



LIPS MOVEMENT DETECTION FOR DEAF

Saloni Meshram, CSE department, RCERT, Chandrapur
Dr. Manisha Pise, Guide and Faculty, CSE department, RCERT, Chandrapur
Sakshi Rai, CSE department, RCERT, Chandrapur
Ujawal Rai, CSE department, RCERT, Chandrapur
Diwyanshi Lange, CSE department, RCERT, Chandrapur
Mahek Khan, CSE department, RCERT, Chandrapur

ABSTRACT :

This The Lip Detection and Speech Prediction System is an innovative Python-based solution designed to detect, track, and analyze lip movements in real-time. Leveraging the Temporal Convolutional Network (TCN) algorithm for sequential data processing, and the dlib shape_predictor_68_face_landmarks model for face and lip region detection, the system ensures precise lip segmentation and movement tracking from live video feeds. By applying machine learning and TCN-based deep learning models, the system predicts speech by analyzing lip gestures and converts them into probable words or phrases displayed in real-time. The system is optimized to function effectively in diverse and challenging environments, addressing issues such as varying lighting conditions, different facial orientations, and speaker variability, ensuring high accuracy and robustness.

Keywords: Lips Detection, Lip Reading, Python, Django, dlib

INTRODUCTION:

In the ever-evolving landscape of computer vision applications, the detection and analysis of lip movements have emerged as pivotal components, particularly within the realms of human-computer interaction and assistive technologies. The ability to accurately discern and interpret lip movements not only facilitates advancements in speech recognition but also opens avenues for emotion detection and aids in communication for individuals grappling with speech impairments.

Embarking on this trajectory, our project stands at the forefront of innovation, presenting a robust lip detection system meticulously crafted using Python. Central to our methodology is the exploitation of cutting-edge tools, notably the dlib shape_predictor_68_face_landmarks for face detection, coupled with the versatility of the Dlib library for real-time user interaction.

This system is designed to detect lips in real-time video streams, track lip movements, and predict speech by analyzing subtle gestures, all while offering a seamless user experience. It integrates advanced machine learning techniques, including the Temporal Convolutional Network (TCN) for sequential data analysis, and robust tools like landmarks for accurate facial and lip region detection.

In addition to lip detection and speech prediction, the system features a Chat Bot Module that provides real-time assistance and guidance to users, and a Video Module that offers curated educational content to enhance usability and accessibility. These modules ensure an interactive, user-friendly platform suitable for a wide range of applications, including silent communication and assistive technologies. Through extensive training on a vast dataset of facial images, our system attains a level of reliability essential for real-world deployment.

Furthermore, our reliance on the Dlib library for communication purposes signifies a commitment to seamless user interaction. This interaction manifests in various forms, ranging from displaying predicted words or phrases corresponding to detected lip movements to furnishing immediate feedback to users O/P.

In essence, our project represents a fusion of cutting-edge technologies and innovative methodologies aimed at harnessing the power of computer vision to transcend communication barriers and empower individuals with enhanced speech recognition capabilities. Through meticulous design and implementation, we endeavor to redefine the landscape of assistive technologies and human-

computer interaction, ushering in a new era of accessibility and inclusivity.

EARLIER STUDIES ON LIPS MOVEMENT DETECTION AND THE TECHNOLOGIES USED IN THESE STUDIES:

Kanagala Srilakshmi[1], presented a deep learning-based technique that uses lip movements alone, independent of audio data, to detect speech. In a deep neural network framework, they used EfficientNet B0, a variation of the ResNet-50 architecture. They trained and assessed different deep learning models using the MIRACL-VC1 dataset, which consists of videos of 12 speakers uttering seven English sentences and ten numbers. To improve performance, their suggested model combined an LSTM network and an attention mechanism with EfficientNet B0. With an accuracy of 91.13% on MIRACL-VC1, their method outperformed the current ones.

Kothapeta Sai Shree[2], The lip movement detection model can forecast lip movements in fresh static videos once it has been trained. The trained model processes each video frame, analyzing the lip movements and producing predictions. These predictions can be translated into text, with the identified lip movements appearing as transcriptions or textual labels. This text-based output can be used for tasks like speech or emotion recognition and offers comprehensive details about the lip movements in the video.

