

# Visual Speech Recognition for Seamless Communication with Hearing Impaired Persons

Md. Laraib Ahmad , Scholar ID: 2322207

## Supervisor

Dr. Debbrota Paul Chowdhury



Department of Computer Science and Engineering  
National Institute of Technology Silchar

May 14, 2025

- ① Introduction
- ② Research Advancement
- ③ Motivation
- ④ Literature Review
- ⑤ Research Gap
- ⑥ Problem Statement
- ⑦ Objective
- ⑧ Available Datasets
- ⑨ Dataset Preparation for Visual Speech Recognition System
- ⑩ Generalised Block Diagram
- ⑪ Methodology
- ⑫ conclusion
- ⑬ Future Work
- ⑭ References
- ⑮ Publication

## What is Lip Reading:



fig no - 1: <https://livingwithhearingloss.com/2016/04/19/lipreading-in-paradise/>

# Why Lip Reading is important?

Communication aid for the deaf and hard-of-hearing

- Primary Communication Tool
- Complement to hearing aids and cochlear implants



<https://www.inc.com/john-boitnott/this-entrepreneur-is-solving-one-of-the-biggest-problems-all-deaf-people-face.html>

<https://www.connecthear.org/post/all-about-cochlear-implants>

cont...

## Laryngeal cancer

- laryngeal cancer can result in loss of speech or changes in the voice, especially if the cancer or its treatment affects the vocal cords or requires the removal of parts of the larynx. In such case speaker can only move their lips to communicate.



<https://utswmed.org/medblog/cold-flu-allergy-hurt-your-voice/>

Time Period	Advancement
<b>Early 20th century</b>	The concept of lip reading as a skill for the deaf and hard-of-hearing began to be formally studied and taught. Schools for the deaf often included lip reading in their curricula.
<b>1970s</b>	Early research in automated lip reading began. Initial efforts focused on understanding the visual aspects of speech and how they could be captured and analyzed by computers.
<b>1980s and 1990s</b>	Advances in computer vision and pattern recognition led to more sophisticated experiments with automated lip reading systems. Algorithms to recognize visual speech elements were developed.
<b>2000s</b>	Deep learning and AI started to significantly improve automated lip reading accuracy. Researchers used statistical models like Hidden Markov Models (HMMs) to analyze visual speech data.
<b>2010s</b>	Deep learning techniques, especially CNNs and RNNs, began to be applied to lip reading, leading to substantial improvements. Large datasets and improved computational power also contributed to advancements.
<b>2016</b>	Google DeepMind's LipNet model demonstrated high accuracy in lip reading by leveraging deep learning techniques.
<b>2020s</b>	Ongoing research continues to refine and improve lip reading technologies, which are now applied in fields such as assistive technology for the deaf and hard-of-hearing, security, and human-computer interaction.

- Accessibility for Hearing-Impaired Individuals
- Speech Enhancement in Noisy Environments
- Multimodal Systems
- Advances in AI and Machine Learning

Author, year	Methodology	Dataset Used	Findings
Freitas <i>et al.</i> , 2016 [1]	<p>LipNet model: spatiotemporal convolutional neural networks (STCNNs) to extract spatial and temporal features, followed by Bidirectional Gated Recurrent Units (Bi-GRUs) to capture temporal dependencies, and employs Connectionist Temporal Classification (CTC) loss for end-to-end training without the need for pre-segmented data</p>	GRID corpus	<ul style="list-style-type: none"><li>Achieving a 95.2% sentence-level accuracy on the GRID corpus.</li><li>The study highlights the effectiveness of combining spatiotemporal convolutions, recurrent networks, and Connectionist Temporal Classification (CTC) for sentence-level prediction, marking a major improvement in automated lipreading</li></ul>

cont...

Author, year	Methodology	Dataset Used	Findings
Zimmermann <i>et al.</i> , 2020[2]	<p>The paper employs two methods: Temporal Conditional GANs (TC-GANs) to generate lip movement videos for unseen utterances and a viseme-concatenation approach to synthesize videos by mapping phonemes to visemes to enable zero-shot learning in visual speech recognition.</p>	OuluVS2 dataset	<p>Using GANs for zero-shot learning significantly improves visual speech recognition accuracy for unseen utterances, effectively addresses the cold-start problem, and generalizes to new languages, with GANs outperforming the viseme-concatenation approach.</p>

cont...

