

Autonomous IoT-Enabled Hazard Detection and Communication System for Deaf Drivers.

24-25J-132

Our Team



**Dr. Kapila
Dissanayaka**
(Supervisor)



**Ms. Ishara
Weerathunga**
(Co-Supervisor)



**Rathnayake
R.M.S.N**
IT21388316



**Fernando
W.T.R.P**
IT21278280



**Thatsara
S.M.K**
IT21219566



**Iroshan
G.H.M**
IT21229084

Introduction

Deaf drivers face challenges due to the inability to perceive auditory cues like emergency sirens and car honking, leading to delayed responses and increased accident risk. Current vehicle safety systems often overlook the needs of deaf drivers, resulting in a lack of reliable alerts and communication methods during critical situations.

This research aims to develop an "Autonomous IoT Enabled Hazard Detection and Communication System for Deaf Drivers" which will:

- Provide real-time visual and haptic alerts.
- Use advanced sensor technology, machine learning algorithms, and a mobile application.
- Bridge the communication gap.
- Enhance driving safety for deaf drivers.



Research Question



- How can an autonomous IoT-enabled system effectively detect emergency vehicle sirens and other critical sounds in real-time for deaf drivers?
- What types of sensor technologies and machine learning algorithms are most effective for identifying and processing auditory cues on the road?
- How can visual and haptic feedback mechanisms be optimized to ensure deaf drivers receive timely and clear notifications?
- What are the challenges in integrating such a system with existing vehicle infrastructure and how can they be overcome?
- How can a mobile application be designed to support real-time hazard detection and provide a seamless communication interface for deaf drivers?
- How can the system ensure reliable performance in diverse driving environments and conditions?

OBJECTIVES

MAIN OBJECTIVE

The main objective of this research is to develop a robust and accessible hazard detection and communication system specifically designed for deaf drivers.

SUB OBJECTIVE



Emergency Vehicle Detection and Notification System



Continued Vehicle Horn Detection and Alert System



Comprehensive Driver Behavior Monitoring and Reporting System



Integrated Emergency Support System for Deaf Drivers

Commercialization

Market Potential:

Partnerships with automotive manufacturers.

Targeted marketing towards the deaf community.

Integration with existing vehicle safety systems.

Potential for government or organizational funding support.





Rathnayake R.M.S.N

IT21388316

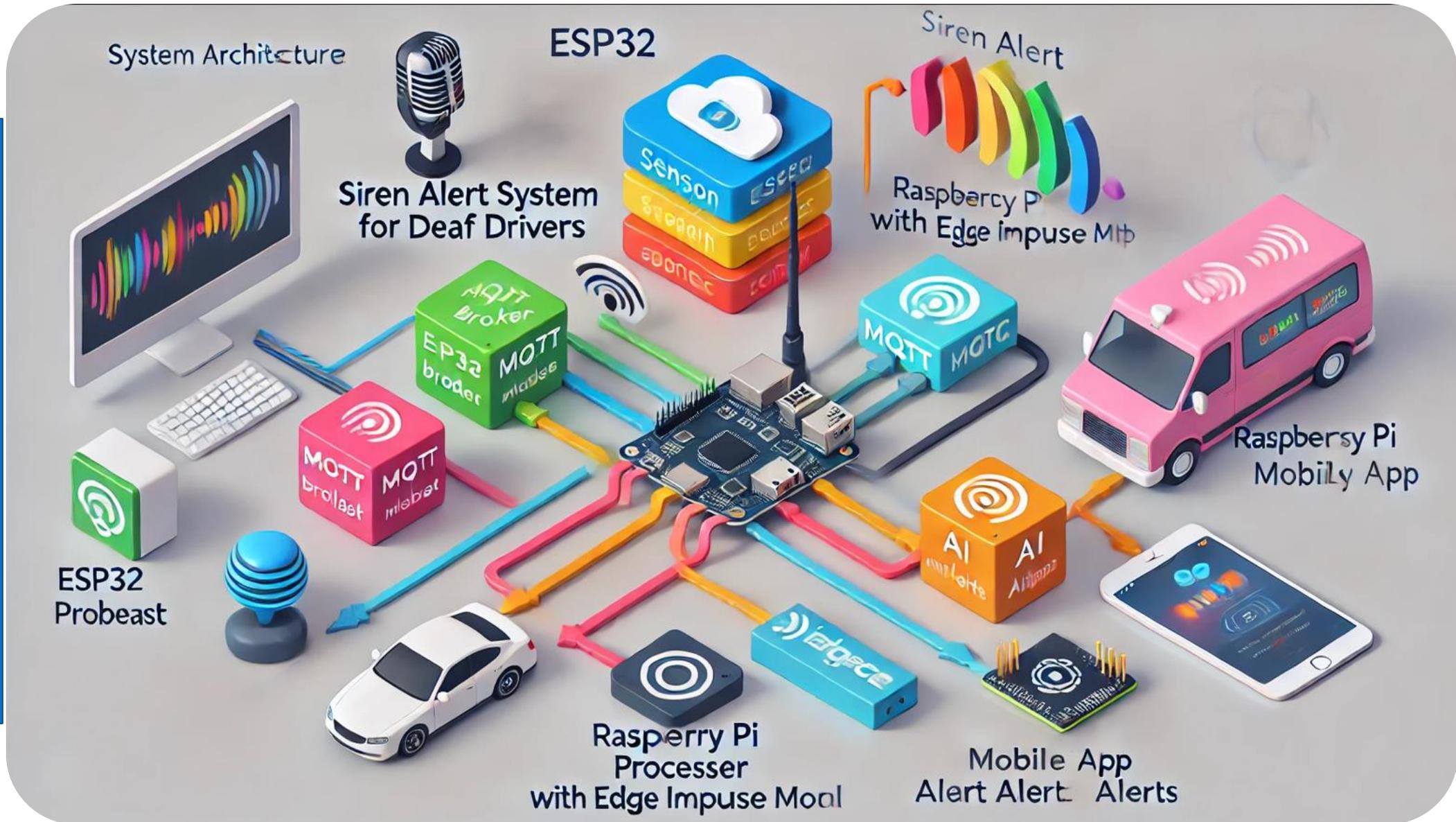
SOFTWARE ENGINEERING



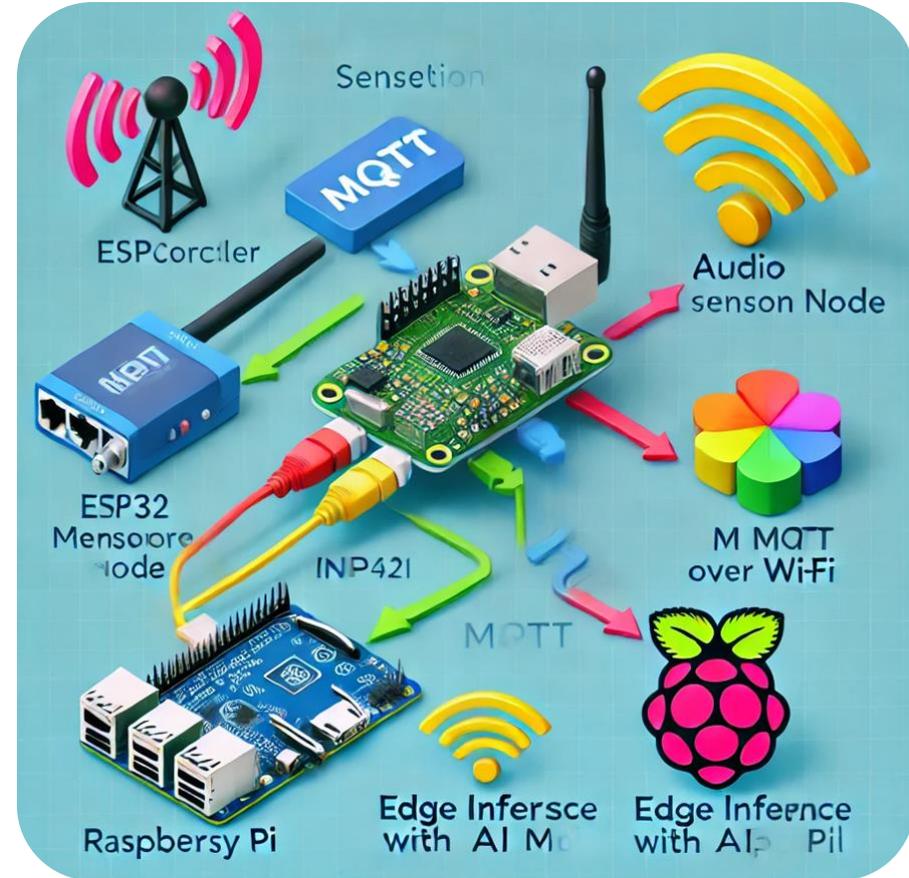
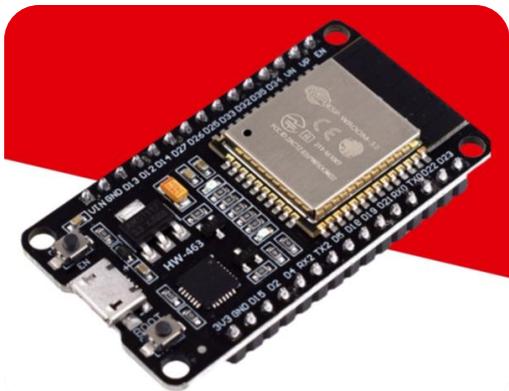


