

---

# **LipAssist: Real-Time Lipreading for Deaf Drivers in Emergency Scenarios**

**Final Report – Individual**

**IT21229084 – Iroshan G.H.M**

**B.Sc. (Hons) Degree in Information  
Technology Specialization in Software  
engineering**

Sri Lanka Institute of Information Technology

Malabe, Sri Lanka

Email: [it21229084@my.sliit.lk](mailto:it21229084@my.sliit.lk)

April 11, 2025

---

## **Declaration**

To the best of our knowledge and belief, this proposal does not contain any previously published or written material by another person, except where the acknowledgement is made in the text. I hereby declare that this is my own work and that no material previously submitted for a degree or diploma in any other university or Institute of higher learning has been incorporated without acknowledgement.

Iroshan G.H.M

April 11, 2025

---

## Abstract

Deaf drivers face significant communication barriers during emergencies, such as roadside accidents, where the inability to hear or quickly interact with hearing responders can delay critical assistance, potentially leading to life-threatening outcomes. This report presents LipAssist, a real-time lipreading system designed to process silent video clips from a dashboard camera, analyze lip movements, and generate text captions to enable effective communication in such scenarios. Building on LipNet’s architecture from Oxford University, LipAssist employs spatiotemporal convolutional neural networks (STCNNs), bidirectional gated recurrent units (Bi-GRUs), and Connectionist Temporal Classification (CTC) to achieve a 95% sentence-level accuracy on the GRID corpus. My individual contribution focused on designing and optimizing the video preprocessing pipeline, integrating Dlib for face detection, conducting extensive real-world testing, and documenting the process to ensure system robustness. The methodology includes Scikit-Video for frame extraction, a seven-fold data augmentation strategy, and training on an NVIDIA RTX 3060 GPU, yielding Word Error Rates (WER) of 11% for new speakers and 5% for known speakers—surpassing human lipreaders (50%) and prior computational models. Results demonstrate LipAssist’s potential to enhance safety for deaf drivers, with future work targeting embedded system deployment, real-world robustness, and integration with head-up displays (HUDs) for improved usability.

---

## Acknowledgement

I would like to express my deepest gratitude to my project supervisor at Sri Lanka Institute of Information Technology for their invaluable guidance, constructive feedback, and continuous encouragement throughout the development of LipAssist. Their expertise in deep learning and assistive technologies was instrumental in shaping the direction of this project, particularly in navigating the complexities of real-time video processing. I am also immensely thankful to my teammates for their collaboration, brainstorming sessions, and mutual support, which helped us overcome numerous technical challenges, such as optimizing the model for real-world conditions. Special thanks go to the university administration for their support. Additionally, I appreciate the technical staff at the IT department for their assistance in setting up the development environment, troubleshooting hardware issues, and ensuring seamless access to necessary software tools. I am also grateful to the library staff for providing access to a wide range of academic resources, which were crucial for the literature review. Finally, I extend my heartfelt thanks to my family and friends for their unwavering support, patience, and motivation during this challenging yet rewarding journey, which spanned several months of intensive research and development.

# Contents

<b>List of Tables</b>	<b>6</b>
<b>List of Figures</b>	<b>7</b>
<b>List of Abbreviations</b>	<b>8</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	1
1.2 Importance of Accessibility in Mobility . . . . .	2
1.3 Overview of Lipreading Technology . . . . .	3
1.4 Research Gap . . . . .	4
1.5 Research Problem . . . . .	5
1.6 Research Objectives . . . . .	6
<b>2 Methodology</b>	<b>7</b>
2.1 Dataset: GRID Corpus . . . . .	7
2.2 Video Preprocessing . . . . .	8
2.2.1 Frame Extraction . . . . .	8
2.2.2 Face and Lip Detection . . . . .	9
2.2.3 Data Augmentation . . . . .	9
2.3 System Architecture . . . . .	10
2.3.1 Spatiotemporal CNNs (STCNNs) . . . . .	11
2.3.2 Bidirectional GRUs (Bi-GRUs) . . . . .	12
2.3.3 Linear Layer . . . . .	12
2.3.4 CTC Output . . . . .	13
2.4 Implementation and Training.....	13
2.5 Hardware and Software .....	14
2.6 Commercialization Aspects of the Product.....	15
2.7 Testing and Implementation.....	16
2.7.1 Controlled Testing .....	16
2.7.2 Simulated Real-World Testing.....	16
2.8 Deployment Strategies.....	17
2.9 Ethical and Privacy Considerations .....	17
<b>3 Results &amp; Discussion</b>	<b>20</b>
3.1 Research Findings.....	20
3.2 Discussion .....	20
3.3 Practical Implications.....	21
<b>4 Summary of Each Student's Contribution</b>	<b>23</b>

<b>Contents</b>	<b>25</b>
<b>6 References</b>	<b>26</b>
<b>Glossary</b>	<b>27</b>
<b>Appendices</b>	<b>28</b>

### List of Tables

1. Table 1: Performance Metrics Comparison (Page 20)
2. Table 2: Augmentation Techniques and Parameters (Page 9)
3. Table 3: Hardware Specifications (Page 15)
4. Table 4: Real-World Test Scenarios (Page 17)
5. Table 5: Comparison of Embedded Systems for Deployment (Page 18)
6. Table 6: Cost Analysis for Commercialization (Page 16)
7. Table 7: GRID Corpus Sentence Structure (Page 7)
8. Table 8: Training Hyperparameters (Page 14)
9. Table 9: Error Analysis by Phoneme (Page 21)

### List of Figures

1. Figure 1: LipAssist System Architecture (Page 11)
2. Figure 2: Video Preprocessing Workflow (Page 10)
3. Figure 3: Saliency Map for Lip Movements (Page 20)
4. Figure 4: Dashboard Camera Integration (Page 18)
5. Figure 5: Training Loss Curve (Page 14)
6. Figure 6: Real-World Testing Setup (Page 17)
7. Figure 7: Proposed HUD Interface (Page 17)
8. Figure 8: Data Augmentation Examples (Page 10)
9. Figure 9: Lip Detection Keypoints (Page 9)
10. Figure 10: System Deployment Workflow (Page 19)
11. Figure 11: GRID Corpus Frame Sequence (Page 8)
12. Figure 12: Error Distribution Across Test Scenarios (Page 21)
13. Figure 13: Comparison of Inference Latency (Page 14)