Author, year	Methodology	Dataset Used	Findings
XIAO <i>et al.</i> , 2020[3]	<p>The methodology involves preprocessing video frames to extract lip regions, using a spatial-temporal CNN to generate features, applying a transformer-based model to classify visemes, and converting visemes to words through perplexity analysis for sentence prediction.</p>	BBC LRS2 dataset.	<p>The paper finds that the proposed viseme-based lip reading system significantly improves word accuracy with a 15% reduction in Word Error Rate (WER), achieves a Viseme Error Rate (VER) of 4.6%, and demonstrates robustness to varying lighting conditions, though further optimization is needed in converting visemes to words</p>

cont...

Author, year	Methodology	Dataset Used	Findings
Xie et al., 2024[4]	The methodology involves multi-scale lip motion video extraction, dynamic augmentation, and an end-to-end VSR system with multi-system fusion using diverse encoders for optimal visual speech recognition performance.	The paper uses the CN-CVS dataset for training, along with the development sets of CNVSR-Single/Multi datasets from the Chinese Continuous Visual Speech Recognition Challenge (CNVSR) 2023.	The paper finds that the proposed multi-system VSR approach with E-Branchformer encoder and ROVER fusion achieves leading performance with 34.76% CER in the Single-Speaker Task and 41.06% CER in the Multi-Speaker Task, securing first place in all three CNVSR 2023 tracks.

cont...

Author, year	Methodology	Dataset Used	Findings
Pingchuan et al., 2022[5]	<p>The methodology involves enhancing VSR performance through prediction-based auxiliary tasks, hyperparameter optimization, data augmentation (like time-masking), and pre-training/fine-tuning across multiple languages.</p>	<p>The paper uses the LRS2, LRS3, CMLR (Mandarin), and CMU-MOSEAS (Spanish) datasets for training and evaluation, with a focus on publicly available datasets for achieving state-of-the-art VSR performance across multiple languages. Additionally, the LRW and AVSpeech datasets are used in some experiments for further improvements.</p>	<p>The paper finds that careful model design, including prediction-based auxiliary tasks, data augmentation, and hyperparameter optimization, can significantly improve visual speech recognition performance, even surpassing models trained on much larger datasets.</p>

cont...

Author, year	Methodology	Dataset Used	Findings
GUO <i>et al.</i> , 2020[6]	<p>The methodology involves using a viseme-to-word conversion system with perplexity analysis, where visual speech input is processed through word lookup, chunkification, and iterative beam search to identify the most likely word sequences based on a pre-trained language model.</p>	<p>The paper uses two datasets for experimentation:</p> <ul style="list-style-type: none"><li>• OuluVS Dataset: This consists of short phrases like "hello," "excuse me," "I am sorry," etc.</li><li>• BBC LRS2 Dataset: This contains longer and more varied sentences from BBC videos, making it more challenging due to a wide range of speakers and vocabulary</li></ul>	<p>The findings show that the model effectively predicts short phrases with 100% accuracy and performs reasonably well on longer sentences using perplexity analysis, though it struggles with increased errors when word boundaries are unknown.</p>

cont...

Author, year	Methodology	Dataset Used	Findings
Ro et al., 2024[7]	<p>The proposed methodology adapts a pre-trained lip reading model to individual speakers by incorporating speaker-specific visual prompts and low-rank adaptation in the visual encoder, along with input prompt tuning in the language decoder to capture both visual and linguistic speaker-specific traits.</p>	<p>VoxLRS-SA , which is specifically designed for speaker-adaptive lip reading in real-world, sentence-level scenarios.</p>	<p>The study finds that combining vision- and language-level adaptations using lightweight methods like LoRA and prompt tuning significantly improves sentence-level lip reading accuracy for unseen speakers, even with limited speaker-specific data.</p>

cont...

Author, year	Methodology	Dataset Used	Findings
Pantic <i>et al.</i> , 2022[8]	<p>The methodology combines advanced temporal models (like DC-TCN), effective data augmentations (Time Masking, mixup), and training strategies (self-distillation, word boundary indicators) to systematically enhance lip-reading performance.</p>	LRW dataset	<p>The paper finds that combining advanced temporal models (DC-TCN), Time Masking augmentation, and training strategies like self-distillation and word boundary indicators achieves state-of-the-art performance in lip-reading, particularly improving recognition of difficult words.</p>

cont...

Author, year	Methodology	Dataset Used	Findings
Jakob Uszkoreti et al., 2017[9]	The paper introduces the Transformer, a sequence transduction model based entirely on self-attention mechanisms without recurrence or convolution using an encoder-decoder architecture to efficiently capture long-range dependencies in data.	WMT 2014 English-German, WMT 2014 English-French	The paper demonstrates that the Transformer model, relying solely on self-attention, achieves state-of-the-art results on translation tasks with significantly faster training and better parallelization compared to RNN- or CNN-based models.