Emergency Vehicle Detection and Notification System

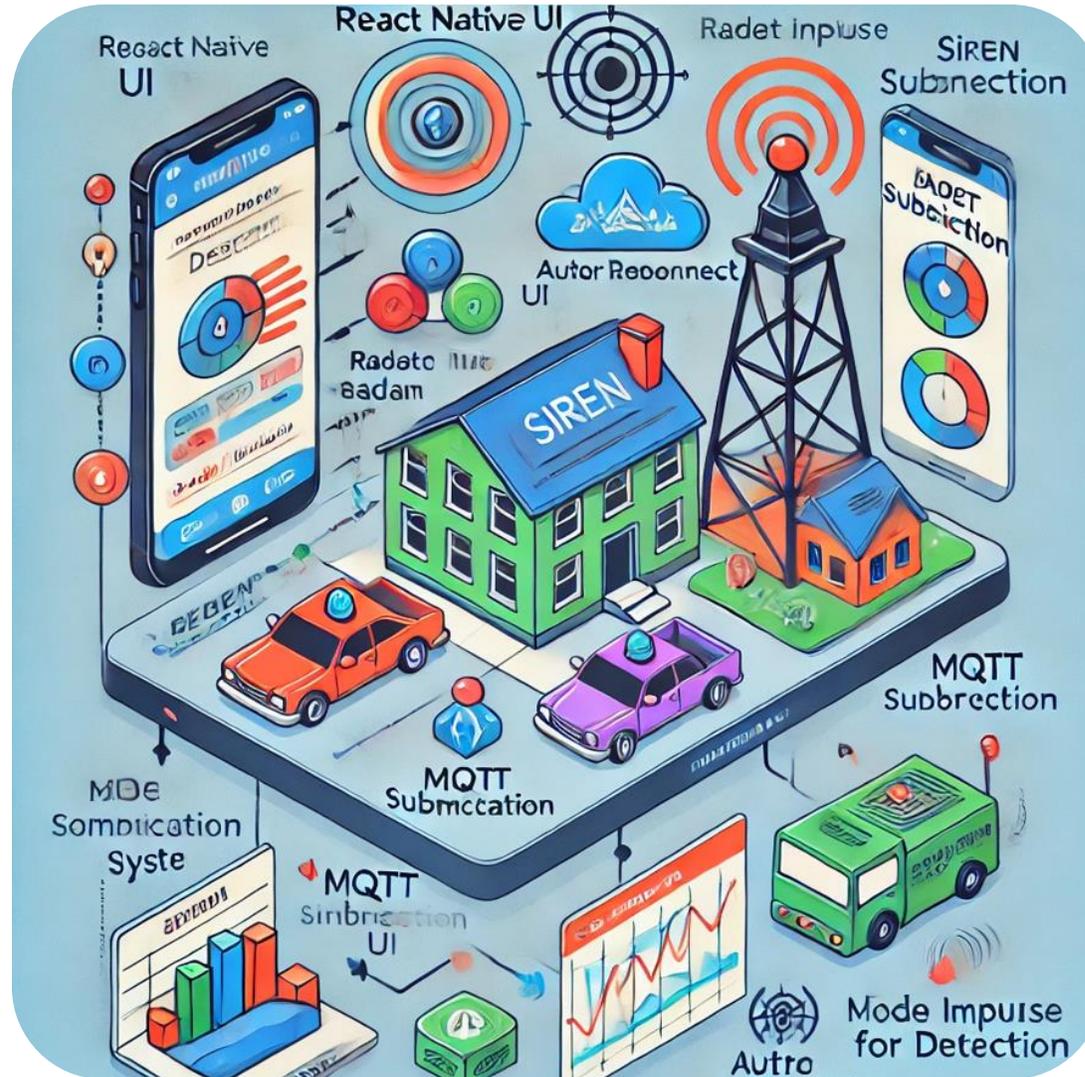
System Overview Diagram



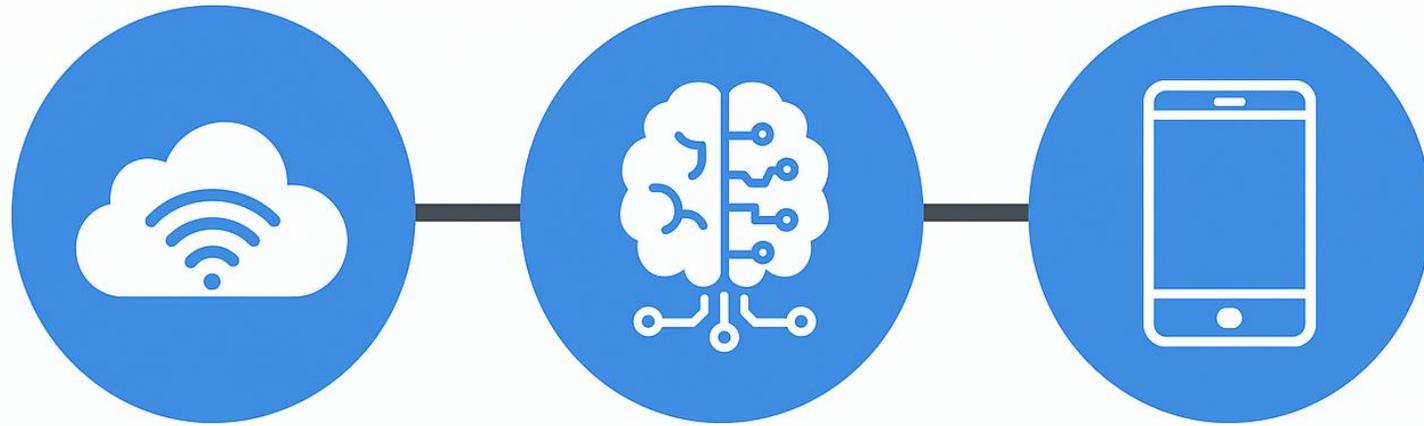
Hardware Architecture



Software Architecture



Covered Areas



IoT

AI

**Mobile App
Development**

Sources of Audio Data

- **Pre-existing Datasets**

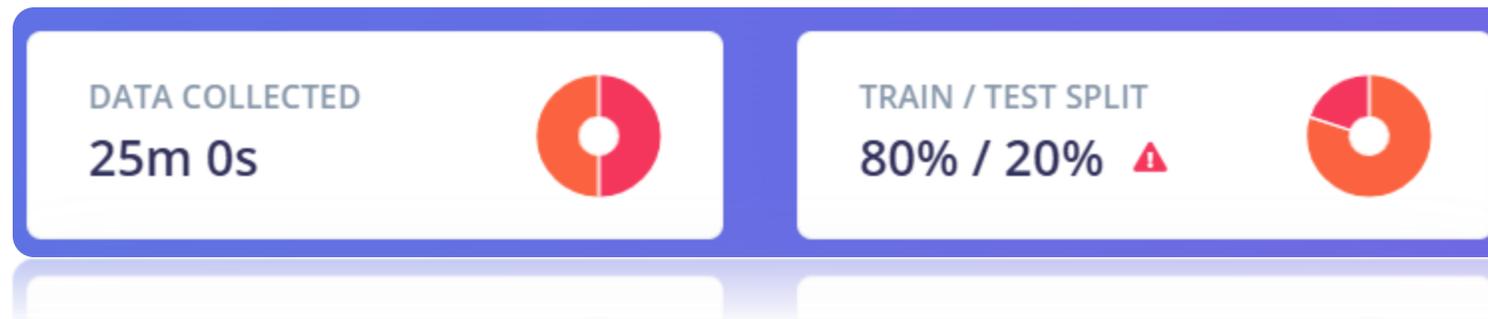
UrbanSound8K: Contains environmental sound clips, including sirens.

- **Manual Collection**

Emergency vehicles like ambulances, police cars, and fire trucks.
Different scenarios (urban, rural, heavy traffic).

- **Publicly Available Recordings**

Extract sounds from YouTube videos.

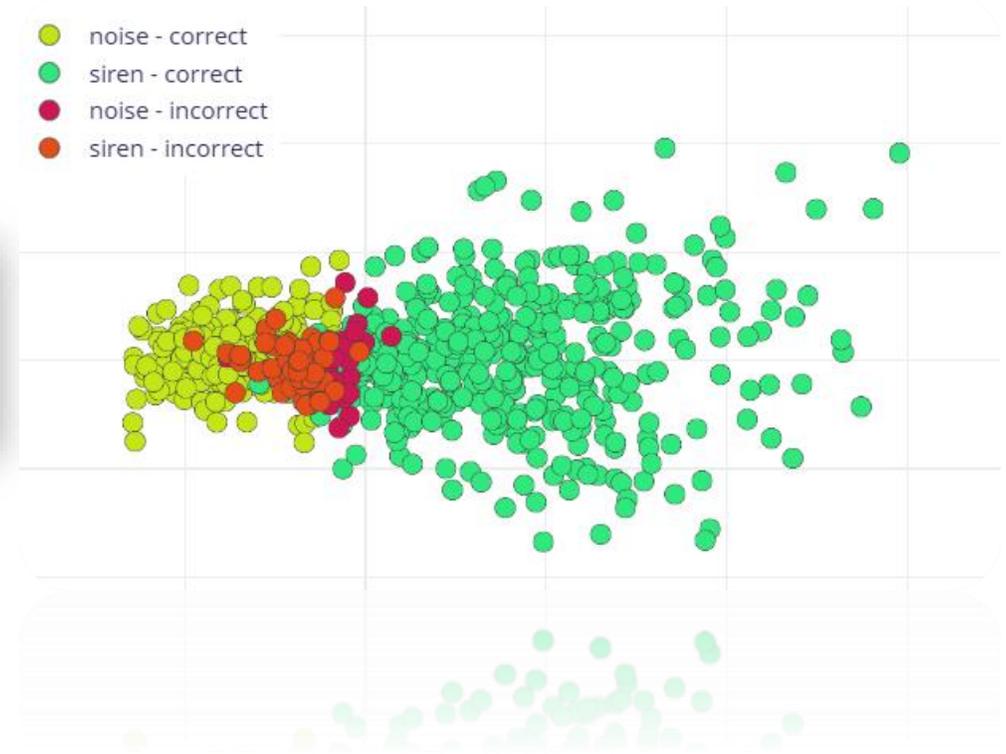
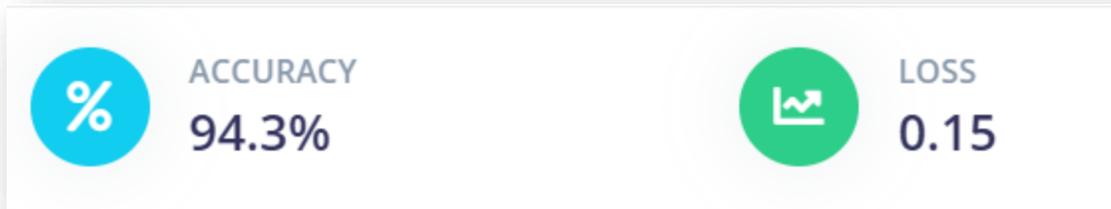