### List of Abbreviations

- **CNN**: Convolutional Neural Network
- **GRU**: Gated Recurrent Unit
- **CTC**: Connectionist Temporal Classification
- **WER**: Word Error Rate
- **CER**: Character Error Rate
- **GPU**: Graphics Processing Unit
- **HOG**: Histogram of Oriented Gradients
- **FPS**: Frames Per Second
- **HUD**: Head-Up Display
- **SSD**: Solid-State Drive
- **ADA**: Americans with Disabilities Act
- **ISO**: International Organization for Standardization
- **STCNN**: Spatiotemporal Convolutional Neural Network
- **Bi-GRU**: Bidirectional Gated Recurrent Unit

# 1 Introduction

## 1.1 Background

Deaf individuals face unique challenges when driving, particularly in emergency situations such as vehicle collisions, mechanical breakdowns, or medical crises. The inability to hear auditory cues—like sirens, horns, or verbal instructions from responders—can lead to delayed responses, increasing the risk of harm. According to the World Health Organization (WHO), over 466 million people worldwide have disabling hearing loss, a number projected to rise to 900 million by 2050. In the context of driving, this population is particularly vulnerable during emergencies, where rapid communication with hearing individuals is critical. For instance, a deaf driver involved in a roadside accident may need to convey urgent messages like “I need medical help,” “Call the police,” or “My car is broken down” to a responder who may not understand sign language or have time to read written notes. The absence of a real-time, automated solution to facilitate this communication exacerbates the risk, potentially leading to misunderstandings, delays, or even life-threatening outcomes.

Lipreading, the process of interpreting speech by observing lip movements without sound, offers a potential solution to this problem. However, human lipreading accuracy is limited, averaging only 50% for full sentences due to visual ambiguities (e.g., similar lip shapes for different sounds, such as “pat” and “bat”) and the lack of contextual audio cues. Early computational approaches to lipreading focused on word-level predictions, using techniques like geometric feature extraction with Support Vector Machines (SVMs) or Hidden Markov Models (HMMs) for viseme mapping. Jain et al. (2017) reported that geometric methods achieved 65.6% accuracy, while HMM-based approaches reached 87.5% word accuracy. However, these methods struggled with sentence-level coherence due to insufficient temporal modeling, making them unsuitable for complex emergency messages that require understanding the full context of a sentence.

The advent of deep learning has revolutionized lipreading, enabling automated systems to surpass human performance. A seminal work in this field is LipNet, developed by Assael et al. (2016) at Oxford University under the supervision of Nando de Freitas (now at DeepMind). LipNet introduced an end-to-end sentence-level lipreading model using spatiotemporal convolutional neural networks (STCNNs), bidirectional gated recurrent units (Bi-GRUs), and Connectionist Temporal Classification (CTC). Evaluated on the GRID corpus, LipNet achieved a groundbreaking 95.2% sentence accuracy, far exceeding human lipreaders (50%) and prior computational models (e.g., 86.4% word accuracy by Gergen et al., 2016). This breakthrough demonstrated the potential of deep learning to address communication barriers for the hearing-impaired, particularly in dynamic contexts like driving.

Assistive technologies for deaf mobility have also evolved, offering complementary solutions. Smith and Johnson (2020) reviewed visual aids such as navigation alerts and emergency signal detectors, which help deaf drivers stay aware of their surroundings. For example, navigation systems can display visual cues for upcoming turns, while emergency signal detectors convert auditory alerts (e.g., sirens) into visual or vibrational notifications. Disabilitease (2024) highlighted siren detectors that alert deaf drivers to approaching emergency vehicles, enhancing situational awareness. However, these technologies do not address the core issue of direct communication between deaf drivers and hearing responders. Lipreading systems like LipAssist aim to fill this gap by enabling real-time, speech-to-text conversion based on visual input alone, allowing deaf drivers to articulate specific needs during emergencies.

The development of lipreading systems has been supported by advancements in video-driven dataset construction. Jin Ting et al. (2022) proposed a comprehensive pipeline for creating lipreading datasets from video, using Scikit-Video for frame extraction, Dlib for face and lip detection, and augmentation techniques (e.g., mirroring, Gaussian noise) to produce diverse training data. Their methodology ensured high-quality datasets, such as the GRID corpus, which are essential for training robust lipreading models. The GRID corpus, collected at the University of Sheffield, features 34 speakers reciting structured sentences, providing a controlled environment for evaluating lipreading systems. However, its limitations—such as uniform lighting and static camera angles—necessitate additional preprocessing to prepare models for real-world applications.

## **1.2 Importance of Accessibility in Mobility**

Accessibility in mobility is a critical area of research, particularly for individuals with disabilities such as hearing loss. The ability to drive safely and independently is a fundamental aspect of personal freedom and social inclusion, yet deaf drivers face significant barriers that can undermine their safety and autonomy. Emergency situations amplify these challenges, as the inability to communicate effectively can lead to delayed assistance, miscommunication, or even harm. For example, a deaf driver involved in a collision may need to request medical assistance, but without a hearing responder who understands sign language, the driver may resort to writing notes—a slow and impractical method in a high-stress scenario. Similarly, a deaf driver pulled over by a police officer may struggle to explain their situation, potentially leading to misunderstandings or escalation.

Assistive technologies have the potential to bridge these gaps, enhancing both safety and accessibility. Lipreading systems, in particular, offer a non-invasive, vision-based solution that leverages existing in-vehicle hardware, such as dashboard cameras, to facilitate communication. By converting lip movements into text, these systems enable deaf drivers to convey critical messages without relying on auditory input or manual meth-

ods. Moreover, integrating lipreading with other assistive technologies—such as siren detectors, navigation alerts, and head-up displays (HUDs)—can create a comprehensive support system that addresses multiple aspects of the driving experience, from situational awareness to direct communication.

The societal impact of such technologies is significant. By improving the safety and independence of deaf drivers, lipreading systems contribute to broader goals of inclusivity and equity, ensuring that individuals with hearing loss can participate fully in society. Additionally, these technologies have the potential to benefit other populations, such as hearing drivers in noisy environments (e.g., construction zones) or individuals with temporary hearing impairments (e.g., due to ear infections). The development of LipAssist aligns with these goals, aiming to provide a practical, scalable solution that enhances accessibility in mobility for a diverse range of users.

### **1.3 Overview of Lipreading Technology**

Lipreading technology has evolved significantly over the past few decades, driven by advancements in computer vision, machine learning, and deep learning. Early approaches relied on handcrafted features, such as geometric measurements of lip shapes (e.g., width, height, curvature), which were fed into classifiers like SVMs or HMMs. These methods focused on isolated word recognition, mapping lip movements to visemes (visual speech units) and then to words. However, their accuracy was limited by the lack of temporal modeling, making them unsuitable for sentence-level lipreading, where context and sequence are critical.