- Accuracy is not very high for word prediction.
- No dataset available for different accents.
- Viseme-based Challenges
- Cross-lingual Transfer Learning

Despite significant progress, visual speech recognition faces challenges such as variability in lip movements, diverse speaking styles, and need for large and labeled datasets for training. Ongoing research aims to address these challenges and further refine technology making it more accurate, reliable and widely applicable.

- To create a dataset of different accents.
- To develop an algorithm that can recognize words correctly independent of accents.

Dataset	Type	Description
<b>GRID Corpus</b>	Sentence-level	Contains 34 speakers uttering structured sentences with fixed vocabulary (1000 unique sentences).
<b>LRS2 (Lip Reading Sentences 2)</b>	Sentence-level	Contains over 224 hours of data from BBC programs with spoken sentences for audio-visual speech recognition.
<b>LRS3 (Lip Reading Sentences 3)</b>	Sentence-level	Larger version of LRS2 with over 475 hours of videos for lip reading in challenging conditions.
<b>TCD-TIMIT</b>	Continuous Speech	Phonetically balanced dataset with 59 speakers reading 98 sentences, suitable for continuous speech recognition.
<b>AVLetters</b>	Alphabet-level	Dataset with speakers uttering letters A-Z multiple times for isolated letter recognition.
<b>LRW (Lip Reading in the Wild)</b>	Word-level	Contains over 500 different words spoken by various speakers extracted from TV broadcasts.
<b>OuluVS2</b>	Phrase-level	Contains 53 speakers saying 10 phrases, repeated 6 times per phrase, for small-scale phrase recognition.

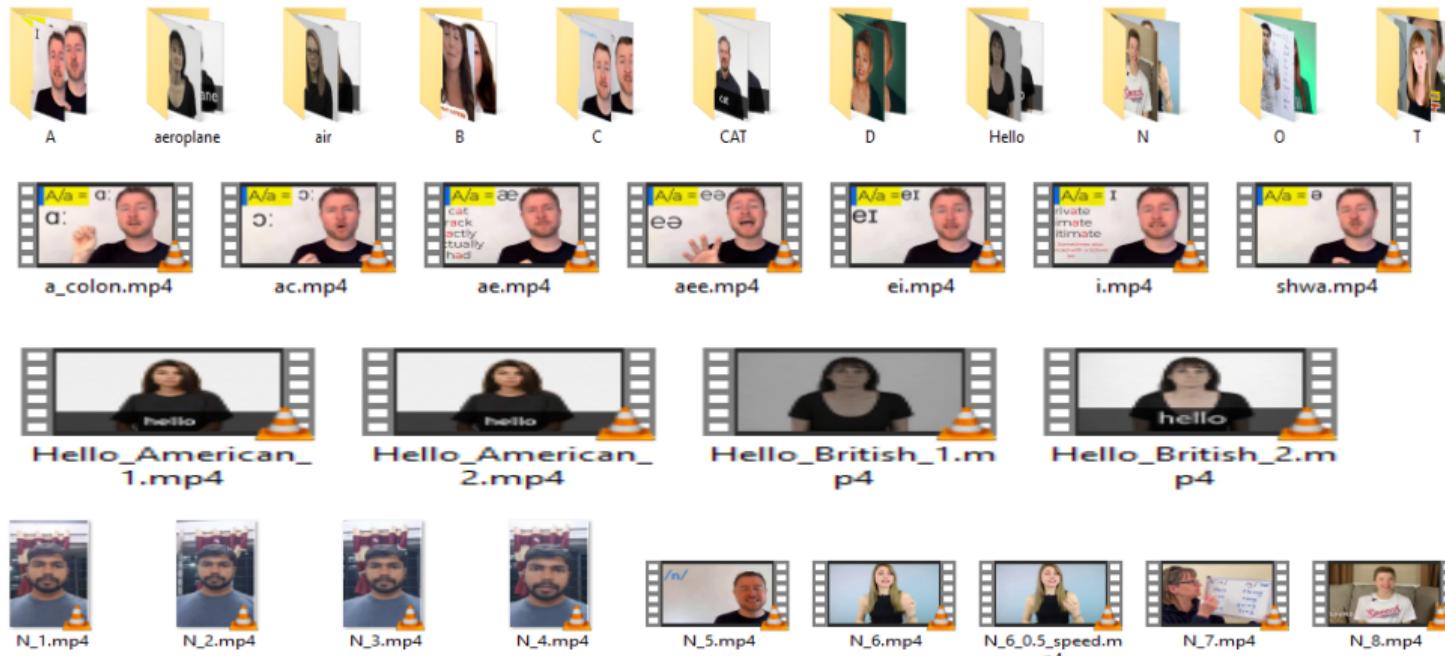
Table 1: Available Datasets for Visual Speech Recognition by Lip Reading

