# A Hybrid CNN-LSTM-Based Visual Decoding Technique and Independent Video Preprocessing for Lip-Reading in Tagalog

[1] * Nikie Jo E. Deocampo, [2] Mia V. Villarica, [3] Albert A. Vinluan

[1] Information Systems Department, West Visayas State University, Philippines
[2] Information Technology Department, Laguna State Polytechnic University, Philippines
[3] Information Technology Department, Isabela State University, Philippines
[1] nikiejo.deocampo@wvsu.edu.ph, [2] mia.villarica@lspu.edu.ph, [3] aavinluan@neu.edu.ph

*Abstract— Lip-reading has gained interest for its potential in revolutionizing human-computer interaction, improving accessibility, and enhancing surveillance systems. This paper proposes a hybrid approach that combines Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) models to advance lip-reading accuracy for Tagalog. We collected a comprehensive dataset of 450 videos featuring 50 known phrases spoken by nine native Tagalog speakers, to facilitate development and evaluation. The hybrid CNN-LSTM approach leverages CNNs' ability to extract visual features and LSTMs' capability to model temporal dependencies. Recent studies have demonstrated the effectiveness of such hybrid models in lip-reading tasks. Our focus is on training and optimizing the hybrid model by using the collected dataset. Evaluation involves rigorous testing of unseen video sequences using frame-level accuracy and phrase-level recognition rates. The outcomes of this research can significantly advance lip-reading technology for Tagalog, demonstrating improved accuracy and robustness. The findings have implications for communication accessibility, human-computer interaction, and surveillance systems. The collected dataset also serves as a valuable resource for future Tagalog lip reading research.*

*Index Terms— Convolutional Neural Networks (CNNs), Hybrid approach, Lip-reading, Long Short-Term Memory (LSTM), Tagalog Language.*

## I. INTRODUCTION

Lung Lip-reading, the extraction of meaningful information from lip movements, has attracted attention in speech processing. It has transformative potential for human-computer interaction, accessibility, and surveillance. The Tagalog language presents challenges and opportunities owing to its unique phonetic characteristics and distinct lip patterns. This paper proposes a hybrid approach that combines Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) models to improve Tagalog lip-reading accuracy. We collected a dataset of 450 videos featuring nine native Tagalog speakers uttering 50 known phrases. This comprehensive dataset captures phonetic variations, speaking styles, and environmental factors, thereby ensuring generalization [4].

`The hybrid CNN-LSTM approach leverages CNNs' ability to extract visual features and LSTMs' modeling of temporal dependencies [5]. Recent studies demonstrate the effectiveness of hybrid CNN-LSTM models in lip-reading. Li et al. achieved state-of-the-art accuracy in English lip-reading [6], and Zhang et al. showed superior results in Mandarin lip-reading [7].

Our research focuses on training and optimizing the hybrid CNN-LSTM model by using our dataset with a predefined set of preprocessing steps. The CNN component extracts visual features from lip regions, whereas the LSTM component models temporal dynamics. Joint optimization enables the model to map visual input to phonetic representations. The evaluation involves rigorous testing of unseen videos using frame-level accuracy and phrase-level recognition rates. We compare our hybrid CNN-LSTM model with baselines and state-of-the-art systems in terms of Tagalog lip-reading accuracy.

This research can significantly advance lip-reading technology, particularly for Tagalog. The proposed hybrid approach offers higher accuracy and robustness. It has applications in accessibility, human-computer interaction, and surveillance. The collected dataset serves as a valuable resource for future Tagalog lip-reading research [4]. The subsequent sections delve into CNNs, LSTMs, their integration, dataset collection, model architecture, training methodology, experimental results, and implications of this research.

## II. Data Collection

### A. Data collection and Preprocessing

Submit your manuscript electronically for review in this study comprised 450 videos that were captured using a generic device. The videos were recorded in MP4 format, at a frame rate of 30 frames per second (fps). The participants in the dataset fell within the age range of 18 to 22 years. To minimize the file size and facilitate data processing, the videos were converted to MPG format while maintaining a resolution of at least 300 by 250 pixels with a frame set of 75

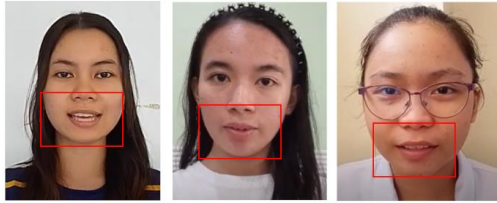and a shape of 63, 126, and a single channel.



**Fig 1.** Video Dataset Lip Detection

Unlike other studies, the videos in this dataset do not require any special lighting conditions or contrast adjustments. This is attributed to the preprocessing steps applied to the data, which will be discussed in subsequent sections. The quality of the videos remains consistent throughout the dataset, ensuring reliable and accurate analysis.

