

    return tf.keras.backend.ctc_batch_cost(y_true, y_pred, input_length, label_length)
```

3.3 Training

Training the LipReaderAI model involves several steps designed to maximize the accuracy and efficiency of the model:

Hyperparameters	Values
Checkpoint Directory	Saved_LipNet_model_10_speakers/checkpoints
Checkpoint monitor	loss
Save Weights Only	True
Save Best Only	True
Verbose Output	1
Latest Checkpoint	Loaded if exists
Epochs	30
Callbacks	Checkpoint_callback, schedule_callback
Optimizer	Adam
Learning Rate	0.0001
Loss Function	CTCLoss
Learning Rate Schedule	Exponential decay after epoch 15

Figure 4: Table Summarizing all the Hyper parameters for Training the Model.

- **Data Preparation:** The video and label data are preprocessed to ensure that each video frame is synchronized with its corresponding label. This involves aligning the video frames with the time stamps provided in the data set’s alignment files.
- **Model Compilation:** The model is compiled using the Adam optimizer, a popular choice for deep learning tasks due to its efficiency in handling sparse gradients and its adaptive learning rate capabilities.
- **Model Fitting:** The model is trained on the prepared data set, using the CTC loss function to optimize the alignment of predicted text with the ground truth. Training parameters such as batch size, number of epochs, and validation splits are chosen based on preliminary experiments to balance training speed and model accuracy.

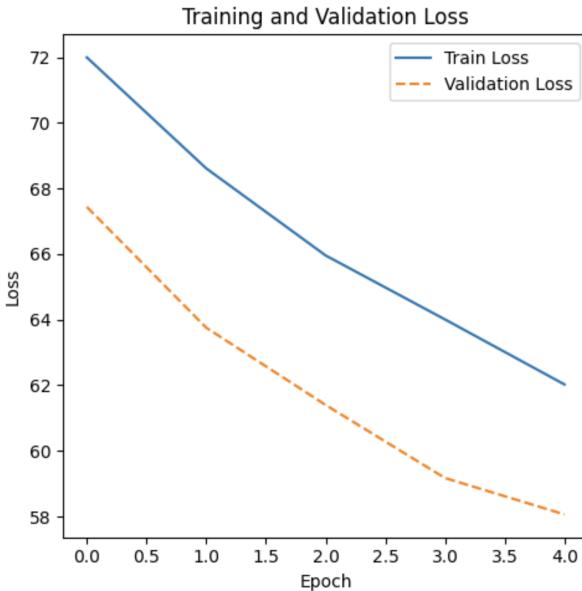


Figure 5: Training loss variation over four epochs. The plot illustrates the decreasing trend in losses. Only four epochs were recorded due to kernel interruptions during training.

4 Model Evaluation

This section explains the methodologies employed to assess the performance and reliability of the LipReaderAI model. We discuss benchmarking methods, optimization techniques used during training, and detailed testing and validation procedures.

4.1 Bench marking Methods

Bench marking for LipReaderAI involves comparison with baseline models and previously established standards in lipreading technology. Performance metrics such as accuracy, precision, and computational efficiency were compared.

- **Model Checkpointing:** Checkpoints were used to save and retrieve the best model configuration during training, ensuring that only the most effective model is used for bench marking. The ModelCheckpoint callback was configured to monitor the loss, saving weights of the model showing the lowest loss.
- **Performance Metrics:** We used accuracy and loss metrics to evaluate and benchmark the model performance during training and testing phases.

4.2 Optimization Techniques

The training process incorporated several optimization techniques to enhance model performance, focusing on both accuracy and efficiency.

- **Checkpointing and Model Loading:** The training setup included saving model checkpoints only when there was an improvement in loss, ensuring efficient use of storage and quick recovery of the best model state for further training or evaluation. The latest checkpoints were loaded at the beginning of the model training sessions to resume training from the last best state.
- **Verbose Output:** Verbose settings in callbacks provided detailed logs during training, helping to monitor the training progress and debug issues effectively.

4.3 Testing and Validation

Testing and validation were critical to ensure the model's robustness and reliability across different scenarios and datasets.

- **Validation Split:** We employed a training-validation split, where the model was trained on a subset of the data and validated on a separate set to monitor and prevent over fitting. Thus, we have trained the model on all the 10010 videos with 75 frames each and the corresponding alignments by making a split of 80/20 training and testing split.

- **Custom Data set Testing:** The model’s generalization capabilities were tested using custom data sets, including videos that were not part of the training data set. This testing was crucial for evaluating the models performance in real-world scenarios.
- **Decoding and Text Prediction:** After making predictions, CTC decoding was used to translate the output probabilities into textual predictions. This step was essential for comparing the predicted text against actual labels to assess the accuracy.