# Dataset Names and Images for Visual Speech Recognition

Dataset	Image
GRID Corpus	
LRS2 (Lip Reading Sentences 2)	
LRS3 (Lip Reading Sentences 3)	
TCD-TIMIT	

Dataset	Image
AVLetters	 
LRW (Lip Reading in the Wild)	  
OuluVS2	 

# Dataset Preparation for Visual Speech Recognition System



# Generalized block diagram of visual speech recognition system

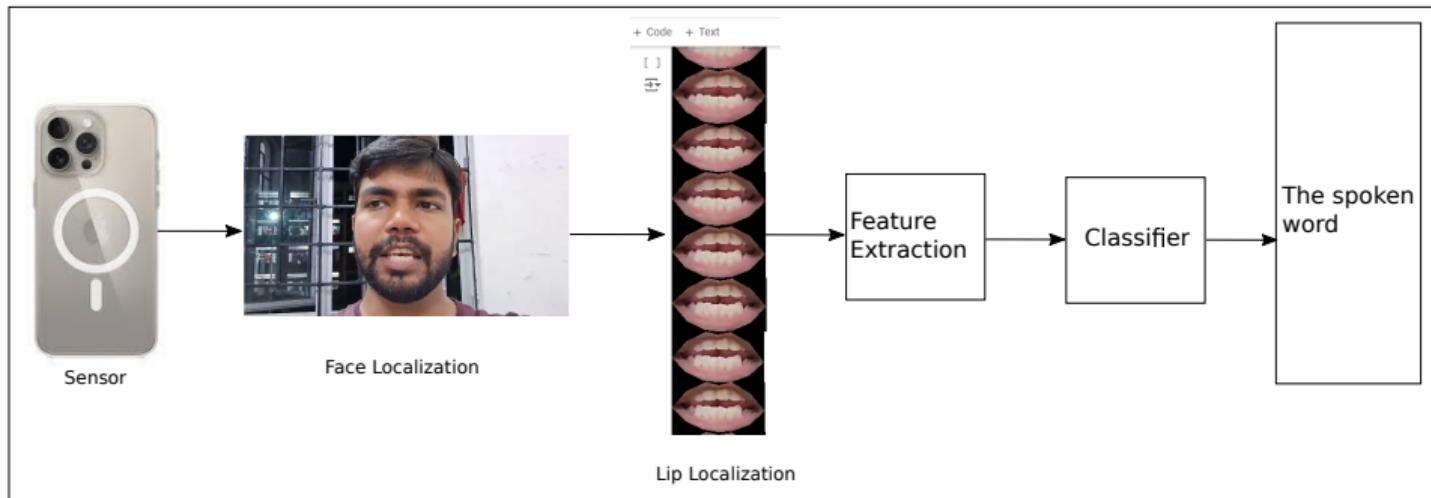


Figure 1: Generalized block diagram of visual speech recognition system

Lip Points with Numbers



Figure 2: lip feature extraction

cont...