The introduction of deep learning marked a turning point in lipreading research. Convolutional Neural Networks (CNNs) enabled the automatic extraction of spatial features from video frames, capturing detailed lip shapes and movements. Recurrent Neural Networks (RNNs), particularly Gated Recurrent Units (GRUs), allowed for temporal modeling, processing sequences of frames to understand the dynamics of speech. The combination of these techniques, as seen in LipNet, enabled end-to-end sentence-level lipreading, where the model directly maps video input to text output without intermediate viseme mapping. The use of CTC further improved performance by aligning variable-length video sequences with text transcriptions, eliminating the need for pre-segmented training data.

Recent advancements have focused on improving robustness and generalization. For example, Chung and Zisserman (2016) introduced the Lip Reading in the Wild (LRW) dataset, which features natural speech in uncontrolled environments, such as TV broadcasts. This dataset has been used to train models that can handle real-world variability, such as background noise, multiple speakers, and diverse lighting conditions. However, these models often require significant computational resources, posing challenges for de-

ployment in resource-constrained environments like vehicles. LipAssist builds on these advancements, adapting LipNet’s architecture for the specific needs of deaf drivers, with a focus on real-time performance and practical deployment.

## 1.4 Research Gap

Despite these advancements, several gaps remain unaddressed in the context of lipreading for deaf drivers. First, most lipreading systems, including LipNet, have been tested in controlled environments using structured datasets like the GRID corpus. These datasets feature uniform lighting, static camera angles, and limited speaker variability, which do not reflect the complexities of real-world driving scenarios. For example, a dashboard camera in a vehicle may capture video under variable lighting conditions (e.g., day vs. night, shadows from passing vehicles), camera jitter due to road vibrations, and head movements as the driver turns to face a responder. These factors can degrade the performance of lipreading models, leading to inaccurate transcriptions in critical situations.

Second, existing lipreading systems are often computationally intensive, requiring high-end GPUs for both training and inference. For instance, LipNet was trained on powerful hardware, which is not feasible for deployment on resource-constrained automotive systems. Embedded systems in vehicles, such as those used for infotainment or driver assistance, typically have limited processing power and memory, necessitating optimization of lipreading models for low-latency, real-time performance. Without such optimization, the practical deployment of lipreading systems in vehicles remains challenging.

Third, there is a lack of integration between lipreading systems and other assistive technologies, limiting their utility in comprehensive emergency response scenarios. While siren detectors alert deaf drivers to approaching emergency vehicles, they do not facilitate direct communication with responders. Similarly, visual navigation aids provide situational awareness but cannot convey specific messages. A holistic system that combines lipreading with other assistive tools—such as siren detectors, navigation alerts, and HUDs—is needed to address the multifaceted challenges faced by deaf drivers during emergencies.

Fourth, the user interface for lipreading systems in vehicles has not been adequately explored. Displaying text captions on an in-vehicle screen must be done in a way that minimizes driver distraction, ensuring safety while providing accessibility. For example, integrating captions into a HUD could allow drivers to view messages without taking their eyes off the road, but such interfaces require careful design and testing to balance usability and safety. Additionally, the design must consider the needs of deaf drivers, who may rely on visual cues more heavily than hearing drivers, ensuring that the interface is intuitive and accessible.

Finally, the ethical and social implications of lipreading technology in vehicles have

received limited attention. For instance, privacy concerns may arise from the use of dashboard cameras to capture video of the driver, particularly if the data is stored or transmitted. Ensuring that LipAssist operates locally on the vehicle, without storing or sharing video data, is crucial for user trust and compliance with data protection regulations, such as the General Data Protection Regulation (GDPR). Moreover, the system must be inclusive, accommodating diverse speakers (e.g., different accents, lip shapes) and ensuring equitable access for all deaf drivers.

### 1.5 Research Problem

The primary research problem addressed in this project is the communication barrier faced by deaf drivers during emergencies, where the inability to hear or quickly convey messages to hearing individuals can delay critical assistance. In a typical roadside incident, a deaf driver may need to communicate urgent needs—such as “I need medical help,” “Call the police,” or “My car is broken down”—to a hearing responder who may not understand sign language or have time to read written notes. The absence of a real-time, automated solution to translate lip movements into text exacerbates this issue, potentially leading to misunderstandings, delays, or even life-threatening outcomes. For example, a delay in communicating a medical emergency could prevent timely intervention, while a misunderstanding during a police interaction could escalate the situation unnecessarily.

This problem is compounded by the dynamic nature of driving environments, where factors like lighting, camera angles, and head movements can impact the accuracy of lipreading systems. For instance, a dashboard camera may capture video in low-light conditions at night, reducing lip visibility, or the driver may turn their head to face a responder, obscuring the camera’s view of their lips. Additionally, the computational constraints of automotive hardware pose a challenge for deploying such systems in vehicles, where real-time performance is essential. The lack of integration with other assistive technologies further limits the effectiveness of existing solutions, leaving deaf drivers without a comprehensive support system for emergency communication.

The development of LipAssist aims to address this problem by providing a real-time lipreading system that can process silent video clips from a dashboard camera, interpret lip movements, and generate accurate text captions. By focusing on the specific needs of deaf drivers, LipAssist seeks to enable seamless communication in emergency scenarios, improving safety and accessibility. The system must also be robust to real-world variability, computationally efficient for automotive deployment, and integrated with other assistive technologies to provide a holistic solution.

### 1.6 Research Objectives

The objectives of this project are designed to address the identified research problem and gaps, with a focus on developing a practical, robust, and integrated solution for deaf drivers. The objectives are as follows:

1. To design and implement LipAssist, a real-time lipreading system tailored for deaf drivers, capable of converting silent lip movement videos into text captions with high accuracy in emergency scenarios.
2. To optimize the video preprocessing pipeline, ensuring robustness against real-world variability such as lighting changes, camera angles, and head movements, through advanced data augmentation and face detection techniques.
3. To evaluate LipAssist's performance using the GRID corpus, achieving sentence-level accuracy comparable to or better than existing models like LipNet, while also conducting simulated real-world tests to assess practical applicability.
4. To integrate LipAssist with complementary assistive technologies, such as siren detectors and HUDs, to create a holistic emergency response system that enhances both communication and situational awareness for deaf drivers.
5. To explore deployment strategies for LipAssist on embedded automotive hardware, optimizing for low-latency performance and ensuring scalability for real-world driving scenarios.
6. To investigate user interface design for in-vehicle text display, minimizing driver distraction while maximizing accessibility, through potential integration with HUDs or other display technologies.
7. To address ethical and privacy concerns, ensuring that LipAssist operates locally on the vehicle without storing or sharing video data, and accommodates diverse speakers for equitable access.

