

International Journal of Research Publication and Reviews

Journal homepage: www.ijrpr.com ISSN 2582-7421

LIP READING AI: HARNESSING AI FOR SILENT COMMUNICATION

R Sameera¹, R Sanjeev Dasu², M Chandu³, P Manohar Naidu⁴

¹RSameera, Information Technology, GMRIT, Rajam, India
²RSanjeevDasu, Information Technology, GMRIT, Rajam, India
³MChandu,Information Technology, GMRIT, Rajam, India
⁴PManohar Naidu,Information Technology, GMRIT, Rajam, India

ABSTRACT :

Lip reading AI, or automated visual speech recognition (AVSR), uses advanced algorithms to interpret spoken language by analysing lip movements. This technology integrates computer vision, deep learning, and natural language processing to enable applications in accessibility, speech enhancement, and silent communication systems. Leveraging convolutional neural networks (CNNs), recurrent neural networks (RNNs), and transformer architectures, recent advancements have significantly improved accuracy. Additionally, techniques like generative adversarial networks (GANs) for occlusion handling and attention mechanisms for temporal modelling enhance the robustness of lip-reading systems. This paper presents a project focused on developing a Lip-Reading AI system, specifically designed to assist individuals with speech and hearing disabilities, such as those who are mute or deaf. By incorporating state-of the-art methodologies, including hybrid and multi-modal approaches, the project aims to enhance communication accessibility and independence for differently-abled individuals. It explores innovations in real-time applications, multilingual support, and practical implementations to improve accuracy, robustness, and adaptability across different environments and user groups.

Key Words: Automated visual speech recognition (AVSR), convolutional neural networks (CNNs), recurrent neural networks (RNNs).

1. INTRODUCTION

Lip reading, also known as visual speech recognition (VSR), is an advanced artificial intelligence (AI) technique that enables computers to interpret speech by analysing++92 lip movements. Traditional speech recognition systems rely primarily on audio signals, but in noisy environments or for individuals with hearing impairments, lip reading technology offers a crucial alternative for accurate speech interpretation. With the advancements in deep learning and computer vision, modern AI-powered lip-reading models leverage convolutional neural networks (CNNs), recurrent neural networks (RNNs), transformers, and encoder-decoder architectures to improve speech prediction accuracy. These models extract spatial and temporal features from video sequences, allowing them to recognize words and sentences effectively. The objective of this project is to develop an AI-based lip-reading system that can accurately transcribe speech from silent video footage. This will involve training deep learning models on large scale lip-reading datasets, integrating feature extraction techniques, and optimizing performance for real-time applications.

This AI-driven approach has wide-ranging applications, including: Assistive technology for the hearing impaired, Enhanced communication in noisy environments, Security and surveillance for speech-based authentication, Human-computer interaction and virtual assistants By leveraging state-of-theart AI architectures, this project aims to push the boundaries of silent speech recognition, improving both accuracy and usability in real-world scenarios

2. LITERATURE SURVEY

This survey employs a systematic review approach to analyze deep learning advancements in automated lip reading. The methodology includes a structured examination of feature extraction techniques, model architectures, and dataset utilization. The survey evaluates recent contributions, emphasizing the role of transformers, self-supervised learning, and multi-modal fusion in enhancing accuracy. To address challenges such as speaker variations and occlusions, the paper explores solutions like domain adaptation and data augmentation, suggesting a comparative analysis of their effectiveness. The study further reviews benchmark datasets, including GRID, LRW, and LRS2, highlighting their characteristics and relevance to lip-reading tasks.

Ai and Fang propose a novel methodology for lip-reading that integrates multi-motion-informed contexts into a cross modal language model. Their approach extracts diverse lip-motion features from multiple subspaces using stacked convolutional layers, capturing different aspects of lip dynamics. These extracted features are then fed into a cross-modal language model equipped with a source-target attention mechanism, allowing effective alignment between visual inputs (lip movements) and textual representations. This integration ensures that the model leverages both spatial and temporal variations in lip motion to improve recognition accuracy. To enhance learning efficiency, the authors employ a piece-wise pre-training

strategy inspired by multi-task learning. This involves separately training the visual feature extraction module and the decoder responsible for text generation, ensuring that both components are optimized before full model integration. The proposed methodology is evaluated on benchmark datasets such as LRS2, LRS3, LRW, and GRID, demonstrating superior performance compared to existing models. By combining multi-motion-informed representations with cross-modal learning, their approach significantly improves automated lip-reading, making it more robust to variations in speaker articulation and environmental conditions.

3. PROPOSED SYSTEM

Collect high-quality datasets of videos with labelled transcriptions, such as LRW (Lip Reading in the Wild), GRID, LRS2, and LRS3. • Ensure the dataset covers various speakers, accents, and lighting conditions for generalization.



Fig 1: Flowchart of LipReading

1. Data Collection:

The process begins with Collected and analysed key lip-reading datasets like GRID, LRW, LRS2, and LRS3. Highlighted differences in vocabulary size, resolution, speaker diversity, and sentence structure.

2. Preprocessing:

Applied face detection and alignment to extract and standardize mouth regions. Used 3D-CNNs and optical flow for spatiotemporal feature extraction.

3. Model Architectures:

Reviewed deep learning models like CNNs, RNNs, and Transformers for lip-reading. Highlighted hybrid approaches combining spatial and temporal modelling.