## Deep Learning & Geometrical-based Approach

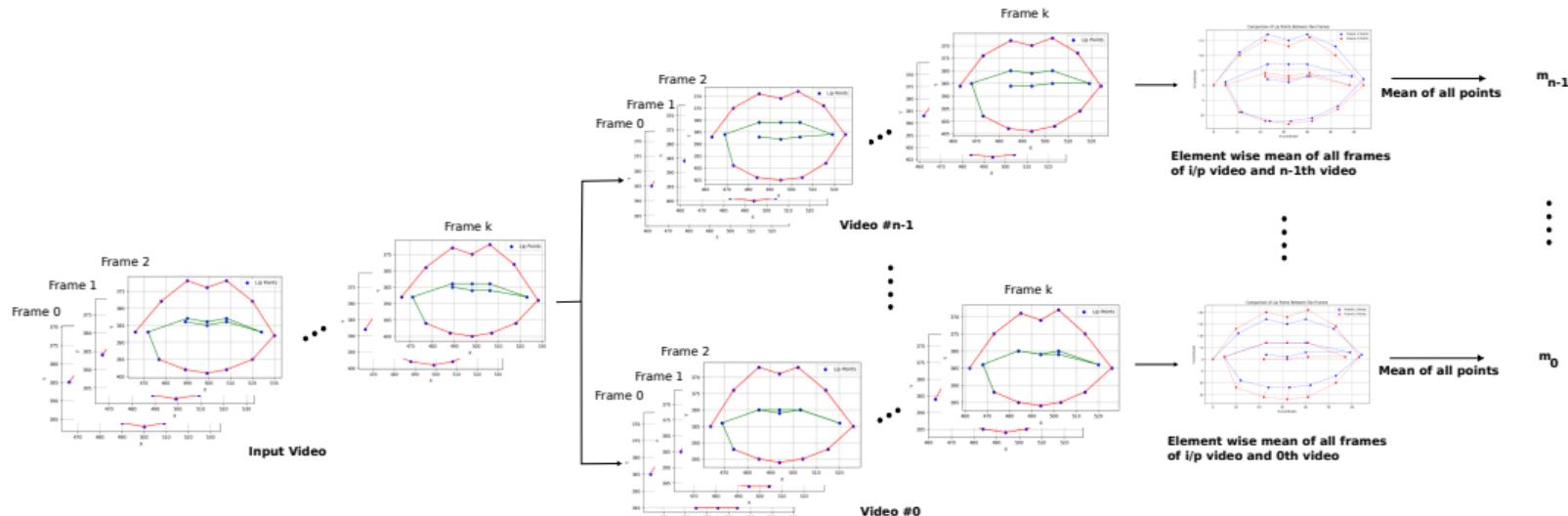


Figure 3: Calculating Euclidean Distance for the entire Dataset

$$\text{Predicted Word} = \min\{m_0, m_1, m_3, \dots, m_{n-1}\}$$

cont...

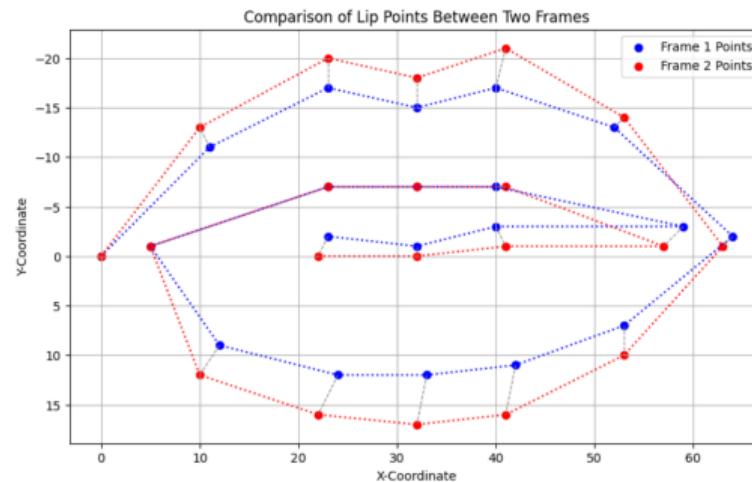


Figure 4: lip comparision

# Detection of whether speaker is speaking or silent

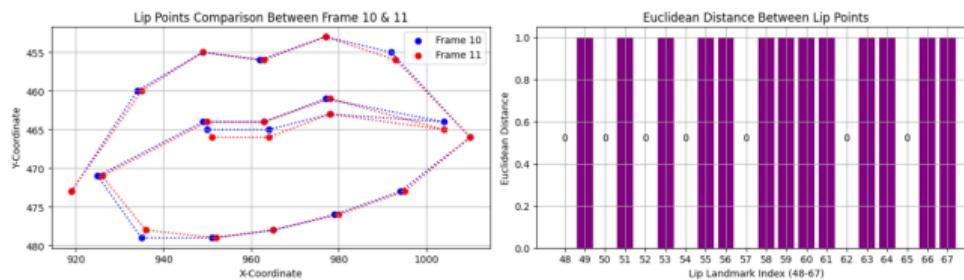


Figure 5: frame 10 and 11 lip comparision

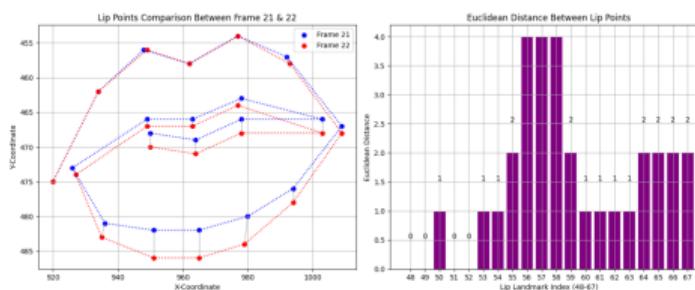


Figure 6: Beginning of Speaking Frame frame 21 and 22

cont...