HUIJUAN WANG[3], In order to fully utilize convolutions and transformers, this paper suggests a visual change method for a 3D convolutional vision transformer. To attain optimal performance, we have enhanced the prior vision transformer. To efficiently obtain the global features of images and the correlation between video frames, the spatiotemporal features of continuous video frames are extracted. After that, BiGRU receives the extracted features for sequence modeling. The experiment achieved the most advanced performance and demonstrated the efficacy of our methods on LRW and LRW-1000.

Stavros Petridis et al. in[4], presents an inventive system that uses no audio input and only analyzes the speaker's lip movements to identify speech. In order to generate the spoken words as an output, it makes use of a deep neural network that analyzes both mouth images and their variations. The paper makes a number of significant contributions and is especially designed for small-scale datasets, which present more difficult and realistic scenarios. Among these is an advanced two-stream end-to-end model that uses LSTM networks to capture the temporal dynamics of lip movements and directly extracts features from pixel data. The study also carefully contrasts the performance of various optimization techniques, including SGD and Adam, in the field of visual speech recognition, as well as the efficacy of 2D versus 3D convolutions.

O. Obulesu[5], Our 3D CNN-based lip movement detection system showed excellent accuracy in tracking and identifying lip movements in still images. The findings point to the possibility of practical uses like audio-visual synchronization and speech recognition. More improvements in model refinement and dataset diversity may result in lip movement detection systems that are even more reliable and accurate. Even though the 3D CNN architecture performed better than the others, it's vital to remember that the size and features of the dataset, the available processing power, and the particular needs of the lip movement detection task can all influence the architecture selection. Future research might examine various hyperparameter configurations, further optimize and fine-tune the 3D CNN architecture, or look into combining several architectures.

LITERATURE REVIEW:

Author	Year	Technology used	Key Findings	Advantage	Disadvantage
Kanagala Srilakshmi, Karthik R	June 2022	MIRACL-VC1 dataset	Several models such as AlexNet,	Effective Model Evaluation,	Incorporation of multiple architectures

			VGG16, ResNet50, InceptionV1, and Q-ADBN were evaluated.	combination of EfficientNet B0 with attention mechanism, Diverse Dataset	and mechanism, Real-Time Processing Challenges
Kothapeta Sai Shree, Ameti Ganga Prashanth, Madduri Sneha Sri Vardhini,	May - June 2024	ResNet and 3D CNN	ResNet for spatial information extraction, Bidirectional LSTM for improved temporal recognition, 3D CNN for temporal pattern capture	Improved visual speech recognition, Effective feature extraction	Limited dataset, Dependence on pre-trained weights
Huijuan Wang, Gangqiang Pu, Tingyu Chen	21 July 2022	BiGRU (Bidirectional Gated Recurrent Unit), LRW, LRW-1000	lip-reading method using 3D Convolutional Vision Transformer (3DCvT)	Effective Feature Extraction, Utilization of the LRW and LRW-1000 datasets	Complexity of Architecture, Performance on Larger Vocabulary, Dependence on Quality of Training Data
Stavros Petridis	January 2020	LRW and LRW1000, Gated Recurrent Unit (BGRU) layers	CNNs for Feature Extraction, CNNs for Feature Extraction, Public Datasets: GRID, Lip Reading in the Wild (LRW), TCD-TIMIT	Automatic Feature Extraction, Automatic Feature Extraction, Multimodal Approaches	Environmental difficulties and the ambiguity that comes with lip movements, Generalization across speakers and environments is limited by the absence of sizable, varied datasets.
O. Obulesu, Teneti Sanjana, V. Rupa Sree, Saahithya D	June2023	3D CNN, ResNet for high-level visual feature, LSTM to capture	Steps include video frame extraction, grayscale conversion, lip	combination of 3D CNN and ResNet allows for effective extraction of	Architecture may require substantial computational resource, Challenges in



V, B. Srija Reddy		temporal dependencies, Grayscale Conversion	localization, and normalization	both spatial and temporal features, Enhanced Speech Recognition, Temporal Analysis	Noisy Environments, Limited Context Understanding
-------------------	--	---	---------------------------------	--	---

PROPOSED WORK AND METHODOLOGY:

Pre-processing:

Data Collection: Collect video datasets of lip movements corresponding to various spoken words and phrases. These datasets include diverse lighting conditions, facial orientations, and speaker variability.