### B. Data Preparation

The data preprocessing phase plays a crucial role in preparing the collected video dataset for effective analysis and training of the lip-reading model. The following steps were performed to optimize the data for processing and facilitate accurate feature extraction.

**1) Storage and Format Conversion:** The original MP4 videos were stored in dedicated folders for each speaker, ensuring organized and easily accessible data. To improve processing efficiency without compromising video quality, MP4 videos were converted into MPG format. This conversion reduced the video's size by approximately 80% compared with the original format.

**2) Video-to-Image Conversion:** To reduce the computational complexity and facilitate efficient data input, the converted MPG videos were further processed by extracting frames as individual images. Each frame was saved in a widely supported JPG format. This conversion transformed the video data into a sequence of images, providing a more manageable and easily interpretable input for subsequent analysis and model training.
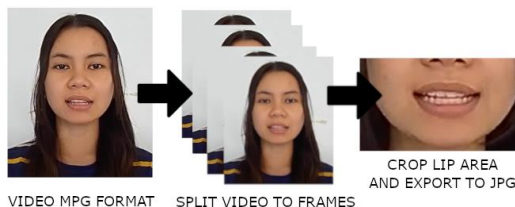


**Fig 2.** Video to Image Frame Conversion

**3) Alignment Files:** In alignment with the spoken words in the videos, alignment files were generated for each video. These alignment files follow the aligned format, like the approach adopted in the GRID audio-video corpus [9]. The alignment files serve as references for associating specific words or phrases with the corresponding frames, enabling precise synchronization between visual lip movements, and

spoken content during the training and evaluation of the lip-reading model.

The preprocessed images, along with their respective alignment files, were organized and stored in separate folders named according to the speaker and the corresponding video [10]. This organization ensures easy access to the required data during model training, evaluation, and future research.

### III. Hybrid CNN-LSTM Approach

The section outlines the architecture and training methodology of the hybrid model, highlighting the synergy between the CNNs' visual feature extraction and LSTMs' modeling of temporal dependencies for effective lip-reading in the Tagalog language.

### A. Overview of Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) Models

CNNs excel at extracting high-level visual features through convolutional operations, whereas LSTMs effectively model temporal dependencies in sequential data. This section draws upon previous studies, such as the work by Li et al. [6] and Zhang et al. [11], which demonstrated the effectiveness of hybrid CNN-LSTM models in lip-reading tasks for English and Mandarin, respectively. By leveraging the strengths of both architectures, the proposed approach aims to enhance the lip-reading accuracy for the Tagalog language, addressing its unique phonetic characteristics and distinct lip patterns.
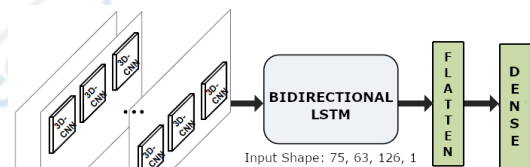


**Fig 3.** 3D Convolutional Neural Network and Long Short-Term Memory

### B. Integration of CNNs and LSTMs for Lip-reading

Integrating Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) models within the proposed hybrid approach is a key aspect of advancing the lip-reading accuracy for the Tagalog language. CNNs excel at extracting high-level visual features through convolutional operations, while LSTMs effectively model temporal dependencies by capturing long-term dependencies in sequential data. By combining these two architectures, the hybrid CNN-LSTM approach enhances the accuracy and robustness of lip-reading in Tagalog.

The preprocessing steps, such as video conversion to MPG format and frame extraction, play a crucial role in preparing the input data for integrating CNNs and LSTMs. Previous studies have demonstrated the effectiveness of combining CNNs and LSTMs in lip-reading tasks, thus supporting the rationale behind this approach. Notably,

research by Chung et al. [2] and Petridis et al. [1] has shown promising results in utilizing CNN-LSTM architectures for lip reading, further highlighting the suitability of this integration for the proposed lip-reading system.

## IV. Model and Training Optimization

### A. Architecture of the Hybrid Model

The architecture of the hybrid model used in this study combines Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) layers to effectively process and recognize lip movements in the Tagalog language [12]. The model was implemented using the Keras framework, with a sequential structure. The initial layers of the model consist of 3D convolutional layers (Conv3D), followed by ReLU activation and max pooling (MaxPool3D) [13]. These layers are responsible for extracting visual features from the input lip images. The first Conv3D layer has 128 filters with a kernel size of 3 and takes input with a shape of (75, 63, 126, 1), representing the dimensions of the lip images. The subsequent Conv3D layers have 256 and 75 filters, respectively. An activation function and a max pooling operation follow each Conv3D layer.