- **Similarity Calculation:** Custom similarity scores were calculated to quantitatively assess how closely the predicted text matched the actual text. This method provided a direct measure of the transcription accuracy.

```
1/1 [=====] - 0s 62ms/step
Actual Text: bin white at x seven again
Predicted Text: bin white at k seven again
Custom Similarity Score: 0.9629629629629629
```

Figure 6: Custom Dataset Testing

- **Average Similarity Score:** Across multiple test videos, the average similarity score was computed to provide an overall measure of the model’s effectiveness on unseen data. For our model, it comes out to be 0.78 when tested on a unknown speaker with 1001 videos.

These methods, implemented as described, ensured that LipReaderAI was rigorously tested and optimized for both performance and reliability, addressing the challenges of real-time lipreading with high accuracy and efficiency.

5 Results

This section details the outcomes of the tests performed on the LipReaderAI model using various metrics to evaluate its performance. The results include direct comparisons of predicted text against actual text, as well as computed similarity scores which provide a quantitative measure of the model’s accuracy.

5.1 Performance Metrics

The performance of the LipReaderAI model was evaluated through several tests on a custom data set not seen during the training phase. These tests are crucial for assessing the real-world applicability of the model in accurately transcribing spoken words from lip movements.



1/1 [=====] - 2s 2s/step
Actual Text: bin white by n three soon
Predicted Text: lay white by s tw soon
Custom Similarity Score: 0.5

(a) Validation on Speaker 01



1/1 [=====] - 0s 62ms/step
Actual Text: bin white at x seven again
Predicted Text: bin white at k seven again
Custom Similarity Score: 0.9629629629629629

(b) Validation on Speaker 23

Figure 7: Validation Results

- **Individual Test Example:** In one of the test runs, the model was tasked with predicting the text from a video where the actual spoken phrase was "bin white at x seven again". The model predicted "bin white at k seven again". This resulted in a Custom Similarity Score of approximately 0.9629. This high score indicates that the model can accurately capture most of the phonetic elements from the visual input, though it sometimes confuses similar-sounding words or letters.
- **Average Performance Across Tests:** The overall performance of the model was further quantified by calculating the average similarity score across a series of tests conducted on different videos from the custom dataset. The average similarity score obtained was about 0.78. This lower than expected average score may be due to the limited diversity in the training dataset, as the model was initially trained on videos from only 10 speakers. This restricted training set may not have provided sufficient variability to effectively differentiate between a broader range of speakers in the custom dataset.
- **Implications and Future Training Plans:** The limited diversity in the training data has been identified as a key area for improvement. The initial training on just 5 speakers may have constrained the model's ability to generalize effectively across new, unseen speakers. To enhance performance and robustness, the dataset will be expanded to include a larger and more diverse array of speakers. This expansion is expected to provide the model with a broader phonetic and visual range to learn from, potentially increasing its accuracy and performance on custom datasets.

These results underscore the importance of continuous refinement and testing of the LipReaderAI model to ensure its robustness and reliability across a wide range of real-world applications. Further investigations into the causes of variability in performance scores will be crucial for targeted improvements in the model architecture and training processes.

6 Discussion

6.1 Interpretation of the Results

The results from the testing and evaluation of the LipReaderAI model show a promising ability to accurately interpret lip movements into text. The high similarity scores in certain tests demonstrate the model's potential in environments where audio cannot be relied upon. However, the variation in performance across different tests, particularly the lower average similarity score, highlights the complexity of visual speech recognition and the influence of diverse variables such as speaker differences, lighting conditions, and lip movement idiosyncrasies.

6.2 Limitations and Challenges

One of the primary limitations encountered in this project was the relatively small and homogeneous training dataset. Training the model on just ten speakers limited its ability to generalize to new speakers, which was evident from the performance on the custom dataset. Furthermore, visual speech recognition inherently deals with subtle nuances in lip movements, which can vary drastically not only across different speakers but also due to the speaker's emotions, speech pace, and environmental factors.

6.3 Future Directions

To overcome the limitations and enhance the model's robustness and accuracy, several steps are proposed:

- **Expanding the Dataset:** Increasing the number of speakers in the training dataset and incorporating more diverse speaking styles, accents, and non-standard speech patterns will help improve the model's generalization capabilities.
- **Incorporating Multimodal Data:** Future versions of the model could integrate audio data where available, adopting a multimodal approach to improve accuracy and robustness under varied conditions.

- **Advanced Modeling Techniques:** Exploring more complex neural network architectures or newer technologies like transformers in the visual speech recognition domain could potentially yield better results.
- **Real-World Application Testing:** Extending testing to more practical applications, such as in noisy environments or real-time communication systems, to assess and refine the model’s performance in real-world scenarios.