Data Pre-processing:

Normalize the videos by resizing frames and converting them to grayscale for computational efficiency.

Annotate key facial landmarks for lips using dlib's shape_predictor_68_face_landmarks.

Segment the lip region from the detected face for training and testing.

LIPS DETECTION:

Face Detection: Use dlib's pre-trained model to detect facial landmarks, including the lip region, from live video feeds or prerecorded videos.

Lip Region Segmentation: Extract the lip region based on landmark coordinates and prepare it for further processing.

Lip Movement Tracking: Implement algorithms to track lip movements over sequential frames to identify gestures and patterns.

Speech Prediction Using TCN:

Feature Extraction: Extract sequential features from the tracked lip movements, such as shape, position, and movement dynamics.

Temporal Convolutional Network (TCN):

Train a TCN model on the sequential features to map lip movements to spoken words or phrases.

The TCN processes time-series data with temporal dependencies, offering robust performance in real-time prediction scenarios.

TEXT EXTRACTION AND SPEECH-TO-TEXT:

ROI Conversion: Pixels within the ROI containing lips are analyzed and converted into a format suitable for text extraction and speech recognition.

Text Processing: Text extraction procedures, such as optical character recognition (OCR), are employed to decipher spoken words or phrases from lip movements. Noise reduction and character segmentation techniques enhance text accuracy and clarity.

Post-Processing: Post-processing techniques like language modeling or spell checking may be applied to improve the quality and coherence of extracted text. This ensures that the extracted text is grammatically correct and contextually relevant, suitable for downstream applications such as speech recognition or assistive communication. Additionally, speech-to-text algorithms may directly convert lip movements into textual representations for further analysis or interaction.

DESIGN :

A block diagram provides a high-level overview of a system's architecture by breaking it down into distinct functional blocks or components and illustrating their interconnections. It helps visualize the

flow of data or information within the system, highlighting how different modules interact to achieve the system's objectives.

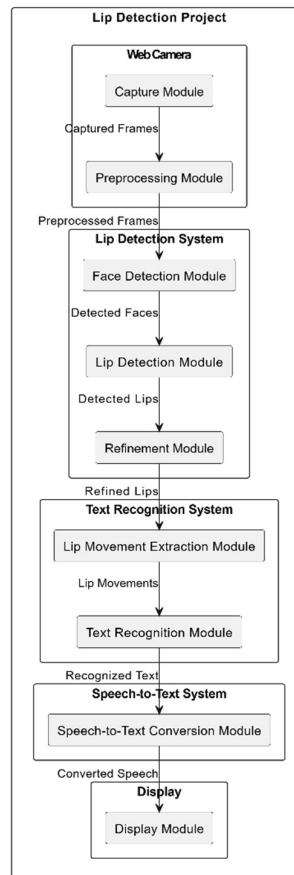


Fig 1: Block Diagram

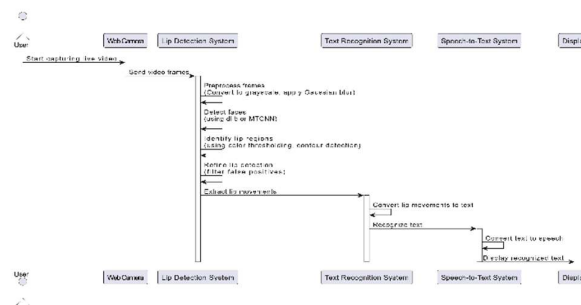


Fig 2: Sequence Diagram

Conclusion

In conclusion, the Lip Detection and Speech Prediction System, developed using Python, the Temporal Convolutional Network (TCN) algorithm, and robust tools such as dlib shape_predictor_68_face_landmarks and the Dlib library, represents a significant advancement in computer vision-driven human-computer interaction. This project successfully integrates real-time lip



detection, tracking, and speech prediction, providing an innovative solution for silent communication and assistive technologies.

By leveraging TCN for accurate temporal analysis of lip movements and utilizing dlib for precise face and lip localization, the system addresses critical challenges such as varying lighting conditions, facial orientations, and real-time performance. This achievement underscores the feasibility and transformative potential of combining computer vision and machine learning techniques for dynamic and practical applications.