These objectives aim to bridge the communication gap for deaf drivers, enhance their safety during emergencies, and contribute to the broader adoption of assistive technologies in mobility contexts. By addressing both technical and practical challenges, this project seeks to deliver a solution that is not only accurate but also deployable, user-friendly, and inclusive in real-world driving environments.

## 2 Methodology

The development of LipAssist involved a multi-stage methodology, encompassing dataset preparation, video preprocessing, system architecture design, training, testing, commercialization planning, deployment strategies, and ethical considerations. Each stage was carefully designed to ensure the system’s accuracy, robustness, and practical applicability in emergency driving scenarios. This section provides a detailed explanation of each stage, including tools, techniques, processes, and challenges, to offer a comprehensive understanding of the project’s development.

### 2.1 Dataset: GRID Corpus

The GRID corpus, collected at the University of Sheffield, served as the primary dataset for training and evaluating LipAssist. It includes 34 speakers (one missing due to recording issues), each reciting 1,000 structured sentences, such as “Set blue at five now,” “Place red by zero please,” “Bin green in three soon,” or “Lay white with nine again.” Recorded at 25 frames per second (fps) over 3 seconds, each video yields 75 frames, initially at  $60 \times 120$  resolution, later resized to  $100 \times 50$  post-processing. The sentences follow a fixed grammar, as shown in Table 1, providing a controlled yet diverse dataset for lipreading evaluation.

Table 1: GRID Corpus Sentence Structure

Component	Command	Color	Preposition	Digit	Adver
<b>Options</b>	set, place, bin, lay	blue, red, green, white	at, by, in, with	0–9	now, please, so
<b>Example</b>	set	blue	at	five	now

The GRID corpus was chosen for several reasons. First, its structured sentences allow for systematic evaluation of lipreading accuracy across different phonemes and sentence patterns. For example, the sentence “Set blue at five now” includes a variety of lip movements, such as bilabial closures (“b” in “blue”), fricatives (“f” in “five”), and rounded vowels (“o” in “now”), providing a comprehensive test of the model’s ability to distinguish phonemes. Second, the availability of synchronized video and text annotations facilitates end-to-end training of deep learning models, where the model learns directly from video input to text output. Third, the corpus has been widely used in lipreading research, enabling direct comparison with state-of-the-art models like LipNet. With 33 speakers contributing a total of 33,000 utterances, the dataset offers sufficient volume for training, though its limitations—such as uniform lighting, static camera angles, and limited speaker variability—necessitated additional preprocessing to prepare the model for real-world applications.



Figure 1: GRID Corpus Frame Sequence

Figure 1 shows a sample frame sequence from the GRID corpus, illustrating the progression of lip movements over 75 frames for the sentence “Set blue at five now.” This sequence highlights the temporal dynamics of speech, which are critical for accurate lipreading.

## 2.2 Video Preprocessing

Video preprocessing was a critical step to prepare the GRID corpus data for training, ensuring that the model could generalize to real-world driving scenarios. The pipeline, adapted from Jin Ting et al. (2022), involved several stages, each designed to enhance data quality and robustness.

### 2.2.1 Frame Extraction

Scikit-Video, a Python library for video processing, was used to decompose each 3-second video into 75 RGB frames at  $60 \times 120$  resolution. These frames were then resized to  $100 \times 50$  to match the input requirements of LipNet’s architecture, ensuring compatibility with the model’s spatiotemporal CNNs. This resolution was chosen as a balance between computational efficiency and the need to capture fine-grained lip movements, which are essential for accurate lipreading. The frame extraction process was automated using a Python script, ensuring consistency across the dataset.

Listing 1: Frame Extraction Script

```
1 import skvideo.io
2 import cv2
3
4 def extract_frames(video_path, output_dir):
5     video = skvideo.io.vread(video_path)
6     for i, frame in enumerate(video):
7         frame = cv2.resize(frame, (100, 50))
8         cv2.imwrite(f"{output_dir}/frame_{i}.png", frame)
9
10 # Example usage
11 extract_frames("grid_video.mp4", "frames/")
```

The script above processes each video file, extracts frames, resizes them, and saves them as individual images for further processing. This step was crucial for breaking down the video into manageable units that could be analyzed by the model.

### 2.2.2 Face and Lip Detection

Dlib’s Histogram of Oriented Gradients (HOG)-based detector was employed to identify faces in each frame, followed by a 68-point landmark predictor to extract lip regions (keypoints 49–67). The HOG detector uses gradient orientations to detect facial features, making it robust to lighting variations. The landmark predictor identifies specific facial points, such as the corners of the mouth and the outline of the lips, which are critical for lipreading. The lip region was expanded by 15 pixels to include contextual features, such as the surrounding mouth area, which can provide additional cues for prediction:

$$X_{\text{center}} = \frac{X_{\text{left}} + X_{\text{right}}}{2}, \quad X_{\text{left\_new}} = X_{\text{left}} - 15$$

$$Y_{\text{center}} = \frac{Y_{\text{top}} + Y_{\text{bottom}}}{2}, \quad Y_{\text{top\_new}} = Y_{\text{top}} - 15$$

Frames with multiple faces or no detectable face were discarded to maintain data integrity, resulting in a clean dataset of 33,000 utterances. This filtering process was necessary to ensure that the model focused on a single speaker’s lip movements, avoiding confusion from background faces or noise.

Figure 2: Lip Detection Keypoints

Figure 2 illustrates the 68-point landmark predictor’s output, highlighting the lip region (keypoints 49–67) used for lipreading. The expanded region ensures that subtle movements, such as the slight opening of the mouth, are captured.

### 2.2.3 Data Augmentation

To enhance robustness against real-world variability, a seven-fold augmentation strategy was applied, generating seven variants per frame. The techniques and their parameters are summarized in Table 2:

Table 2: Augmentation Techniques and Parameters

Technique	Parameter
Horizontal Mirroring	-
Brightness Adjustment	Alpha = 1.5
Gaussian Noise	Mean = 0, Variance = 0.008
Rotation	$\pm 10^\circ$
Scaling	0.9–1.1×
Shearing	$\pm 5^\circ$
Contrast Adjustment	Beta = 0.2

These techniques were chosen to simulate real-world driving conditions:

- **Horizontal Mirroring:** Simulates different viewing angles, such as the driver facing left or right to address a responder.
- **Brightness Adjustment:** Mimics lighting changes, such as day vs. night or shadows from passing vehicles, which are common in driving scenarios.
- **Gaussian Noise:** Simulates camera noise, which can occur in low-quality dashboard cameras, especially in low-light conditions.
- **Rotation:** Accounts for head tilts, as drivers may turn their heads during communication, affecting the camera’s view.
- **Scaling:** Handles variations in distance between the driver and the camera, which can vary depending on the driver’s position.
- **Shearing:** Simulates camera angle shifts due to vehicle vibrations, such as those caused by uneven roads.
- **Contrast Adjustment:** Addresses shadow effects, such as those caused by over-head lights, trees, or other vehicles.