Data Acquisition



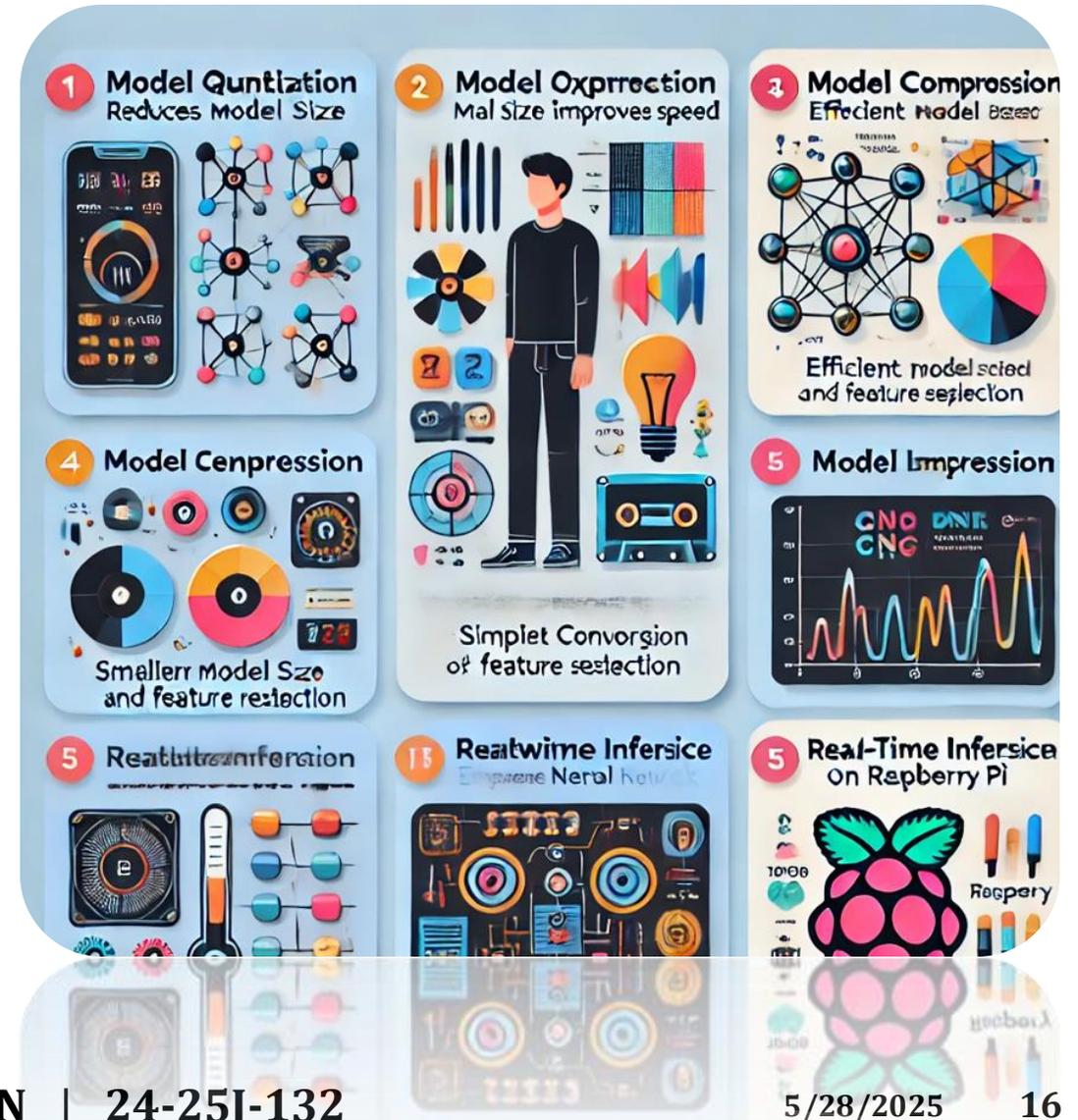
Deep Learning Model

- **Architecture** : Convolutional Neural Network (CNN) architecture
- **Number of training cycles(epochs)** : 300
- **Training processor** : CPU
- **Add Data augmentation**



Model Optimization Techniques

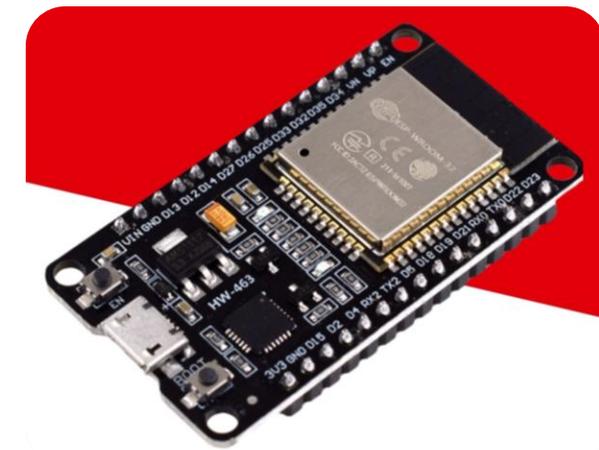
- **Quantization** — Reduces size, improves speed
- **used dropout** — to prevent overfitting
- **MFCC Features** — Efficient audio representation



ESP32 Nodemcu Board Processing

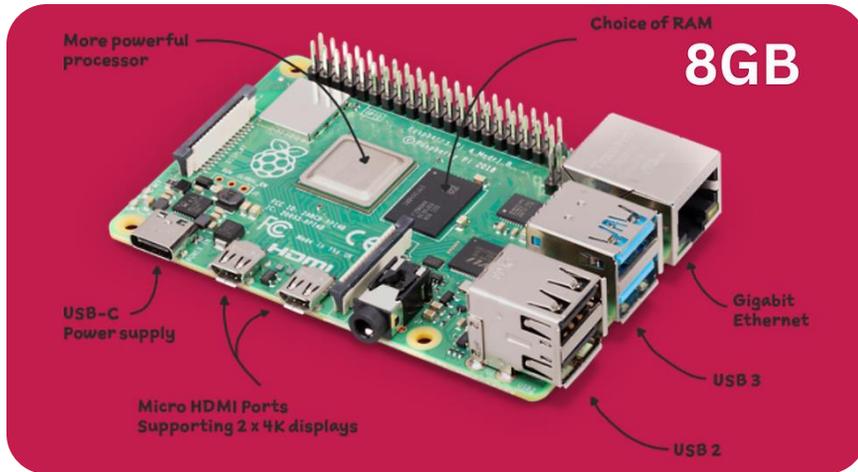
```
sketch_mar4a | Arduino IDE 2.3.4
File Edit Sketch Tools Help
ESP32 Dev Module
sketch_mar4a.ino
1 #include <WiFi.h>
2 #include <PubSubClient.h>
3 #include <driver/i2s.h>
4
5 #define I2S_WS 15
6 #define I2S_SCK 14
7 #define I2S_SD 32
8
9 const char* ssid = "Dialog 4G 815";
10 const char* password = "7e6B37f3";
11 const char* mqtt_server = "192.168.8.161";
12 const int mqtt_port = 1883;
13
14 // Optional: MQTT Authentication (if required)
15 const char* mqtt_username = ""; // Leave empty if no username
16 const char* mqtt_password = ""; // Leave empty if no password
17
18 WiFiClient espClient;
19 PubSubClient client(espClient);
20
21 // Generate a unique Client ID for each connection
22 String clientId = "ESP32_Client_" + String(random(1000, 9999));
23
24 // I2S Audio Configuration
25 void setupI2S() {
26     i2s_config_t i2s_config = {
27         .mode = (i2s_mode_t)(I2S_MODE_MASTER | I2S_MODE_RX),
28         .sample_rate = 16000,
29         .bits_per_sample = I2S_BITS_PER_SAMPLE_16BIT,
```

```
✓ WiFi Connected!
📶 Connecting to MQTT... ✓ Connected to MQTT Broker (Anonymous)
✓ Audio Chunk Sent Successfully
```



Raspberry pi Board Processing

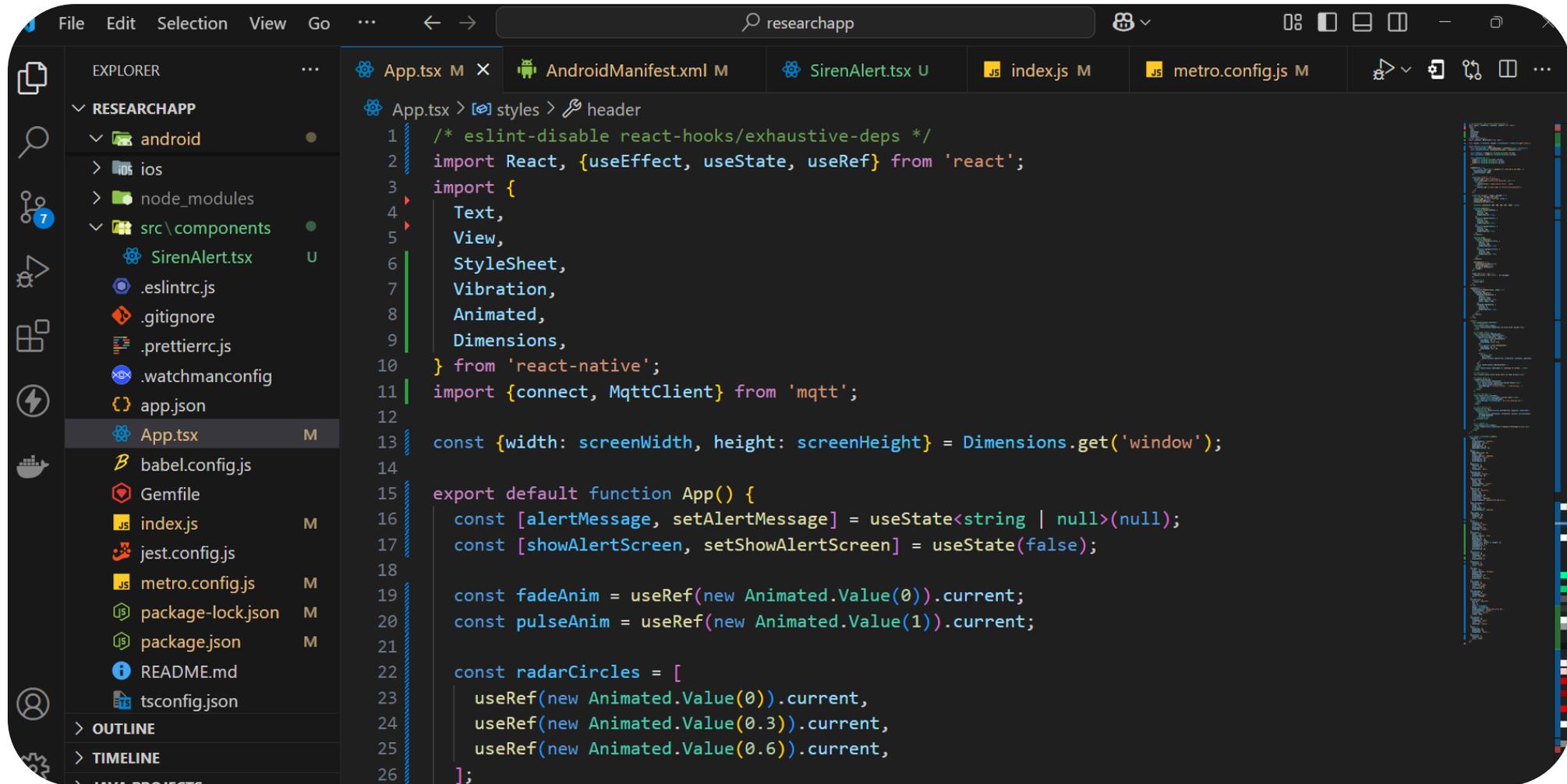
Siren_detection.py



```
File Edit Tabs Help
pi@raspberrypi ~
?? Received 64 samples... Total buffer size: 11136
?? Received 64 samples... Total buffer size: 11200
?? Received 64 samples... Total buffer size: 11264
?? Received 64 samples... Total buffer size: 11328
?? Received 64 samples... Total buffer size: 11392
?? Received 64 samples... Total buffer size: 11456
?? Received 64 samples... Total buffer size: 11520
?? Received 64 samples... Total buffer size: 11584
?? Received 64 samples... Total buffer size: 11648
?? Received 64 samples... Total buffer size: 11712
?? Received 64 samples... Total buffer size: 11776
?? Received 64 samples... Total buffer size: 11840
?? Received 64 samples... Total buffer size: 11904
?? Received 64 samples... Total buffer size: 11968
?? Received 64 samples... Total buffer size: 12032
?? Received 64 samples... Total buffer size: 12096
?? Received 64 samples... Total buffer size: 12160
?? Received 64 samples... Total buffer size: 12224
?? Received 64 samples... Total buffer size: 12288
?? Received 64 samples... Total buffer size: 12352
?? Received 64 samples... Total buffer size: 12416
?? Received 64 samples... Total buffer size: 12480
?? Received 64 samples... Total buffer size: 12544
```

```
siren_detection.py x
1 import paho.mqtt.client as mqtt
2 import numpy as np
3 import edge_impulse_linux.runner as ei
4 import time
5
6 # MQTT Configuration
7 MQTT_BROKER = "192.168.8.161"
8 MQTT_PORT = 1883
9 AUDIO_TOPIC = "audio/data"
10 ALERT_TOPIC = "alerts/siren_detected"
11
12 # Siren detection threshold
13 SIREN_THRESHOLD = 0.7 # Adjust if needed
14 COOLDOWN_PERIOD = 5 # 5 seconds cooldown between detections
15
16 # Initialize cooldown timer
17 last_detection_time = 0
18
19 # Initialize Edge Impulse model
20 try:
21     model = ei.ImpulseRunner("/home/pi/siren_model.eim")
22     model.init()
23     print("Edge Impulse Model Initialized!")
24 except Exception as e:
25     print(f"Error initializing model: {e}")
26     exit(1)
27
28 # MQTT Callbacks
29 def on_connect(client, userdata, flags, rc):
30     if rc == 0:
31         print("Connected to MQTT Broker")
32         client.subscribe(AUDIO_TOPIC)
33     else:
34         print(f"Failed to connect, return code {rc}")
35
36 def on_message(client, userdata, message):
37     global last_detection_time
38
39     audio_data = np.frombuffer(message.payload, dtype=np.int16)
40     print(f"Received {len(audio_data)} samples... Total buffer size: {userdata['buffer_size']}")
41
42     # Accumulate samples until we have 16,000 for inference
43     userdata['buffer'] = np.append(userdata['buffer'], audio_data)
44
45 This is Geany 1.38.
46 Setting Spaces indentation mode for /home/pi/siren_detection.py.
47 Setting Spaces indentation mode for /home/pi/siren_detection.py.
48 File /home/pi/siren_detection.py opened (1).
```