7 Conclusion

The LipReaderAI project embarked on developing a model capable of translating lip movements into textual speech with the aim of enhancing communication accessibility for individuals with hearing impairments and in environments where audio is not feasible. Despite the challenges encountered, such as limitations in training data diversity and intrinsic challenges of visual speech recognition, the model demonstrated a substantial potential to fulfill its intended purpose. The insights gained from the extensive testing and evaluation phases have laid a solid foundation for future enhancements. Expanding the training dataset and integrating advanced modeling techniques are identified as critical next steps to achieving a more robust and universally applicable system. This work not only contributes to the field of assistive technologies but also opens avenues for further research in speech recognition and artificial intelligence.

8 Appendix

8.1 Additional Resources

During the development of the LipReaderAI project, numerous resources were consulted to enhance the understanding and implementation of the model. The following are some of the GitHub repositories and research papers that provided significant insights and were instrumental in the development of this project:

- Assael, Y., et al. "LipNet: End-to-End Sentence-level Lipreading." 2016 (Assael et al., 2016).
- Ting, Y., et al. "Comprehensive Review on Lip Reading Techniques." 2022 (Ting, Song, Huang, & Tian, 2022).

- Son, J., et al. "Lip Reading: A New Leap in Speech Recognition." 2017 (Son Chung, Senior, Vinyals, & Zisserman, 2017).
- Pu, Y., et al. "Review on Automatic Lip-Reading in the Wild." 2023 (Pu & Wang, 2023).
- Gutierrez, O., et al. "Lip Reading and Recognition for Improved Speech Systems." 2017 (Gutierrez & Robert, 2017).
- Amitabha, A., et al. "Lip2Text: Harnessing Deep Neural Networks for Mobile Lip-Reading." 2024 (Amitabha et al., 2024).
- Mestri, S., et al. "Analysis of Lip Movement for Speech Recognition and its Applications." 2019 (Mestri, Limaye, Khuteta, & Bansode, 2019).
- Owens, E., et al. "Visually-Aided Speech Enhancement." 2016 (Owens et al., 2016).
- Exarchos, M., et al. "Deep Learning Models for Lip-Reading." 2024 (Exarchos et al., 2024).
- Lu, D., et al. "Automatic Lip Reading Based on Spatial-Temporal Visual Analysis." 2019 (Lu & Li, 2019).
- Kastaniotis, I., et al. "Lip Reading Through Deep Learning: A Review." 2020 (Kastaniotis, Tsourounis, & Fotopoulos, 2020).

These references provide a broad spectrum of methodologies and advancements in the field of lip-reading and speech recognition technologies, contributing significantly to the theoretical and practical development of LipReaderAI.

References

- Amitabha, A. D., et al. (2024). Lip2text: Sentence-level lipreading on english speakers using the deep learning approach.
- Assael, Y. M., Shillingford, B., Whiteson, S., & De Freitas, N. (2016). Lipnet: End-to-end sentence-level lipreading. *arXiv preprint arXiv:1611.01599*.
- Exarchos, T., Dimitrakopoulos, G. N., Vrahatis, A. G., Chrysovitsiotis, G., Zachou, Z., & Kyrodimos, E. (2024). Lip-reading advancements: A 3d convolutional neural network/long

- short-term memory fusion for precise word recognition. *BioMedInformatics*, 4(1), 410–422.
- Gutierrez, A., & Robert, Z. (2017). Lip reading word classification. *Comput Vision-ACCV*.
- Kastaniotis, D., Tsourounis, D., & Fotopoulos, S. (2020). Lip reading modeling with temporal convolutional networks for medical support applications. In *2020 13th international congress on image and signal processing, biomedical engineering and informatics (cisp-bmei)* (pp. 366–371).
- Lu, Y., & Li, H. (2019). Automatic lip-reading system based on deep convolutional neural network and attention-based long short-term memory. *Applied Sciences*, 9(8), 1599.
- Mestri, R., Limaye, P., Khuteta, S., & Bansode, M. (2019). Analysis of feature extraction and classification models for lip-reading. In *2019 3rd international conference on trends in electronics and informatics (icoei)* (pp. 911–915).
- Owens, A., Isola, P., McDermott, J., Torralba, A., Adelson, E. H., & Freeman, W. T. (2016). Visually indicated sounds. In *Proceedings of the ieee conference on computer vision and pattern recognition* (pp. 2405–2413).
- Pu, G., & Wang, H. (2023). Review on research progress of machine lip reading. *The Visual Computer*, 39(7), 3041–3057.
- Son Chung, J., Senior, A., Vinyals, O., & Zisserman, A. (2017). Lip reading sentences in the wild. In *Proceedings of the ieee conference on computer vision and pattern recognition* (pp. 6447–6456).
- Ting, J., Song, C., Huang, H., & Tian, T. (2022). A comprehensive dataset for machine-learning-based lip-reading algorithm. *Procedia Computer Science*, 199, 1444–1449.