4. Speech Recognition Approaches:

Compared CTC-based and sequence-to-sequence models with attention mechanisms. Explored language model integration and end-to-end vs modular systems.

5. Multimodal Fusion & Enhancement:

Discussed combining audio-visual inputs using early/late fusion and attention. Showed that visual input boosts performance in noisy audio conditions.

6. Challenges & Performance Evaluation:

Outlined key challenges like speaker variation, lighting, and dataset limitations. Used WER and CER to evaluate model accuracy and generalizability.

4. RESULT



Fig. 2: Expected Result

5. CONCLUSION

Lip Reading AI has evolved significantly in recent years, transitioning from traditional handcrafted feature extraction methods to powerful deep learning-based models. Early approaches relied on techniques like DCT and optical flow, but these were limited by their sensitivity to noise and variation. Modern systems now leverage Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM) units, and more recently, Transformer-based architectures to automatically learn spatiotemporal features from video data. These advancements, combined with large-scale datasets such as LRW and LRS3, have led to remarkable improvements in word and sentence-level visual speech recognition. Despite these advancements, Lip Reading AI continues to face notable challenges. A core issue lies in the visual ambiguity of phonemes, as multiple speech sounds can appear similar on the lips (viseme confusion). Additionally, variability in speaker a ppearance, lighting, head pose, and occlusion can reduce model robustness. There is also a need for real-time performance, especially in applications involving low-resource or embedded devices. The limited availability of annotated data for training deep models, especially across different languages and dialects, poses further obstacles to generalization and deployment in the wild. Looking ahead, the future of Lip Reading AI is promising. Continued research into multimodal fusion (combining audio and visual data), self-supervised learning, domain adaptation, and model compression is helping to overcome current barriers. Applications are expanding across accessibility tools for the hearing impaired, silent speech interfaces, in-car systems, and even surveillance. As models become more efficient and accurate, Lip Reading AI is set to become a key component in robust, context-aware, and human-centric artificial intelligence systems.

6. LIMITATIONS

Despite significant advancements, the Lip Reading AI system developed in this project faces several limitations. One major challenge is the visual ambiguity of phonemes, where different speech sounds appear similar on the lips, leading to recognition errors. Variability in speaker appearance, such as differences in lip shape, skin tone, and speaking style, further reduces the model's generalization capability. Environmental factors like inconsistent lighting, background clutter, and head movements introduce additional complexities, often degrading prediction accuracy. Although we in this project

developed and deployed a high-performance Lip Reading AI system, it has some limitations. Visual confusability is a wrong visual impression due to the ambiguity of visual aspects in phoneme (some speech sounds appear similar on the lips, thus causing many recognition failures.) Additionally, the speaker can show large variations in appearance including lip shape, skin tone, and speaking manner, which makes the model hardly generalize. Environmental factors such as non-uniform lighting variations, background clutters, and head movement add more challenges to this task and usually lead to a decrease in the prediction performance. A third two reducing the time requirements.

REFERENCES :

- 1. S. Fenghour, D. Chen, K. Guo, B. Li and P. Xiao, "Deep Learning-Based Automated Lip-Reading: A Survey,"2024 in IEEE Access
- X. Ai and B. Fang, "Cross-Modal Language Modeling in Multi-Motion-Informed Context for Lip Reading," in IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol. 31, pp. 2220-2232, 2023, doi: 10.1109/TASLP.2023.3282109.
- P. Ma, Y. Wang, S. Petridis, J. Shen and M. Pantic, "Training Strategies for Improved Lip-Reading," ICASSP 2022 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Singapore, Singapore, 2022, pp. 8472-8476, doi: 10.1109/ICASSP43922.2022.9746706.
- P. Ma, B. Martinez, S. Petridis and M. Pantic, "Towards Practical Lipreading with Distilled and Efficient Models," ICASSP 2021 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Toronto, ON, Canada, 2021, pp. 7608-7612, doi: 10.1109/ICASSP39728.2021.9415063.
- M. Hao, M. Mamut, N. Yadikar, A. Aysa and K. Ubul, "A Survey of Research on Lipreading Technology," in IEEE Access, vol. 8, pp. 204518-204544, 2020, doi: 10.1109/ACCESS.2020.3036865
- 6. Li, Y., Hashim, A. S., Lin, Y., Nohuddin, P. N., Venkatachalam, K., & Ahmadian, A. (2024). AI-based visual speech recognition towards realistic avatars and lip-reading applications in the metaverse. Applied Soft Computing, 164, 111906.
- Kulkarni, A. H., &Kirange, D. (2019, July). Artificial intelligence: A survey on lip reading techniques. In 2019 10th International Conference on Computing, Communication and Networking Technologies (ICCCNT) (pp. 1-5). IEEE..
- Geetha, C., & Reddy, S. S. V. P. (2024, June). AI Lip Reader Detecting Speech Visual Data with Deep Learning. In 2023 4th International Conference on Intelligent Technologies (CONIT) (pp. 1-6). IEEE.
- Juyal, A., Joshi, R. C., Jain, V., & Chaturvedy, S. (2023, November). Analysis of Lip-Reading Using Deep Learning Techniques: A Review. In 2023 International Conference on the Confluence of Advancements in Robotics, Vision and Interdisciplinary Technology Management (IC-RVITM) (pp. 1-6). IEEE.