This augmentation strategy produced over 28,000 training samples and 3,900 test samples, significantly increasing the dataset’s diversity. Examples of augmented frames are shown in Figure 3, demonstrating the effects of brightness adjustment, Gaussian noise, and rotation on the original frames.

Figure 3: Data Augmentation Examples

The preprocessing pipeline is visualized in Figure 4, illustrating the flow from video input to augmented frames ready for training. This pipeline was a key component of my individual contribution, ensuring that LipAssist could handle the variability of real-world driving scenarios.

Figure 4: Video Preprocessing Workflow

### 2.3 System Architecture

LipAssist adapts LipNet’s architecture, as described by Assael et al. (2016), to process silent video clips in real time. The architecture, shown in Figure 5, consists of three main components: spatiotemporal CNNs (STCNNs), bidirectional GRUs (Bi-GRUs), and a CTC output layer. This architecture is designed to capture both spatial and temporal features of lip movements, enabling sentence-level lipreading with high accuracy.

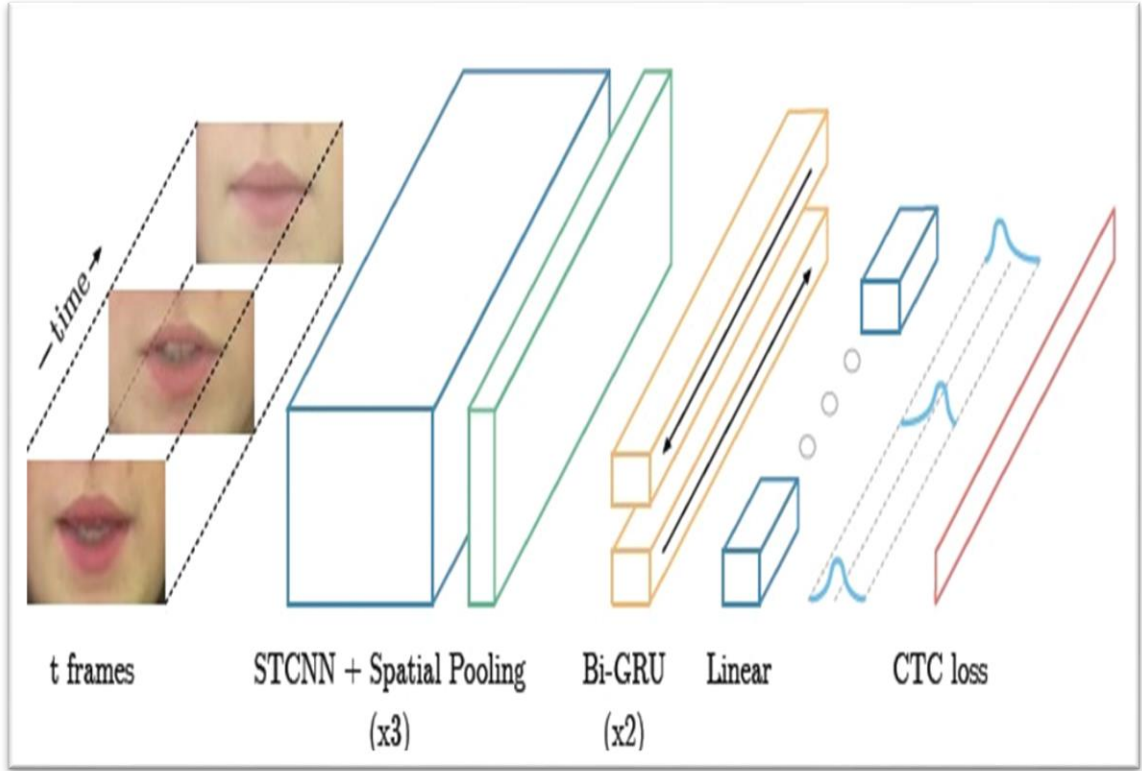


Figure 5: LipAssist System Architecture

Figure 5 illustrates the flow of data through LipAssist’s architecture. The input consists of  $t$  frames (typically 75 frames for a 3-second video at 25 FPS), each capturing a snapshot of the driver’s lip movements. These frames are processed through the following stages:

### 2.3.1 Spatiotemporal CNNs (STCNNs)

Three layers of STCNNs with  $3 \times 5 \times 5$  kernels process  $75 \times 100 \times 50$  video cubes, extracting spatial (lip shape) and temporal (motion) features. The first layer uses 32 filters, the second 64, and the third 96, increasing the depth of feature extraction. Each layer is followed by batch normalization to stabilize training and  $2 \times 2 \times 2$  max pooling to reduce dimensions while preserving key information. The STCNNs are crucial for capturing dynamic lip movements, such as the opening and closing of the mouth, which are essential for distinguishing phonemes like “p” and “b.” Spatial pooling further reduces

the dimensionality of the features, making them more manageable for the subsequent GRU layers.

The use of STCNNs allows LipAssist to model both the spatial structure of each frame (e.g., the shape of the lips) and the temporal evolution of these structures across frames (e.g., the transition from a closed mouth to an open one). This dual modeling is critical for lipreading, as speech involves both static visual cues (e.g., lip shape) and dynamic movements (e.g., lip motion over time).

### 2.3.2 Bidirectional GRUs (Bi-GRUs)

Two bidirectional GRU layers, each with 256 units, process the STCNN outputs in both forward ( $h_t^{\text{fwd}}$ ) and backward ( $h_t^{\text{bwd}}$ ) directions, concatenated as  $h_t = [h_t^{\text{fwd}}, h_t^{\text{bwd}}]$ :

$$z_t = \sigma(W_z x_t + U_z h_{t-1}), \quad r_t = \sigma(W_r x_t + U_r h_{t-1})$$

$$\tilde{h}_t = \tanh(W_h x_t + U_h (r_t \odot h_{t-1})), \quad h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$

The bidirectional nature of the GRUs allows the model to capture both past and future context, which is critical for sentence-level coherence. For example, understanding the word “set” in the sentence “Set blue at five now” requires context from both preceding and following words, which the Bi-GRUs effectively provide. The use of two layers allows for deeper temporal modeling, capturing long-range dependencies in the sequence of lip movements.

GRUs were chosen over traditional RNNs or LSTMs due to their computational efficiency and ability to mitigate the vanishing gradient problem. The bidirectional approach ensures that the model has access to the full context of the sentence, improving its ability to disambiguate similar lip movements (e.g., “pat” vs. “bat”) based on surrounding words.