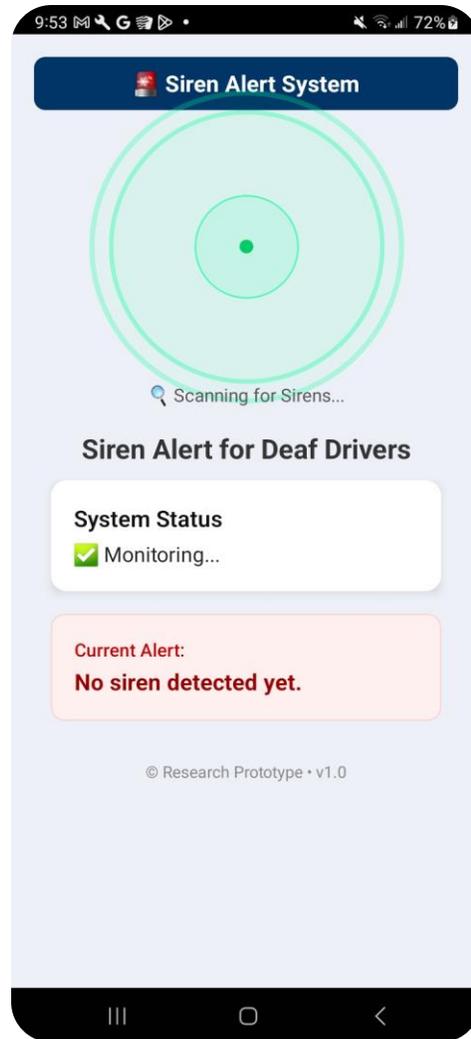
React Native App



The image shows a screenshot of the Visual Studio Code editor interface. The Explorer view on the left displays the project structure for 'RESEARCHAPP', including folders for 'android', 'ios', and 'node_modules', and a 'src/components' folder containing 'SirenAlert.tsx'. The main editor area shows the 'App.tsx' file with the following code:

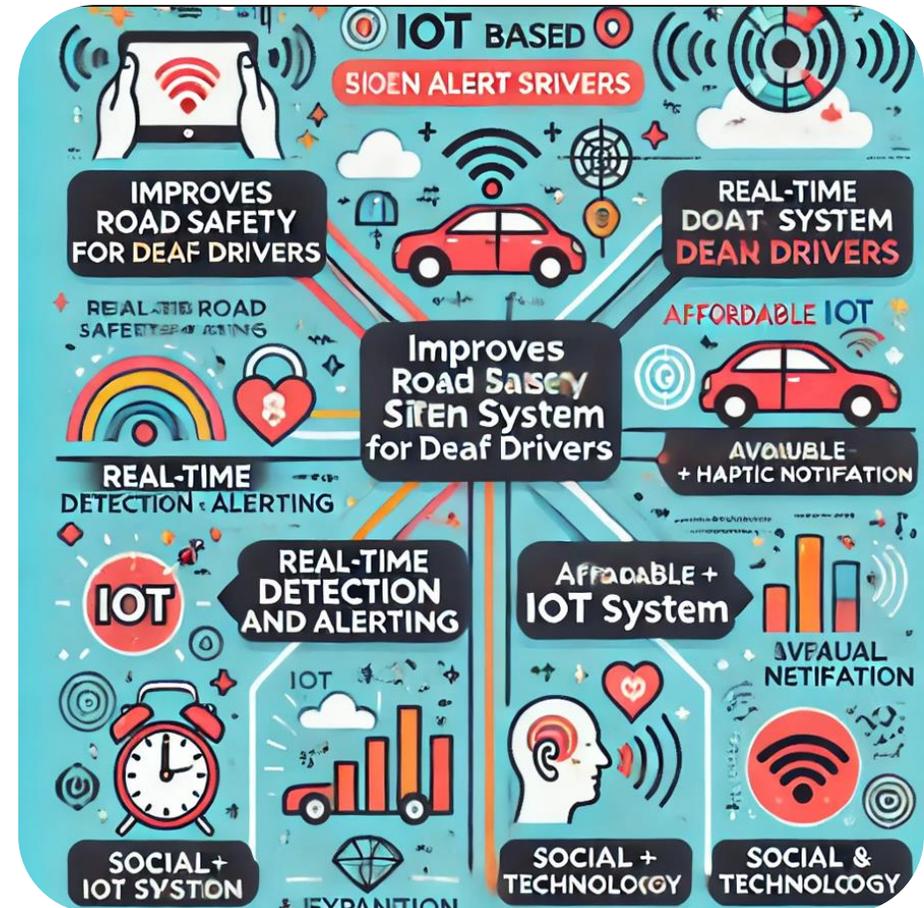
```
1  /* eslint-disable react-hooks/exhaustive-deps */
2  import React, {useEffect, useState, useRef} from 'react';
3  import {
4    Text,
5    View,
6    StyleSheet,
7    Vibration,
8    Animated,
9    Dimensions,
10 } from 'react-native';
11 import {connect, MqttClient} from 'mqtt';
12
13 const {width: screenWidth, height: screenHeight} = Dimensions.get('window');
14
15 export default function App() {
16   const [alertMessage, setAlertMessage] = useState<string | null>(null);
17   const [showAlertScreen, setShowAlertScreen] = useState(false);
18
19   const fadeAnim = useRef(new Animated.Value(0)).current;
20   const pulseAnim = useRef(new Animated.Value(1)).current;
21
22   const radarCircles = [
23     useRef(new Animated.Value(0)).current,
24     useRef(new Animated.Value(0.3)).current,
25     useRef(new Animated.Value(0.6)).current,
26   ];
```

Notification System



Impact & Benefits

- ❑ Improves Road Safety for Deaf Drivers
- ❑ Real-Time Detection and Alerting
- ❑ Affordable IoT-based System
- ❑ Scalable and Expandable
- ❑ Visual + Haptic notification
- ❑ Social and Technological Impact



References

1. T. Harshitha and T. Abhigna, "Emergency Vehicle Sound Detection Systems in Traffic Congestion," Ramachandrapuram, Telangana, 502032, and Jeedimetla, Quthbullapur, Telangana, 500055.
2. J. Smith, R. Brown, and C. White, "Machine Learning for Sound Pattern Recognition," International Journal of Artificial Intelligence, vol. 36, no. 2, pp. 89-101, Feb. 2021.
3. National Safety Council, "Road Safety and Deaf Drivers," National Safety Council Report, 2020. [Online]. Available: <https://www.nsc.org>
4. Tech Innovators, "Advancements in MEMS Microphone Technology," Tech Innovators Report, 2019. [Online]. Available: <https://www.techinnovators.com>
5. National Association of the Deaf, "Deaf and Hard of Hearing Drivers," National Association of the Deaf, 2020. [Online]. Available: <https://www.nad.org>
6. TensorFlow, "TensorFlow: An end-to-end open source machine learning platform," 2023. [Online]. Available: <https://www.tensorflow.org>



Fernando W.T.R.P

IT21278280

Software Engineering



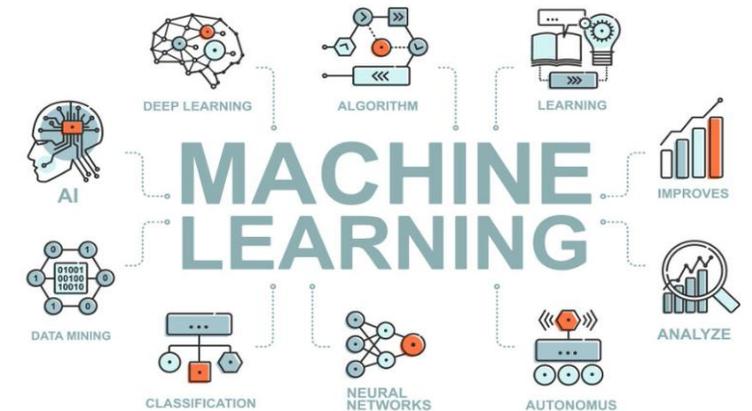
Continued Vehicle Horn Detection and Alert System



Introduction

- Deaf drivers cannot hear vehicle horns, risking their safety.

- Existing systems don't help deaf drivers effectively.



Research Question



How can we detect and classify horn sounds accurately?



What machine learning methods can separate horn sounds from background noise?



How can we determine the direction of the horn sound using microphones?