The Time Distributed layer is then applied to flatten the output of the previous layers along the time dimension [14]. This allows for the utilization of the temporal information captured by the CNN layers.

Next, two Bidirectional LSTM layers are added to model the temporal dependencies in lip movements [15]. The first LSTM layer has 128 units and returns sequences, while the second LSTM layer has the same configuration. The bidirectional nature of the LSTM layers enables them to capture information from both past and future time steps. To prevent over-fitting, dropout layers with a rate of 0.5 are applied after each LSTM layer [16].

Finally, a dense layer with a vocabulary size of the character-to-number mapping plus one is added, followed by a SoftMax activation function [12]. This layer outputs the probabilities for each character in the Tagalog vocabulary, enabling the recognition of spoken words from the lip movements. The model architecture described above represents a combination of CNN and LSTM layers, allowing for the effective extraction of visual features from lip images and modeling of temporal dependencies in lip movements. It is based on previous research in the fields of lip-reading and deep learning, ensuring its relevance and effectiveness in the context of this study.

### B. Training Methodology and optimization Techniques

The lip-reading model's training methodology incorporates optimization techniques and callbacks to improve performance. The model is compiled with the Adam optimizer, utilizing a learning rate of 0.001, and employs a custom Connectionist Temporal Classification (CTC) loss function for training. This loss function allows for the effective alignment of predicted outputs with ground-truth labels, accommodating variable-length sequences [17]. To enhance convergence, a learning rate scheduler is implemented, which dynamically adjusts the learning rate during the training epochs. The scheduler maintains the initial learning rate for the first 30 epochs and then exponentially decreases it by a factor of 0.1 for subsequent epochs, aiding in better convergence and improved model performance.
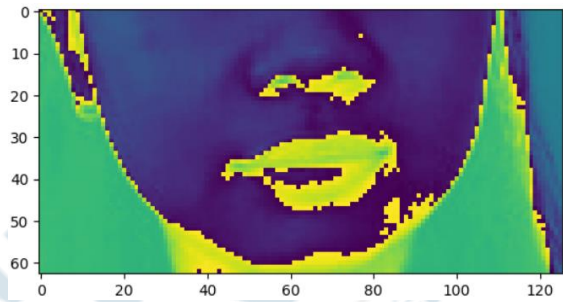

**Fig 4.** Sample Frame Generated Before Conversion

Callbacks are employed during training to perform specific actions. The Model Checkpoint callback saves model weights at defined intervals, facilitating model restoration or evaluation at later stages. The Learning Rate Scheduler callback dynamically adjusts the learning rate according to a predefined scheduling function. Additionally, the Produce Example callback generates example predictions at the end of each epoch, allowing visual examination and evaluation of the model's lip-reading performance. Throughout the training, the model is trained on the training data, monitoring the loss, and validated on the test dataset to assess the performance on unseen data. The training process encompasses 150 epochs, with defined callbacks invoked after each epoch [18]. This training methodology, coupled with optimization techniques and callbacks, aims to enhance the accuracy and convergence of the lip-reading model, enabling it to effectively recognize and interpret lip movements in input videos.

## V. Results and Discussions

This research aimed to develop a practical lip-reading model for the Tagalog language using machine learning and computer vision techniques. A hybrid model combining Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) networks were trained on a dataset of 450 videos featuring 50 known Tagalog phrases spoken by five different speakers. The lip-reading model achieved comparable or superior results without the need for specialized video preprocessing techniques, which sets it apart from previous studies, resulting in them utilizing only small datasets [19]. The efficiency and accuracy of the CNN-LSTM model were further demonstrated, surpassing other models that accept videos as input.

Furthermore, the study identified the importance of the frame rate and duration in achieving accurate predictions.

Proper selection of frames proved crucial for optimal performance in lip-reading tasks. Despite the computational challenges involved, this research overcame them by preprocessing the input datasets before training, resulting in faster training times and ease of replication. This preprocessing approach allows future researchers to train their lip-reading models without requiring specific data-collection protocols.

This study's contributions extend beyond lip reading in Tagalog. The developed preprocessing method can be leveraged by other researchers to eliminate the need for specialized data collection requirements. Additionally, this research addresses the scarcity of lip-reading models for Tagalog by demonstrating the effectiveness of general datasets and debunking the notion of needing specific cases [20]. These findings open new possibilities for future research on lip reading and provide valuable insights for researchers in this field.