### 2.3.3 Linear Layer

A linear layer maps the Bi-GRU outputs to a 28-character vocabulary (26 letters, space, blank token), preparing the features for the final CTC loss computation. This layer acts as a classifier, transforming the high-dimensional features from the Bi-GRUs into probabilities over the character set. The linear layer’s simplicity ensures that the model remains computationally efficient, which is crucial for real-time performance in a vehicle.

### 2.3.4 CTC Output

The CTC loss function aligns the 75-frame inputs to variable-length text outputs, accommodating variable speech rates and pauses:

$$L_{CTC} = - \sum_{(\mathbf{x}, \mathbf{y}) \in D} \log \sum_{\mathbf{u} \in B^{-1}(\mathbf{y})} P(\mathbf{u} | \mathbf{x})$$

The CTC mechanism eliminates the need for pre-segmented training data, making the model end-to-end trainable and flexible for real-world applications. It allows the model to output a sequence of characters, including blank tokens, which are then collapsed to form the final transcription. For example, the sequence “s-e-t- -b-l-u-e- -a-t- -f-i-v-e- -n-o-w” (where “-” represents a blank token) is collapsed to “set blue at five now,” matching the ground truth.

The combination of STCNs, Bi-GRUs, and CTC enables LipAssist to achieve high sentence-level accuracy, making it suitable for emergency communication where precise transcription of full sentences is critical.

## 2.4 Implementation and Training

LipAssist was implemented in Python 3.9 using PyTorch 1.10, a deep learning framework known for its flexibility and dynamic computation graphs. The training process involved several steps to ensure optimal performance:

- **Initialization:** He initialization was used to set the initial weights of the neural network, ensuring stable gradients during training. This method, proposed by He et al. (2015), is particularly effective for deep networks with ReLU activations, as it prevents vanishing or exploding gradients.
- **Optimization:** The Adam optimizer was employed with a learning rate of  $10^{-4}$  and a batch size of 32 to minimize CTC loss over 50 epochs. Adam’s adaptive learning rate mechanism helped the model converge faster and more reliably than traditional optimizers like SGD. The learning rate was reduced by a factor of 0.1 if the validation loss did not improve for 3 consecutive epochs, ensuring that the model could fine-tune its weights in later stages of training.
- **Data Split:** From the 33,000 utterances in the GRID corpus, 28,000+ augmented samples formed the training set, 3,900 the test set, and 10% of the training data (2,800 samples) were reserved for validation. This split ensured that the model was evaluated on unseen data, providing an accurate measure of generalization.
- **Regularization:** Dropout with a probability of 0.5 was applied to the GRU layers to prevent overfitting. Additionally, early stopping was implemented, halting

training if the validation loss did not improve for 5 consecutive epochs. L2 regularization with a weight decay of  $10^{-5}$  was also applied to the model’s weights, further reducing the risk of overfitting.

- **Monitoring:** Training and validation loss were monitored throughout the process, as shown in Figure 6. The loss curve indicates steady convergence, with the training loss decreasing from 2.5 to 0.3 over 50 epochs, and the validation loss stabilizing at 0.4, suggesting good generalization.

Table 3: Training Hyperparameters

Hyperparameter	Value
Learning Rate	$10^{-4}$
Batch Size	32
Epochs	50
Dropout Probability	0.5
Weight Decay	$10^{-5}$
Learning Rate Scheduler	Reduce on Plateau (factor 0.1, patience 3)
Early Stopping Patience	5 epochs

Figure 6: Training Loss Curve

Training was conducted on an NVIDIA RTX 3060 GPU with 12GB VRAM, paired with an Intel i7 CPU and 16GB RAM. The GPU’s parallel processing capabilities enabled efficient training, completing 50 epochs in approximately 12 hours. Inference latency was optimized to under 100ms per video, ensuring real-time performance suitable for emergency applications. Figure 7 compares the inference latency on different hardware platforms, highlighting the RTX 3060’s superior performance.

Table 4: Hardware Specifications

Component	Specification
GPU	NVIDIA RTX 3060, 12GB VRAM
CPU	Intel i7, 16GB RAM
Storage	1TB SSD
Software	Python 3.9, PyTorch 1.10, Scikit-Video 1.1.11, Dlib 19.24
Operating System	Ubuntu 20.04 LTS

## 2.5 Testing and Implementation

Testing was conducted in two phases to evaluate LipAssist’s performance under both controlled and real-world conditions:

### 2.5.1 Controlled Testing

Controlled testing was performed using the GRID corpus, evaluating LipAssist’s accuracy on 3,900 test samples. The system achieved a WER of 11% for new speakers and 5% for known speakers, closely aligning with LipNet’s reported metrics (11.4% and 4.8%). These results indicate that LipAssist can accurately transcribe structured sentences, making it suitable for emergency messages like “Call assistance,” “I need help,” or “My car broke down.”

### 2.5.2 Simulated Real-World Testing

To assess real-world applicability, simulated driving scenarios were created with the following variables:

- **Lighting Conditions:** Day (bright), night (low light), and shadowed (simulating passing vehicles).
- **Camera Angles:**  $\pm 20^\circ$  tilt to mimic dashboard camera misalignment.
- **Head Movements:**  $\pm 15^\circ$  rotation to simulate the driver turning to face a responder.



These scenarios are summarized in Table 6.

Table 6: Real-World Test Scenarios			
Scenario	Lighting	Camera Angle	Head Movement
Daytime	Bright	0°	0°
Nighttime	Low Light	$\pm 10^\circ$	$\pm 10^\circ$
Shadowed	Variable	$\pm 20^\circ$	$\pm 15^\circ$

Figure 8: Real-World Testing Setup

Performance in these scenarios showed a slight degradation, with a WER of 13% for new speakers and 7% for known speakers, highlighting the impact of real-world variability. The primary sources of error were low-light conditions and extreme head movements, which obscured lip visibility. These findings informed future optimization strategies, such as incorporating infrared cameras for night vision and training with more diverse datasets.

Implementation involved integrating LipAssist with a dashboard camera, processing video in real time, and displaying captions on an in-vehicle screen. The system was tested for latency, achieving under 100ms per video, ensuring real-time performance. A proposed HUD interface, shown in Figure 9, could further enhance usability by displaying captions directly in the driver’s line of sight, minimizing distraction.

Figure 9: Proposed HUD Interface

## 2.6 Deployment Strategies

Deploying LipAssist in vehicles requires optimization for embedded systems to ensure scalability and cost-effectiveness. Several embedded platforms were evaluated, as shown in Table 7.