Current Frame	Next Frame	Average Euclidean Distance
1	2	0.41
2	3	0.42
3	4	0.40
4	5	0.44
5	6	0.53
6	7	0.55
7	8	0.67
8	9	0.30
9	10	0.42
10	11	0.71
11	12	0.70
12	13	0.65
13	14	0.39
14	15	0.87
15	16	0.57
16	17	0.71
17	18	1.06
18	19	0.79
19	20	1.38
20	21	0.54
21	22	1.63
22	23	2.49
23	24	3.11
24	25	2.15
25	26	1.40
26	27	1.62
27	28	1.93
28	29	2.30
29	30	2.11
30	31	1.17
31	32	4.54
32	33	2.40
33	34	3.88
34	35	2.79
35	36	2.12
36	37	1.60
37	38	4.02
38	39	0.85
39	40	4.80
40	41	2.04
41	42	0.44
42	43	2.21
43	44	2.41
44	45	0.74
45	46	1.69
46	47	0.57
47	48	0.88
48	49	1.21
49	50	1.29
50	51	0.49
51	52	0.87
52	53	1.42
53	54	1.02
54	55	0.32
55	56	1.10
56	57	0.15
57	58	0.48
58	59	0.35
59	60	0.20

Figure 7: Detection of Speech Activity

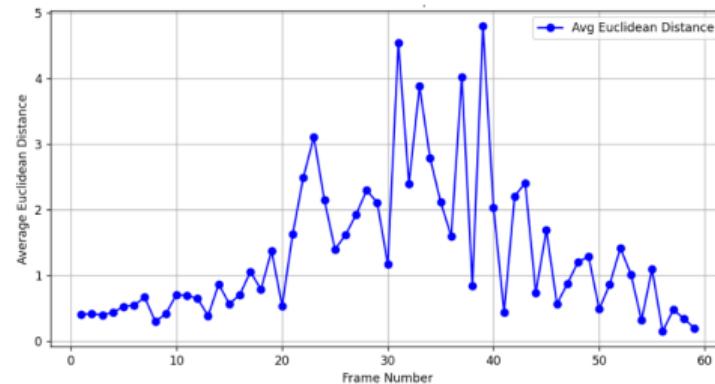


Figure 8: Detection of Speech Activity Graph

Given a set of lip landmark points, we normalize them such that point 48 becomes the origin  $(0,0)$ .

If  $48 \in \text{lip\_points}$ : 
$$\begin{cases} \text{origin} = \text{lip\_points}[48] \\ \text{normalized\_points}[i] = (x_i - x_{48}, y_i - y_{48}) \quad \forall i \in \text{lip\_points} \end{cases}$$

Otherwise, return the original points.

cont...

## Frame-wise Euclidean Distance

$$ED = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}, \quad (1)$$

Where  $(x_1, y_1)$  and  $(x_2, y_2)$  are the coordinates of a corresponding lip point in the input video and the reference video, respectively.

$$d_{ij} = \frac{\sum_{j=48}^{67} ED_{ij}}{20}, \quad (2)$$

$d_{ij}$  is the average distance for reference frame  $i$ .  $ED_{ij}$  is the Euclidean Distance between tow landmarks point. Where j is the lip landmark number and i is the ith frame of a reference video

cont...

## Video-wise Euclidean Distance

$$m_i = \frac{\sum_{j=1}^N d_{ij}}{N}, \quad (3)$$

- $m_i$  is the average distance for reference video  $i$ .
- $d_{ij}$  is the frame-wise average distance for frame  $j$  in reference video  $i$ .
- $N$  is the total number of frames in the video.