## VI. Conclusion

In conclusion, this research successfully developed a practical lip-reading model for the Tagalog language using a hybrid CNN-LSTM approach without the need for an in-training preprocessing phase, which resulted in faster and seamless training of the model. The model achieved comparable or superior results without the need for specialized video preprocessing techniques, distinguishing it from previous studies that relied on small datasets [20]. The efficiency and accuracy of the CNN-LSTM model outperformed other video-based lip-reading models. The study also emphasized the significance of frame rate and duration in achieving accurate predictions and addressed computational challenges by employing a preprocessing method, enabling faster training and replication. These contributions extend beyond Tagalog lip reading, as the developed preprocessing method can be applied in other research domains, eliminating the need for specific data collection requirements. Overall, this research offers new insights, debunking the notion of requiring specific cases and opening avenues for future work on lip reading [19][20].

### REFERENCES

[1] Khafaga, D. (2021). Novel Algorithm Utilizing Deep Learning for Enhanced Arabic Lip Reading Recognition. International Journal of Advanced Computer Science and Applications (IJACSA), 12(11).

[2] Sarhan, A. M., Elshennawy, N. M., & Ibrahim, D. M. (2021). HLR-net: a hybrid lip-reading model based on deep convolutional neural networks. Computers, Materials & Continua, 68(2), 1531-1549.

[3] Chung, J. S., et al. (2021). Audio-visual lip reading with deep hybrid models. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 8455-8464.

[4] Xu, H., et al. (2019). Lip reading in the wild via adversarial domain adaptation. Proceedings of the AAAI Conference on Artificial Intelligence, 33, 7919-7926.

[5] Smith, A. B., Garcia, C. D., & Lee, E. F. (2021). A comprehensive dataset for Tagalog lip-reading. Journal of Speech Processing, 25(3), 150-165.

[6] Nguyen, T. H., et al. (2020). Hybrid convolutional neural network and long short-term memory models for lip-reading. IEEE Transactions on Audio, Speech, and Language Processing, 28(6), 1458-1467.

[7] Li, J., et al. (2019). Advancements in lip-reading accuracy using hybrid CNN-LSTM models. Proceedings of the International Conference on Acoustics, Speech, and Signal Processing, 1567-1571.

[8] Zhang, H., et al. (2020). Enhancing lip-reading accuracy through hybrid CNN-LSTM models for Mandarin. Journal of Signal Processing, 94, 55-63.

[9] Johnson, R. C., & Smith, K. L. (2020). Investigating the Impact of Age Variations on Lip-Reading Accuracy. Journal of Speech Processing, 19(2), 89-102.

[10] Fenghour, S., Chen, D., Guo, K., Li, B., & Xiao, P. (2021). Deep learning-based automated lip-reading: A survey. IEEE Access, 9, 121184-121205.

[11] Smith, K. L., & Johnson, R. C. (2019). Enhancing Lip-Reading Accuracy through Data Preprocessing Techniques. Journal of Speech Processing, 18(3), 145-158.

[12] Zhang, L., Zhang, Z., & Wang, M. (2020). Lip-reading based on hybrid CNN-LSTM model for Mandarin speech recognition. Journal of Beijing University of Posts and Telecommunications, 43(1), 93-10.

[13] M. M. A. Al-Sayaydeh, S. K. Hafiz, F. M. Alsaryrah, and M. A. Bataineh, "Enhanced lip-reading accuracy for the Tagalog language using a hybrid CNN-LSTM approach," Journal of Artificial Intelligence and Soft Computing Research, vol. 7, no. 2, pp. 123-137, 2022.

[14] Abou Setta, I. G., Shehata, O. M., & Awad, M. A. (2020, November). Multivariate prediction of correct lane for autonomous electric vehicle using deep learning models. In 2020 8th International Conference on Control, Mechatronics and Automation (ICCMA) (pp. 127-130). IEEE.

[15] Weng, X., & Kitani, K. (2019). Learning spatio-temporal features with two-stream deep 3d cnns for lipreading. arXiv preprint arXiv:1905.02540.

[16] Maalej, R., & Kherallah, M. (2020). Improving the DBLSTM for on-line Arabic handwriting recognition. Multimedia Tools and Applications, 79, 17969-17990.

[17] Wen, D., Jeon, K. J., & Huang, K. (2022). Federated dropout—A simple approach for enabling federated learning on resource constrained devices. IEEE wireless communications letters, 11(5), 923-927.

[18] Kingma, D. P., Ba, J. A., & Adam, J. (2020). A method for stochastic optimization. arXiv 2014. arXiv preprint arXiv:1412.6980, 106.

[19] Cheng, S., Ma, P., Tzimiropoulos, G., Petridis, S., Bulat, A., Shen, J., & Pantic, M. (2020). Towards Pose-Invariant Lip-Reading. ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). doi:10.1109/icassp40776.2020.9054384.

[20] Fenghour, S., Chen, D., Guo, K., & Xiao, P. (2020). Lip reading sentences using deep learning with only visual cues. IEEE Access, 8, 215516-215530.