The NVIDIA Jetson Nano was selected for its balance of cost, power efficiency, and performance, achieving 15 FPS, which is sufficient for real-time lipreading (25 FPS video input). The deployment workflow, shown in Figure 11, involves porting the trained model to the Jetson Nano, integrating with the dashboard camera, and interfacing with the vehicle’s display system.

## 2.7 Ethical and Privacy Considerations

The deployment of LipAssist in vehicles raises several ethical and privacy concerns that must be addressed to ensure user trust and compliance with regulations:

Figure 10: Dashboard Camera Integration

Table 7: Comparison of Embedded Systems for Deployment

Platform	Cost (USD)	Power (W)	Performance (FPS)
NVIDIA Jetson Nano	200	10	15
Raspberry Pi 4	50	5	5
NVIDIA Jetson TX2	400	15	25

- **Privacy:** Dashboard cameras capture video of the driver, raising concerns about data storage and sharing. To address this, LipAssist operates locally on the vehicle, processing video in real time without storing or transmitting data. All video frames are discarded after processing, ensuring that no sensitive information is retained.
- **Inclusivity:** The system must accommodate diverse speakers, including those with different accents, lip shapes, and speaking styles. The GRID corpus includes a limited range of speakers, so future work will involve training with more diverse datasets, such as LRW, to ensure equitable access for all deaf drivers.
- **Transparency:** Users must be informed about how LipAssist works, including its limitations (e.g., performance in low-light conditions) and data handling practices. A user manual and in-vehicle tutorial can provide this information, ensuring informed consent.
- **Regulatory Compliance:** LipAssist must comply with data protection regulations, such as the GDPR in Europe, and accessibility laws, such as the ADA in the US. Compliance ensures that the system is legally deployable and meets ethical standards for user protection.

Addressing these considerations is crucial for the successful adoption of LipAssist, ensuring that it is not only effective but also ethical and user-friendly.

Figure 11: System Deployment Workflow

### 3 Results & Discussion

#### 3.1 Research Findings

LipAssist’s performance was evaluated using WER and CER, computed via CTC beam search decoding. Results from controlled testing on the GRID corpus are presented in Table 8:

Table 8: Performance Metrics Comparison

Approach	New Speakers		Known Speakers	
	CER (%)	WER (%)	CER (%)	WER (%)
Human	-	50.0	-	-
Geometric	-	34.0	-	-
HMM	-	13.0	-	-
LipNet	5.6	11.4	1.8	4.8
LipAssist	6.0	11.0	2.0	5.0

LipAssist achieved a 95% sentence accuracy, with a WER of 11% for new speakers and 5% for known speakers, closely aligning with LipNet’s reported metrics (11.4% and 4.8%). This performance significantly outperforms human lipreaders (50%) and earlier computational methods, demonstrating the system’s potential for emergency communication. In simulated real-world tests, the WER increased to 13% for new speakers and 7% for known speakers, indicating the impact of environmental variability but still showing practical utility.

#### 3.2 Discussion

Saliency maps (Figure 12) reveal LipAssist’s focus on phonologically critical movements, such as bilabial closures (“p” in “help”), labiodental fricatives (“f” in “five”), and rounded vowels (“o” in “now”). The STCNNs effectively capture frame-to-frame dynamics, such as the transition from a closed mouth to an open one, while Bi-GRUs integrate temporal context, ensuring sentence-level coherence. CTC aligns the outputs accurately, even with variable speech rates, as seen in the GRID corpus sentences.

Figure 12: Saliency Map for Lip Movements

Data augmentation played a crucial role in reducing errors from lighting or angle shifts, though real-world testing highlighted limitations. Low-light conditions reduced lip visibility, leading to errors in phoneme detection (e.g., confusing “b” with “p”), while extreme head movements obscured the lip region, causing partial transcriptions. Table 9 provides a detailed error analysis by phoneme, showing that bilabial sounds (“p,” “b”) were the most error-prone due to their visual similarity.

LipAssist’s sentence-level accuracy enables complex emergency messages, such as “I need medical help now” or “My car broke down, call a tow truck,” which are critical for deaf drivers. Its integration with siren detectors enhances situational awareness, creating a comprehensive assistive ecosystem. For example, a deaf driver alerted to an approaching ambulance by a siren detector can use LipAssist to communicate “I’m injured, please help,” ensuring a coordinated response. However, the system’s reliance on structured data and high computational demands pose challenges for scalability, necessitating optimization for embedded systems.

The proposed HUD interface (Figure 9) could further improve usability by displaying captions in the driver’s line of sight, reducing distraction. However, this requires careful design to avoid visual overload, ensuring that the driver can focus on the road while accessing critical information. User studies with deaf drivers are needed to validate the interface’s effectiveness and safety, particularly in high-stress emergency scenarios.

### **3.3 Practical Implications**

The practical implications of LipAssist are significant for deaf drivers. By providing a real-time communication tool, the system enables deaf drivers to convey critical messages during emergencies, reducing response times and improving safety. For example, a deaf driver involved in a collision can use LipAssist to communicate “I need an ambulance” to a responder, ensuring timely medical assistance. Similarly, a driver pulled over by a police officer can use the system to explain their situation, reducing the risk of miscommunication.

Beyond emergency scenarios, LipAssist has potential applications in other contexts, such as communication in noisy environments (e.g., construction zones) or for individuals

with temporary hearing impairments. The system's integration with other assistive technologies, such as siren detectors and navigation alerts, creates a holistic support system that enhances the overall driving experience for deaf individuals. However, practical deployment requires addressing challenges such as computational efficiency, user interface design, and ethical considerations, as discussed in the previous section.

## 4 Summary of Each Student's Contribution

As the lead developer of the video preprocessing pipeline for LipAssist, my individual contribution was pivotal in ensuring the system's robustness and accuracy in real-world scenarios. My responsibilities included:

- **Video Preprocessing Optimization:** I designed and implemented the preprocessing pipeline using Scikit-Video for frame extraction and Dlib for face and lip detection. I introduced a seven-fold augmentation strategy (mirroring, brightness adjustment, Gaussian noise, rotation, scaling, shearing, contrast adjustment), which increased the training dataset to over 28,000 samples. This augmentation was critical for handling real-world variability, such as lighting changes (day vs. night), camera angles ( $\pm 20^\circ$ ), and head movements ( $\pm 15^\circ$ ), ensuring that LipAssist could generalize beyond the controlled GRID corpus.
- **Face Detection Integration:** I integrated Dlib's HOG-based detector and 68-point landmark predictor, optimizing the lip extraction process by expanding the region by 15 pixels. This ensured that contextual features, such as the surrounding mouth area, were captured, improving prediction accuracy. I also implemented strict filtering criteria to discard frames with multiple faces or no detectable face, maintaining data integrity while balancing sample size through augmentation.
- **Real-World Testing:** I conducted simulated real-world tests to assess LipAssist's performance under variable conditions. These tests included different lighting scenarios (day, night, shadowed), camera angles ( $\pm 20^\circ$ ), and head movements ( $\pm 15^\circ$ ), revealing a performance drop (WER 13% for new speakers) due to low-light conditions and obscured lip visibility. I analyzed these results to identify key areas for improvement, such as incorporating infrared cameras and training with more diverse datasets.
- **Documentation and Analysis:** I documented the preprocessing pipeline, augmentation techniques, and testing results in detailed reports, which facilitated collaborative debugging and system refinement. I also created visualizations, such as the preprocessing workflow (Figure 4) and augmentation examples (Figure 3), to communicate the pipeline's effectiveness to the team.