$$\text{predicted word} = \min(m_0, m_1, m_2, \dots, m_{n-1})$$

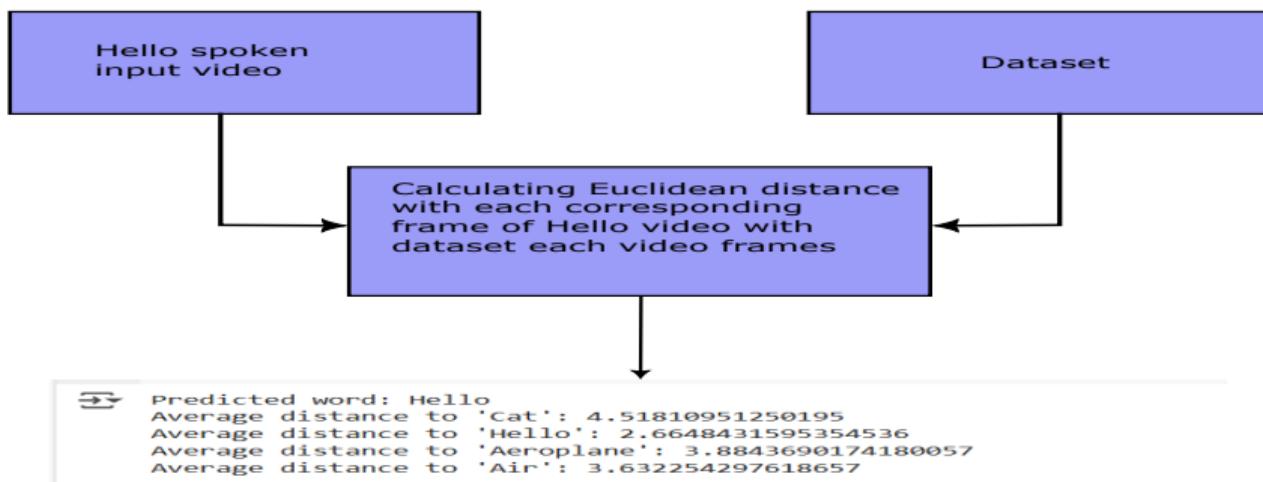


Figure 9: Calculating Euclidean Distance for the Entire Dataset

# Confusion Matrix

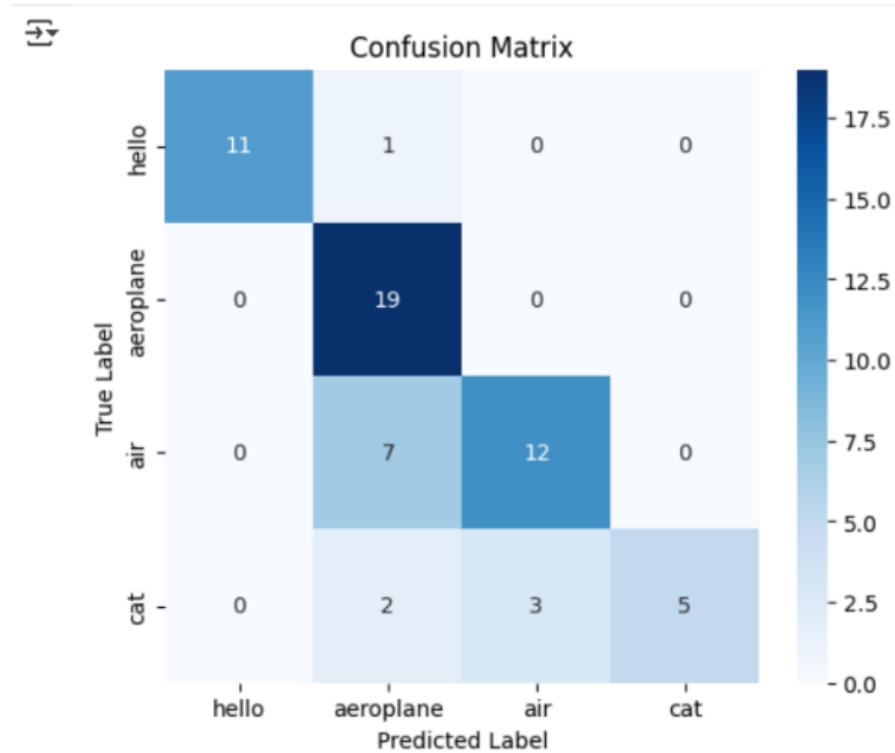


Figure 10: confusion matrix

- Successfully developed an algorithm which can recognize the spoken words by reading lips with partial use of Deep Learning.
- Our algorithm is only applied on Collins word dictionary and a custom dataset, which is created by me and accurately predicts 78.33% of words correctly. However, the accuracy is falling due to similar lip movements of different words(homophones).
- This algorithm can reduce the computational overhead as only deep learning is used for lip extraction from video, not for prediction. For word prediction, we are using our own developed algorithm.

- Acquire videos of the same word with different accents.
- Extract 3D lip feature so that accuracy can be achieved more easily
- Apply our algorithm to our dataset as well as publicly available datasets.
- Reading lips by using Deep Learning and comparing the result with Our Algorithm
- Finding the solution for Homophones
- Compare the results with the state-of-the-art.