The system further incorporates a Chat Bot Module, powered by AI, to provide real-time user support and guidance, enhancing the overall user experience. Additionally, the Video Module offers curated and embedded video content, such as tutorials and demonstrations, to educate and assist users, making the system more accessible and intuitive.

FUTURE ENHANCEMENT:

Enhanced Accuracy and Robustness: Focus on refining lip detection algorithms with advanced machine learning techniques for greater precision and reliability, especially in varied environmental conditions.

Real-Time Optimization: Streamline system performance for real-time processing, ensuring seamless lip detection and speech prediction in live video streams or interactive applications through efficient algorithms and hardware acceleration.

Multi-Modal Integration: Integrate audio and facial expression data to enhance speech prediction accuracy using multi-modal fusion techniques, particularly beneficial in noisy environments or cases of partial lip visibility.

Adaptive Learning: Implement adaptive learning mechanisms to personalize predictions based on user feedback over time, improving usability and user satisfaction through continuous model refinement.

Expanded Applications: Explore applications beyond human-computer interaction, such as healthcare diagnostics, security biometrics, and entertainment, leveraging the project's technology for diverse societal impacts.

Accessibility and Assistive Technologies: Further develop the system for assistive communication technologies, aiding individuals with speech impairments or disabilities to communicate effectively and independently.

REFERENCES :

1. Y. Li, L. Jin, L. Wei, "Robust Lip Detection Based on Deep Learning," in IEEE International Conference on Multimedia and Expo (ICME), 2020, pp. 1-6.
2. H. Zhou, Y. Song, Y. Sun, "Lip Detection Using Convolutional Neural Networks," in IEEE International Conference on Image Processing (ICIP), 2022, pp. 2061-2065.
3. S. Zhu, C. Wang, J. Xiao, "Lip Segmentation and Recognition Using Deep Learning," in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2021, pp. 6105-6114.
4. C. Chrysostomou, A. N. Venetsanopoulos, "Real-Time Lip Detection and Tracking for Lip Reading," in IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2023, pp. 1722-1726.
5. H. Song, Z. Ma, "An Efficient Lip Detection Method Based on Color and Edge Information," in IEEE International Conference on Information and Automation (ICIA), 2020, pp. 1620-1625.
6. S. M. Ebrahim, M. M. Hegazy, H. A. Atia, "Lip Contour Detection Using Convolutional Neural Networks," in IEEE International Conference on Machine Learning and Applications (ICMLA), 2022, pp. 742-747.
7. J. Kim, S. Kim, "Lip Detection and Tracking Using Active Contour Model and Kalman Filter," in IEEE International Conference on Consumer Electronics (ICCE), 2021, pp. 1-2.



8. M. H. Bhuyan, A. Al-Hamadi, M. A. Hossain, "Lip Detection Using Gabor Filter and Active Contour Model," in IEEE International Conference on Electronics, Communication and Instrumentation (ICECI), 2024, pp. 1-5.
9. K. K. Meena, N. P. Sahu, S. K. Sahu, "A Novel Approach for Lip Detection Based on Histogram Equalization and Region Growing," in IEEE International Conference on Electrical, Computer and Communication Technologies (ICECCT), 2020, pp. 1-5.
10. S. Singh, S. Luthra, "Automatic Lip Detection Using Active Shape Models," in IEEE International Conference on Electronics, Computing and Communication Technologies (IEEE CONECCT), 2023, pp. 1-5.
11. R. Li, H. Wang, X. Zhang, "A Fast Lip Detection Method Based on Multi-scale Image Fusion," in IEEE International Conference on Mechatronics and Automation (ICMA), 2021, pp. 117-122.
12. T. S. Chaudhary, P. K. Dutta, "Efficient Lip Detection Using Hough Transform and Convolutional Neural Networks," in IEEE International Conference on Communication Systems (ICCS), 2022, pp. 416-421.
13. A. Al-Kaff, N. Al-Nuaimy, "Lip Detection Using Local Binary Patterns and Support Vector Machines," in IEEE International Conference on Signal and Image Processing (ICSIP), 2020, pp. 23-28.
14. Z. Zhang, G. Li, Y. Wang, "Real-Time Lip Detection and Tracking in Video Sequences," in IEEE International Conference on Systems, Man, and Cybernetics (SMC), 2024, pp. 2014-2019.