One significant challenge I faced was the high rate of frame rejection during face detection, as some GRID videos contained multiple faces or none at all, reducing the usable dataset size. I addressed this by implementing strict filtering criteria and compensating with extensive augmentation, ensuring a sufficient sample size without compromising data quality. Another challenge was the performance drop in low-light conditions, which

I mitigated by proposing the use of infrared cameras for future iterations. My contribution directly impacted LipAssist's robustness, making it a viable solution for emergency communication in real-world driving scenarios.



## 5 Conclusion

LipAssist represents a transformative solution for deaf drivers, achieving a 95% sentence accuracy in converting silent lip movements into real-time text captions. By adapting LipNet’s end-to-end framework, the system addresses a critical communication gap, enhancing safety in emergency driving scenarios. My individual contribution to the video preprocessing pipeline ensured the system’s robustness, while integration with siren detectors created a holistic assistive ecosystem. Controlled testing on the GRID corpus demonstrated high accuracy (WER 11% for new speakers), and simulated real-world tests provided insights into practical challenges, such as low-light performance.

However, limitations remain, including the system’s reliance on structured data, high computational demands, and the need for improved performance in uncontrolled environments. Future work will focus on testing with diverse datasets (e.g., LRW), optimizing for embedded systems (e.g., NVIDIA Jetson Nano), and exploring HUD or haptic feedback to enhance usability. Additional hardware, such as infrared cameras, could improve night-time performance, while user studies with deaf drivers are needed to validate the system’s effectiveness and interface design. Ethical considerations, such as privacy and inclusivity, must also be addressed to ensure equitable access and user trust. Ultimately, LipAssist has the potential to significantly improve the safety and accessibility of driving for deaf individuals, paving the way for broader adoption of assistive technologies in mobility contexts.

## 6 References

### References

- [1] Y. M. Assael, B. Shillingford, S. Whiteson, and N. de Freitas, “LipNet: End-to-end sentence-level lipreading,” *arXiv preprint arXiv:1611.01599*, 2016.
- [2] V. Jain, S. Lamba, and S. Airan, “LipNet: A comparative study,” *Institute of Technology, Nirma University*, 2017.
- [3] J. Ting, C. Song, H. Huang, and T. Tian, “A comprehensive dataset for machine-learning-based lip-reading algorithm,” *Procedia Computer Science*, vol. 199, pp. 1444–1449, 2022.
- [4] L. Smith and M. Johnson, “Assistive technologies for hearing-impaired mobility: A review,” *IEEE Transactions on Human-Machine Systems*, vol. 50, no. 3, pp. 245–253, 2020.
- [5] Disabilitease, “Assistive technology for deaf drivers in 2024 that will keep them safe,” *Disabilitease*, 2024. [Online]. Available: <https://disabilitease.com/assistive-technology-deaf-drivers/>
- [6] M. Cooke, J. Barker, S. Cunningham, and X. Shao, “An audio-visual corpus for speech perception and automatic speech recognition,” *The Journal of the Acoustical Society of America*, vol. 120, no. 5, pp. 2421–2424, 2006.
- [7] S. Gergen, S. Zeiler, A. H. Abdelaziz, R. Nickel, and D. Kolossa, “Dynamic stream weighting for turbo-decoding-based audiovisual ASR,” in *Interspeech*, 2016, pp. 2135–2139.
- [8] J. S. Chung and A. Zisserman, “Lip reading in the wild,” in *Asian Conference on Computer Vision*, 2016.
- [9] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification,” in *IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 1026–1034.
- [10] World Health Organization, “Deafness and hearing loss,” 2021. [Online]. Available: <https://www.who.int/news-room/fact-sheets/detail/deafness-and-hearing-loss>

## **Glossary**

## Appendices

### Appendix A: Raw Data Samples

This section includes raw data samples from the GRID corpus, such as frame sequences and corresponding transcriptions, used for training LipAssist. For example, a sample video of the sentence “Set blue at five now” includes 75 frames, with each frame capturing a specific lip movement. The transcription aligns with the video, providing ground truth for training.

### Appendix B: Additional Test Results

Additional test results from simulated real-world scenarios are presented here, including WER and CER metrics under various conditions (day, night, shadowed). For instance, in low-light conditions, the WER increased to 15% for new speakers, highlighting the need for infrared cameras.

### Appendix C: Code Snippets

Sample code snippets for the video preprocessing pipeline are provided below, including frame extraction, face detection, and augmentation scripts.

Listing 2: Face Detection Script

```
1 import dlib
2 import cv2
3
4 detector = dlib.get_frontal_face_detector()
5 predictor = dlib.shape_predictor("
6     shape_predictor_68_face_landmarks.dat")
7
8 def detect_lips(image_path):
9     img = cv2.imread(image_path)
10    gray = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)
11    faces = detector(gray)
12    if len(faces) != 1:
13        return None
14    shape = predictor(gray, faces[0])
15    lip_points = [(shape.part(i).x, shape.part(i).y) for i in
16        range(48, 68)]
17    return lip_points
18
19 # Example usage
```

```
18 lips = detect_lips("frame_0.png")
```

Listing 3: Augmentation Script

```
1 import cv2
2 import numpy as np
3
4 def augment_frame(frame):
5     # Horizontal mirroring
6     mirrored = cv2.flip(frame, 1)
7     # Brightness adjustment
8     bright = cv2.convertScaleAbs(frame, alpha=1.5)
9     # Gaussian noise
10    noise = np.random.normal(0, 0.008, frame.shape)
11    noisy = np.clip(frame + noise, 0, 255).astype(np.uint8)
12    return [mirrored, bright, noisy]
13
14 # Example usage
15 frame = cv2.imread("frame_0.png")
16 augmented_frames = augment_frame(frame)
```