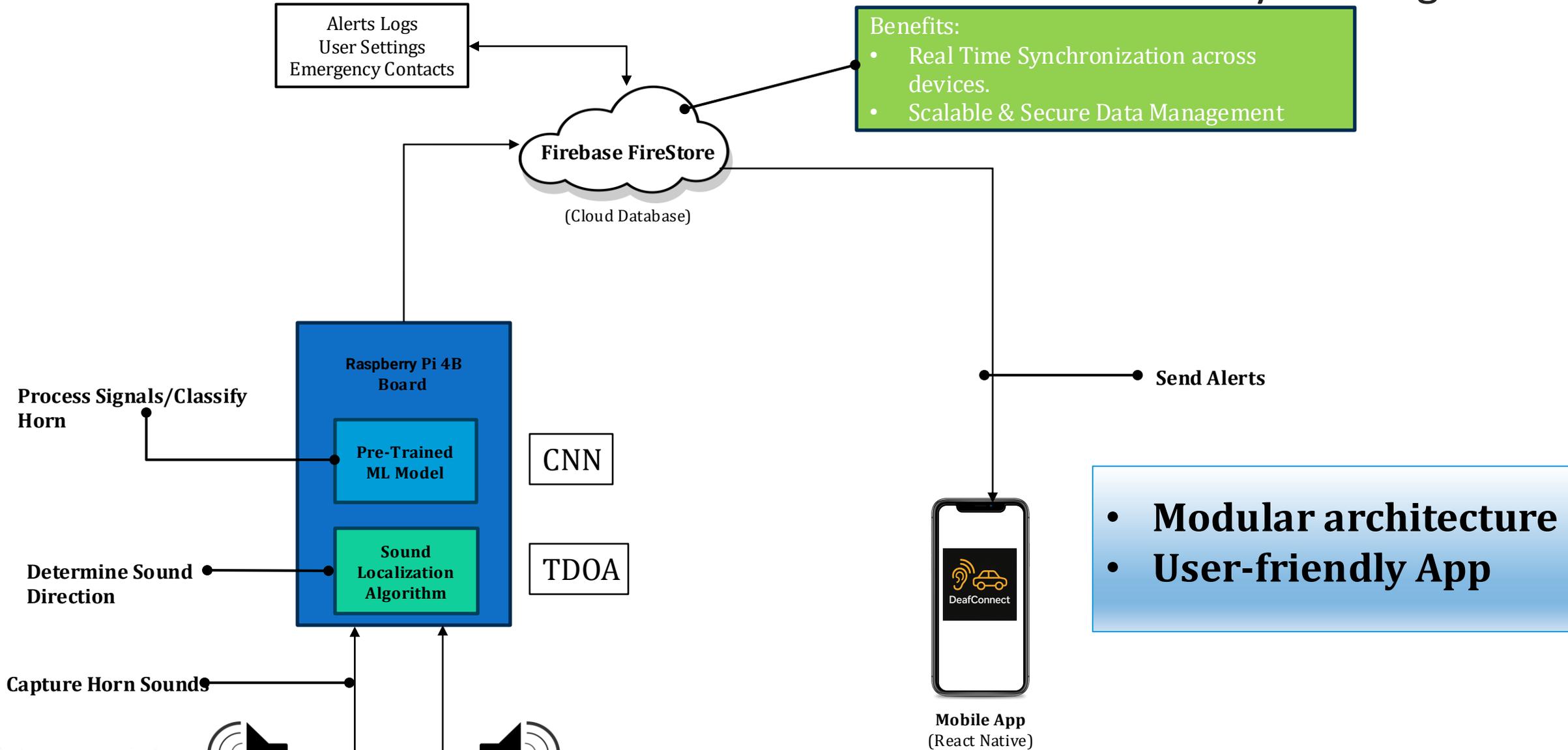
What's the best way to alert deaf drivers using visual & haptic feedback?



Research Gap and Novelty

Research	Real-Time Detection	Localization	Noise Robustness	Deaf Driver Focus	Year	Key Limitation
Beritelli					2021	No directional awareness
Sharma					2024	Predictive focus only
Zhao					2023	Lab-Only testing
My System					2025	Comprehensive Solution

Methodology System Diagram



Specific and Sub Objective

Specific Objective



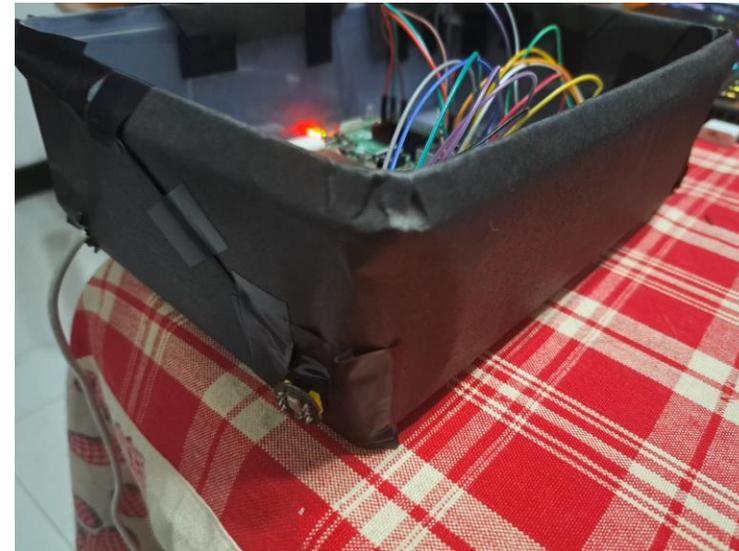
Develop a real-time horn detection system for deaf drivers



Sub Objective

1

Set up Microphones to Capture Horn Sounds



100%

2

Use Machine Learning to Classify Horn Sounds

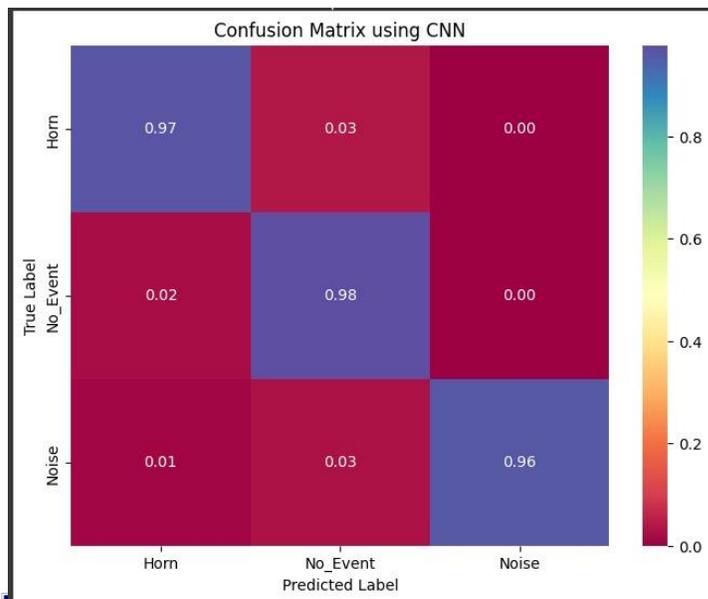


Horn Classification

- A pre-trained TensorFlow Lite model checks if the sound is a horn.



Mel Spectrogram: Convert audio to a frequency-time map (Mel spectrogram) with 2048-point FFT, 512-sample hop, 512 Mel bands, 86 time steps.



Normalizes the spectrogram :

$$S_{\text{norm}} = \frac{S_{\text{dB}} - 1.39 \times 10^{-14}}{0.99999999999953707}$$



Model Prediction:

- Predicted class (e.g., 'horn').
- Confidence (as %).
- Direction angle (θ_{deg}).
- Is it a horn? (Yes if class = 'horn' and confidence > 0.6).
- Left or right? (Left if $\theta_{\text{deg}} < 0$).

Test Accuracy - 97%

Sub Objective

3

Apply TDOA to find the horn's direction

The system figures out if the horn comes from the left or right using the time difference between the microphones.

Cross-Correlation: Compares left and right channel signals to find the time delay

$$\text{time_difference} = \frac{\text{max_lag}}{44100}$$

Angle Calculation:

$$\theta = \arcsin\left(\frac{\text{time_difference} \cdot 343}{0.1}\right)$$

$$\theta_{\text{deg}} = \frac{180}{\pi} \cdot \theta$$

Direction:

- If $\theta_{\text{deg}} < 0$: Sound comes from the left.
- If $\theta_{\text{deg}} > 0$: Sound comes from the right.