## References |

- [1] Yannis M. Assael , Brendan Shillingford, Shimon Whiteson & Nando de Freitas *Lipnet: End-to-end sentence-level lipreading*, arXiv preprint 2016. DOI: 10.48550/arXiv.1611.01599
- [2] Yaman Kumar, Dhruva Sahrawat, Shubham Maheshwari, Debanjan Mahata, Amanda Stent, Yifang Yin, Rajiv Ratn Shah, Roger Zimmermann, *Harnessing GANs for Zero-Shot Learning of New Classes in Visual Speech Recognition*, Proceedings of the AAAI Conference on Artificial Intelligence, 2020, DOI: 10.1609/aaai.v34i03.5649
- [3] SOUHEIL FENGHOUR, (Associate Member, IEEE), DAQING CHEN, (Member, IEEE), KUN GUO2 AND PERRY XIAO *Lip Reading Sentences Using Deep Learning With Only Visual Cues*, IEEE, 2020, DOI: 10.1109/access.2020.3040906
- [4] He Wang, Pengcheng Guo, Wei Chen, Pan Zhou, Lei Xie *The NPU-ASLP-LiAuto System Description for Visual Speech Recognition in CNVSR* 2023,arXiv:2024, DOI: 10.48550/arXiv.2401.06788
- [5] Pingchuan Ma Stavros Petridis, Maja Pantic *Visual Speech Recognition for Multiple Languages in the Wild*, Nature Publishing Group UK London, 2022, DOI: 10.1038/s42256-022-00550-z
- [6] SOUHEIL FENGHOUR, (Associate Member, IEEE), DAQING CHEN, (Member, IEEE), KUN GUO AND PERRY XIAO *DISENTANGLING HOMOPHENES IN LIP READING USING PERPLEXITY ANALYSIS*, arXiv ,2020, DOI: 10.48550/arXiv.2012.07528

- [7] Jeong Hun Yeo, Chae Won Kim, Hyunjung Kim, Hyeongseop Rha, Seunghee Han, Wen-Huang Cheng, Yong Man Ro *Personalized Lip Reading: Adapting to Your Unique Lip Movements with Vision and Language*, arXiv , 2024, DOI:10.48550/arXiv.2409.00986
- [8] Pingchuan Ma, Yujiang Wang, Stavros Petridis, Jie Shen, Maja Pantic *Training Strategies for Improved Lip-Reading by Koohestani and Hadfield*, IEEE, 2022, DOI: 10.48550/arXiv.2209.01383
- [9] Ashish Vaswani 1, Noam Shazeer 2, Niki Parmar 3, Jakob Uszkoreit 4, Llion Jones 5, Aidan N. Gomez 6, Łukasz Kaiser 7, *Attention Is All You Need*, International Conference on Neural Information Processing Systems (NeurIPS), 2017, DOI: 10.48550/arXiv.1706.03762
  - Grid Corpus Dataset - MDPI Journal  
Source: *MDPI Applied Sciences*, 2021.
  - LRS2 Dataset - Oxford VGG Group  
Source: *Visual Geometry Group*, University of Oxford.
  - LRS3 Dataset - ResearchGate  
Source: *ResearchGate*, LRS3 Dataset Overview.
  - TCD-TIMIT Dataset - ResearchGate  
Source: *ResearchGate*, TCD-TIMIT Results Overview.
  - AVLetters Database - ResearchGate  
Source: *ResearchGate*, AVLetters Database Example.

- LRW Dataset - ResearchGate  
Source: *ResearchGate*, LRW Dataset Frames.
- Oulu-VS2 Dataset - ResearchGate  
Source: *ResearchGate*, Oulu-VS2 Dataset Examples.
- Collins Online Dictionary Word Pronunciation Videos  
Source: *Youtube*.

- **Aug 2024 - Sep 27, 2024:** Read research papers.
- **Sep 28 - Dec 13, 2024:** Created block diagram and prepared custom dataset.
- **Jan 2025 - Mar 13, 2025:** Developed the Methodology for the project.
- **Mar 13 - May 14, 2025:** Wrote and prepared paper for submission.

## Communicated

Visual Speech Recognition For Seamless Communication with hearing impaired persons

Md. Laraib Ahmad, Dr. Debbrota Paul Chowdhury,

*Multimedia Tools and Applications*, 2025.

Thank you for listening !

Md. Laraib Ahmad

[mdlaraib\\_pg\\_23@cse.nits.ac.in](mailto:mdlaraib_pg_23@cse.nits.ac.in)