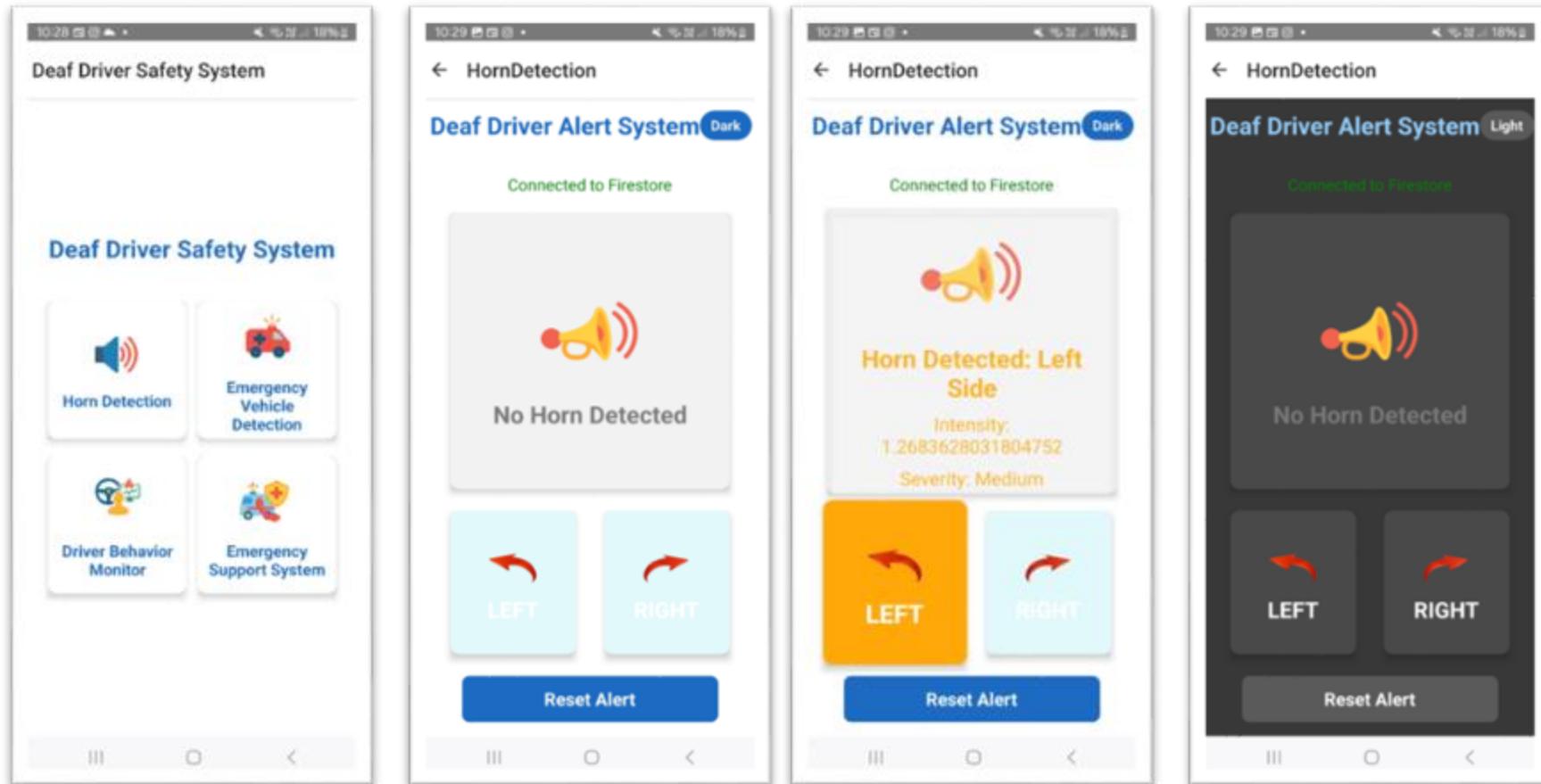
100%



Sub Objective

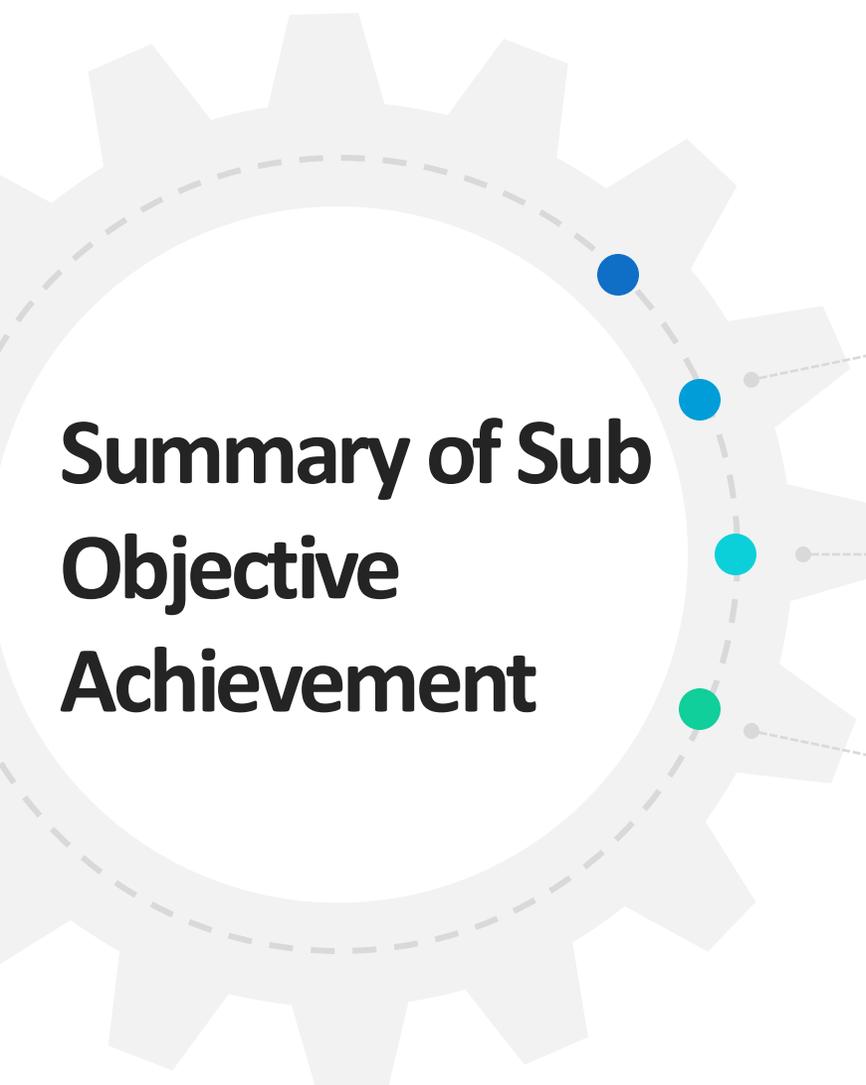
4

Build a Mobile App for Visual & Haptic Alerts



100%

Evidences for Completion



Summary of Sub Objective Achievement



Data Collection and Pre-Processing

```
# Step 1: Label Encoding
label_encoder = LabelEncoder()
encoded_labels = label_encoder.fit_transform(labels)
# Convert to one-hot encoded format for multi-class classification
encoded_labels = to_categorical(encoded_labels, num_classes=len(label_encoder.classes_))

# Step 2: Data Preprocessing

# Normalize the mel-spectrograms (scale to [0, 1])
Mel_features = (Mel_features - np.min(Mel_features)) / (np.max(Mel_features) - np.min(Mel_features))
print("Min = ", np.min(Mel_features))
print("Max = ", np.max(Mel_features))

# Reshape data to include channel dimension for CNN (n_samples, height, width, channels)
Mel_features = Mel_features[... , np.newaxis] # Shape: (n_samples, 512, max_length_frames, 1)

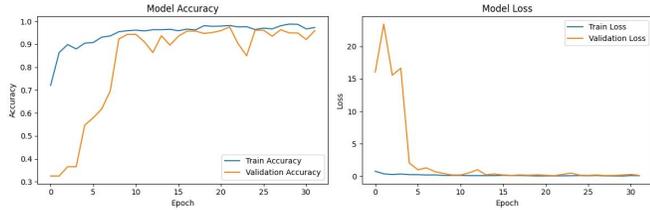
# Split into training and testing sets (80% train, 20% test)
X_train, X_test, y_train, y_test = train_test_split(
    Mel_features, encoded_labels, test_size=0.2, stratify=encoded_labels
)

# Further split training+validation into training and validation (75% train, 25% val of the 80%)
X_train, X_val, y_train, y_val = train_test_split(
    X_train, y_train, test_size=0.3, random_state=52, stratify=y_train
)

# Print shapes to verify
print(f"Training set shape: {X_train.shape}, {y_train.shape}")
print(f"Validation set shape: {X_val.shape}, {y_val.shape}")
print(f"Test set shape: {X_test.shape}, {y_test.shape}")
```



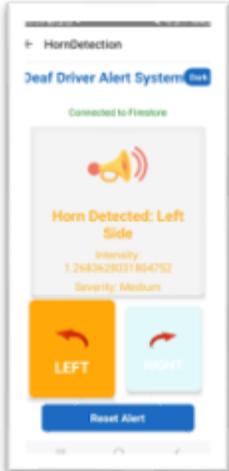
Training the Model



Accuracy Measures

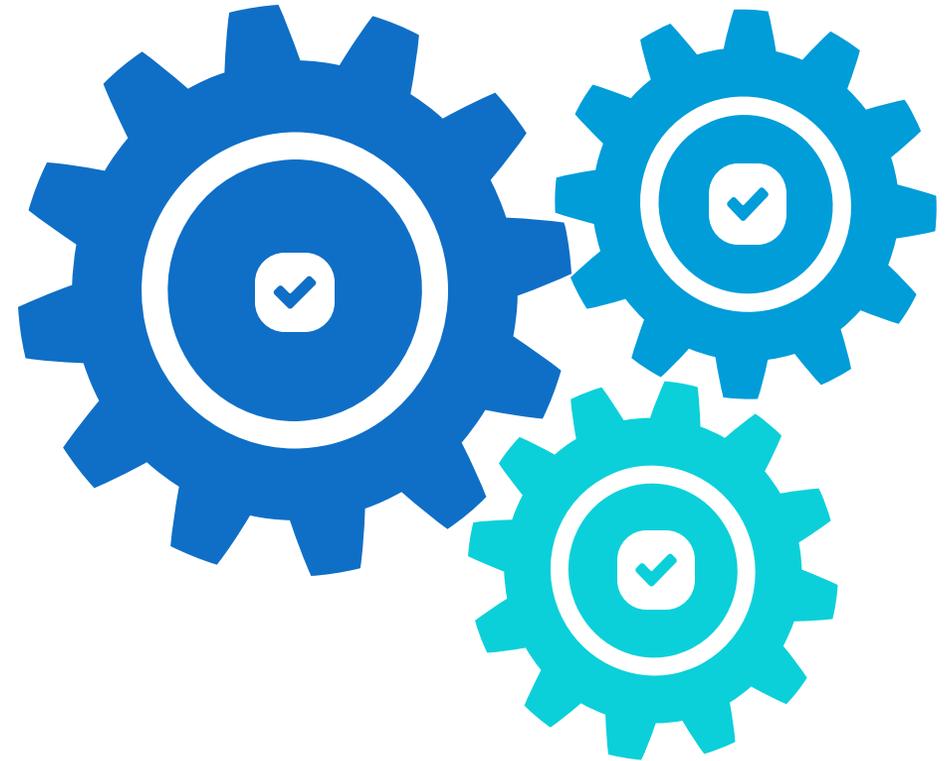
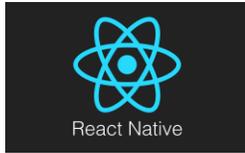


Testing using an App



100% Completion

Technologies, Techniques & Algorithms



Non-Functional Requirements

Standards/
Best Practices

Real-Time Response

Usability for Deaf drivers

Reliability and accuracy

Modular Design

Simple UI for Deaf Users

Error Handling

Open-source Frameworks

Ethical, Legal, and Social Issues



- **Respects user privacy.**
- **Does not store personal data.**
- **Designed for safety and accessibility.**
- **Supports equal access for deaf drivers.**

Data Preprocessing

```
# Step 1: Label Encoding
label_encoder = LabelEncoder()
encoded_labels = label_encoder.fit_transform(labels)
# Convert to one-hot encoded format for multi-class classification
encoded_labels = to_categorical(encoded_labels, num_classes=len(label_encoder.classes_))

# Step 2: Data Preprocessing

# Normalize the mel-spectrograms (scale to [0, 1])
Mel_features = (Mel_features - np.min(Mel_features)) / (np.max(Mel_features) - np.min(Mel_features))
print("Min = ", np.min(Mel_features))
print("Max = ", np.max(Mel_features))

# Reshape data to include channel dimension for CNN (n_samples, height, width, channels)
Mel_features = Mel_features[..., np.newaxis] # Shape: (n_samples, 512, max_length_frames, 1)

# Split into training and testing sets (80% train, 20% test)
X_temp, X_test, y_temp, y_test = train_test_split(
    Mel_features, encoded_labels, test_size=0.2, random_state=40, stratify=encoded_labels
)

# Further split training+validation into training and validation (75% train, 25% val of the 80%)
X_train, X_val, y_train, y_val = train_test_split(
    X_temp, y_temp, test_size=0.3, random_state=62, stratify=y_temp
)

# Print shapes to verify
print(f"Training set shape: {X_train.shape}, {y_train.shape}")
print(f"Validation set shape: {X_val.shape}, {y_val.shape}")
print(f"Test set shape: {X_test.shape}, {y_test.shape}")
```

Proof of Completion

Build the CNN Model

```
# Step 3: Build the CNN Model
def create_cnn_model(input_shape, num_classes):
    model = models.Sequential([
        # First Convolutional Block
        layers.Conv2D(32, (3, 3), activation='relu', padding='same', input_shape=input_shape),
        layers.BatchNormalization(),
        layers.MaxPooling2D((2, 2)),
        layers.Dropout(0.25),

        # Second Convolutional Block
        layers.Conv2D(64, (3, 3), activation='relu', padding='same'),
        layers.BatchNormalization(),
        layers.MaxPooling2D((2, 2)),
        layers.Dropout(0.25),

        # Third Convolutional Block
        layers.Conv2D(128, (3, 3), activation='relu', padding='same'),
        layers.BatchNormalization(),
        layers.MaxPooling2D((2, 2)),
        layers.Dropout(0.25),

        # Flatten and Dense Layers
        layers.Flatten(),
        layers.Dense(256, activation='relu'),
        layers.BatchNormalization(),
        layers.Dropout(0.5),
        layers.Dense(num_classes, activation='softmax')
    ])
    return model

# Define input shape and number of classes
input_shape = (Mel_features.shape[1], Mel_features.shape[2], 1) # (512, max_length_frames, 1)
num_classes = len(label_encoder.classes_) # 3 classes

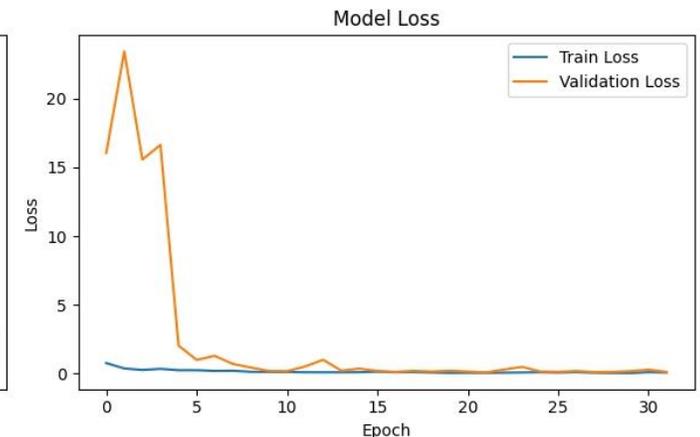
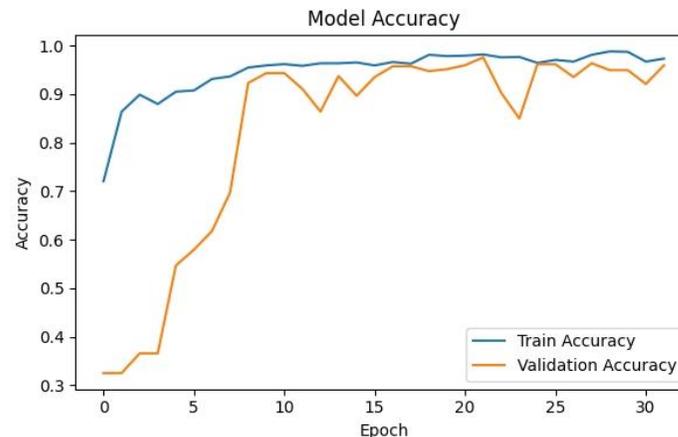
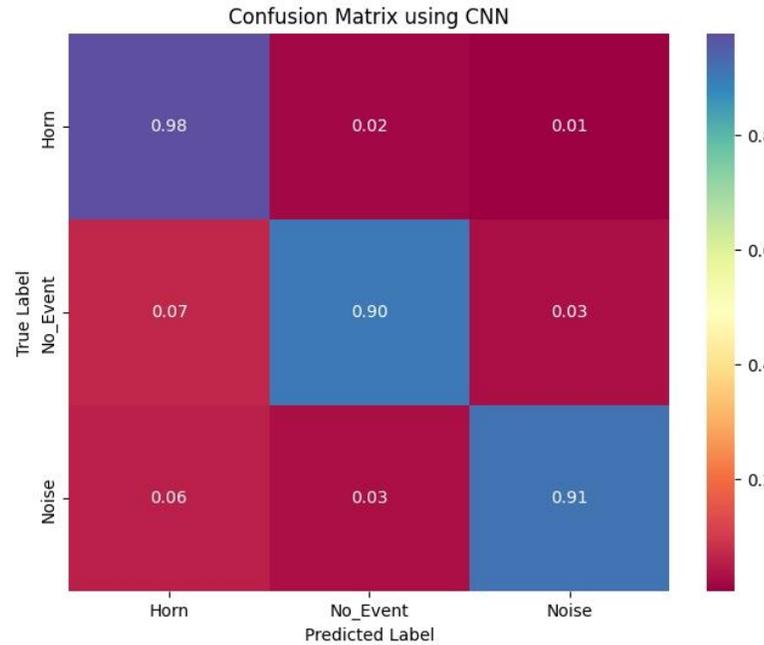
# Create and compile the model
model = create_cnn_model(input_shape, num_classes)
model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])

# Print model summary
model.summary()
```

Proof of Completion

Layer (type)	Output Shape	Param #
conv2d (Conv2D)	(None, 512, 86, 32)	320
batch_normalization (BatchNormalization)	(None, 512, 86, 32)	128
max_pooling2d (MaxPooling2D)	(None, 256, 43, 32)	0
dropout (Dropout)	(None, 256, 43, 32)	0
conv2d_1 (Conv2D)	(None, 256, 43, 64)	18,496
batch_normalization_1 (BatchNormalization)	(None, 256, 43, 64)	256
max_pooling2d_1 (MaxPooling2D)	(None, 128, 21, 64)	0
dropout_1 (Dropout)	(None, 128, 21, 64)	0
conv2d_2 (Conv2D)	(None, 128, 21, 128)	73,856
batch_normalization_2 (BatchNormalization)	(None, 128, 21, 128)	512
max_pooling2d_2 (MaxPooling2D)	(None, 64, 10, 128)	0
dropout_2 (Dropout)	(None, 64, 10, 128)	0
flatten (Flatten)	(None, 81920)	0
dense (Dense)	(None, 256)	20,971,776
batch_normalization_3 (BatchNormalization)	(None, 256)	1,024
dropout_3 (Dropout)	(None, 256)	0
dense_1 (Dense)	(None, 3)	771

Total params: 21,067,139 (80.36 MB)
 Trainable params: 21,066,179 (80.36 MB)
 Non-trainable params: 960 (3.75 KB)

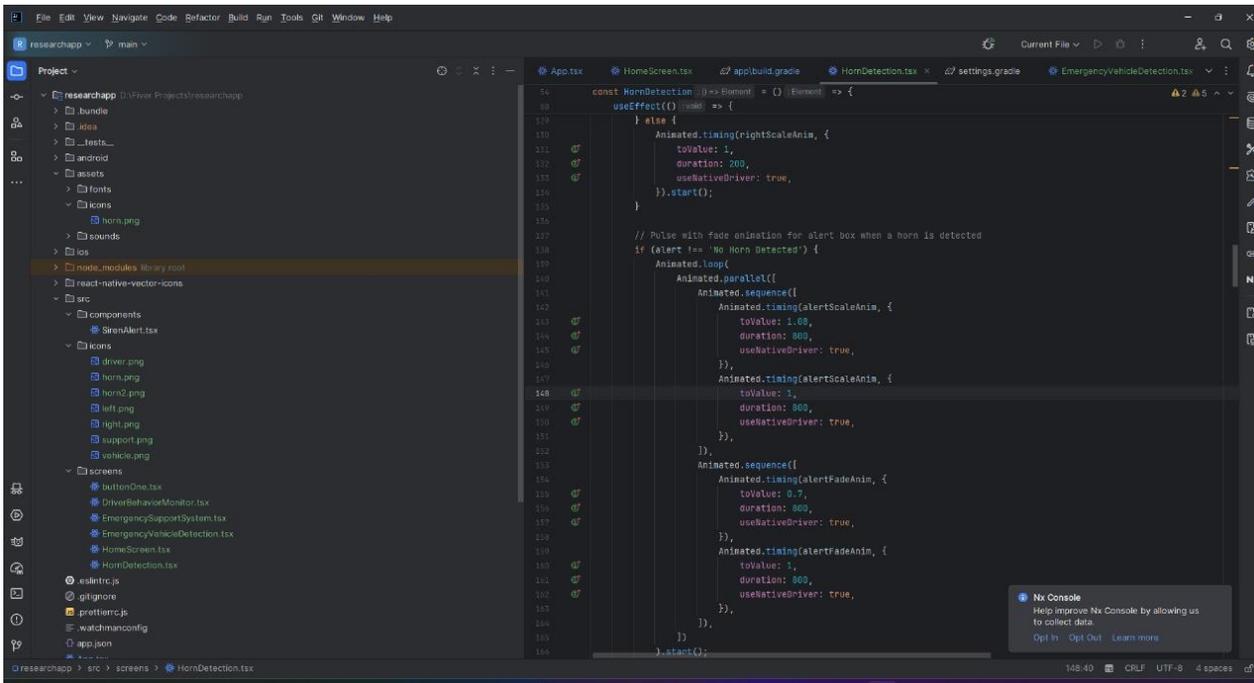


Test Accuracy - 97%

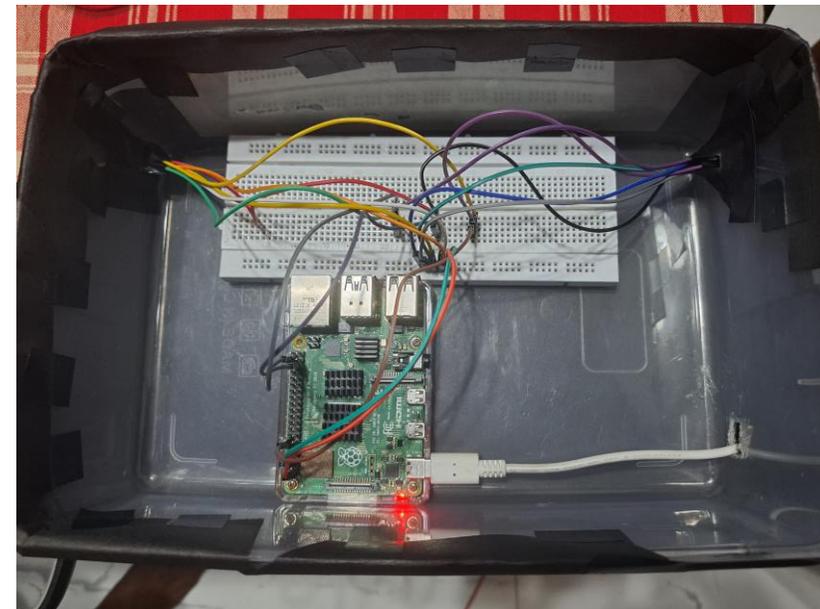
Test Accuracy: 0.9732

Proof of Completion

Evidence of Completion



```
const HornDetection = () => {
  useEffect(() => {
    } else {
      Animated.timing(rightScaleAnim, {
        toValue: 1,
        duration: 200,
        useNativeDriver: true,
      }).start();
    }
  }, [alert]);
  // Pulse with fade animation for alert box when a horn is detected
  if (alert !== 'No Horn Detected') {
    Animated.loop(
      Animated.sequence([
        Animated.parallel([
          Animated.timing(alertScaleAnim, {
            toValue: 1.05,
            duration: 800,
            useNativeDriver: true,
          }),
          Animated.timing(alertScaleAnim, {
            toValue: 1,
            duration: 800,
            useNativeDriver: true,
          }),
        ]),
        Animated.sequence([
          Animated.timing(alertFadeAnim, {
            toValue: 0.7,
            duration: 800,
            useNativeDriver: true,
          }),
          Animated.timing(alertFadeAnim, {
            toValue: 1,
            duration: 800,
            useNativeDriver: true,
          }),
        ]),
      ]),
    ).start();
  }
});
```



References

- [1] Gomez, R., & Lee, S. (2023). "Real-Time Acoustic Signal Processing for Assistive Technologies." *Journal of Assistive Technology and Accessibility Research*, 22(4), 401–419. Retrieved from <https://doi.org/10.1016/j.jatar.2023.04.015>.
- [2] Chen, H., & Kumar, P. (2023). "Machine Learning Approaches for Sound Recognition in Autonomous Vehicles." *IEEE Transactions on Intelligent Transportation Systems*, 25(5), 2450–2462. Retrieved from <https://doi.org/10.1109/TITS.2023.3056125>.
- [3] Garcia, M., & Rodriguez, F. (2023). "IoT Applications for Enhancing Road Safety: A Focus on Hearing Impairments." *Sensors and Actuators: A Physical*, 350, 113789. Retrieved from <https://doi.org/10.1016/j.sna.2023.113789>.
- [4] Smith, J., & Brown, A. (2023). "Assistive Technologies for Deaf Drivers: A Review." *Journal of Transportation Research*, 56(3), 201–215. Retrieved from <https://doi.org/10.1016/j.jtr.2023.03.007>.
- [5] Wang, L., & Zhang, Y. (2022). "IoT-Based Sound Detection Systems for the Hearing Impaired." *IEEE Transactions on Automation Science and Engineering*, 19(2), 1125–1134. Retrieved from <https://doi.org/10.1109/TASE.2022.3142501>.
- [6] Williams, R. (2021). *Intelligent Systems for Sound Recognition: From Theory to Implementation*. Berlin: Springer. ISBN 978-3-030-56789-0. Retrieved from <https://doi.org/10.1007/978-3-030-56790-6>.
- [7] Anderson, J. (2023). *Assistive Technology in Modern Transportation: Challenges and Innovations*. London: Wiley. ISBN 978-1-118-92783-1.
- [8] Davis, L. (2020). *Machine Learning Techniques for Real-Time Sound Classification*. London: Academic Press. ISBN 978-0-12834-5678.
- [9] National Institute on Deafness and Other Communication Disorders (NIDCD). (2023). "Driving and Transportation Options for People with Hearing Loss." Accessed August 5, 2024. <https://www.nidcd.nih.gov/health/driving-hearing-loss>.
- [10] IEEE Spectrum. (2023). "How IoT is Revolutionizing Road Safety for Hearing-Impaired Drivers." Accessed August 5, 2024. <https://spectrum.ieee.org/iot-road-safety-hearing-impaired>.
- [11] Transportation Research Board (TRB). (2023). "Advancements in Assistive Technologies for Drivers with Disabilities." Accessed August 5, 2024. <https://www.trb.org/assistive-technologies-drivers-disabilities>.



Thathsara S. M. K.
IT21219566

Software Engineering



Comprehensive Driver Behavior Monitoring System.



Background

❖ Multimodal Approaches in Driver Behavior Analysis:

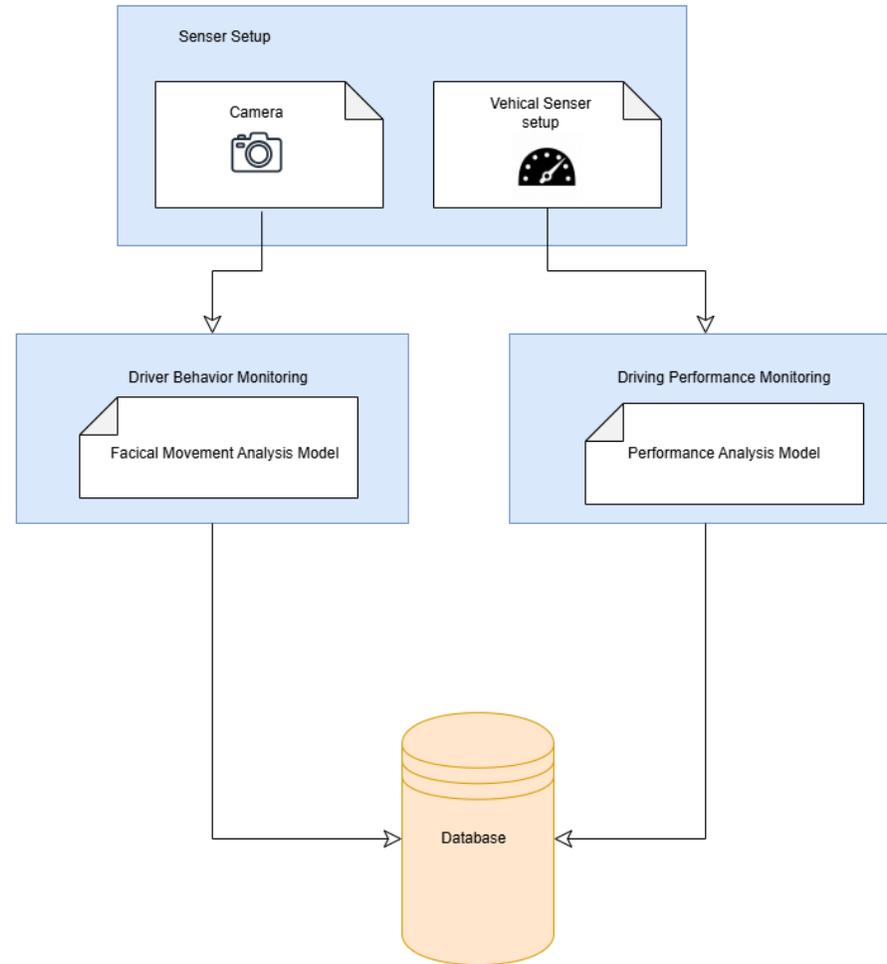
1. CNN for Visual Inputs

- **Purpose:** Process spatial information from video or image data.
- **Model Architecture:** CNNs for video sequences to capture spatiotemporal features.

2. RNN for Temporal Sensor Data

- **Purpose:** Handle sequential data such as time-series signals.
- **Model Architecture:** LSTMs (Long Short-Term Memory) for long-range temporal dependencies.

System Diagram



Data Collection

- Driver behavior data was collected from Google, Kaggle.

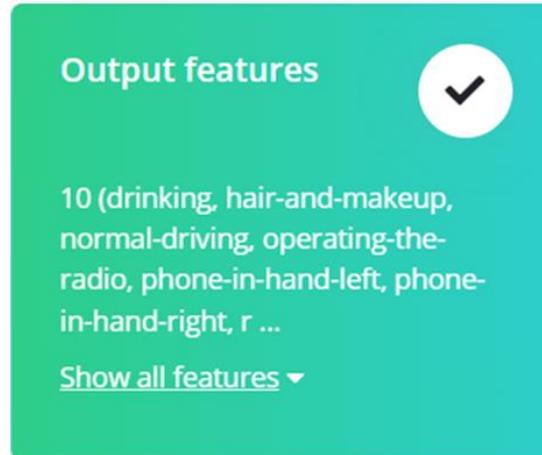
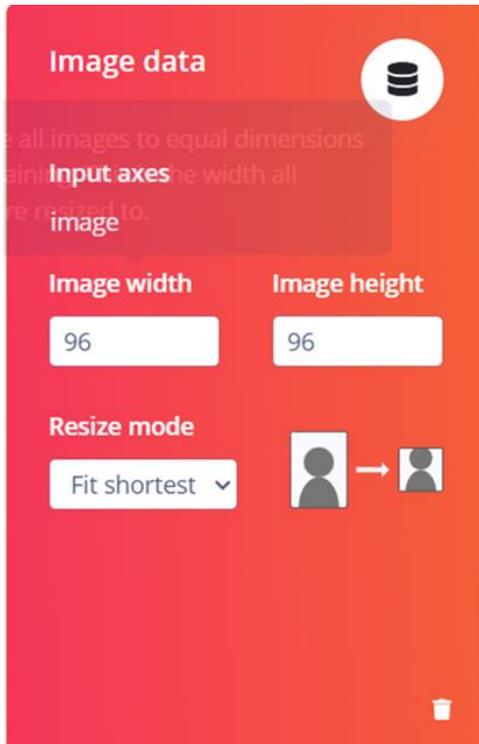
<https://www.kaggle.com/competitions/state-farm-distracted-driver-detection>

The screenshot shows the Kaggle dataset page for 'state-farm-distracted-driver-detection'. At the top, it displays 'DATA COLLECTED 22,429 items' and 'TRAIN / TEST SPLIT 79% / 21%'. Below this, there's a 'Dataset' section with a table of samples. The table has columns for 'SAMPLE NAME', 'LABEL', and 'ADDED'. The first row is highlighted in blue and shows 'img_99979' with the label 'phone-in-hand...' and 'Yesterday, 01:...'. To the right of the table, there's a 'Collect data' section with a 'Connect a device' button. Below that, there's a 'RAW DATA' section showing a video frame of a driver holding a phone while driving, labeled 'img_99979'.

SAMPLE NAME	LABEL	ADDED
img_99979	phone-in-hand...	Yesterday, 01:...
img_99751	phone-in-hand...	Yesterday, 01:...
img_99633	phone-in-hand...	Yesterday, 01:...
img_99822	phone-in-hand...	Yesterday, 01:...
img_9974	phone-in-hand...	Yesterday, 01:...
img_99452	phone-in-hand...	Yesterday, 01:...
img_99429	phone-in-hand...	Yesterday, 01:...
img_99379	phone-in-hand...	Yesterday, 01:...

Data Preprocessing

□ 10 Data classes



1. Normal driving
2. Texting on right hand
3. Talking on the phone right hand
4. Texting on left hand
5. Talking on the phone left hand
6. Operating the radio
7. Drinking / eating
8. Reaching behind
9. Hair and makeup
10. Talking to the passenger

Model Training

Neural Network settings

Training settings

Number of training cycles [?]

Use learned optimizer [?]

Learning rate [?]

Training processor [?]

Data augmentation [?]

Advanced training settings

Validation set size [?] %

Split train/validation set on metadata key [?]

Batch size [?]

Auto-weight classes [?]

Profile int8 model [?]



Model Accuracy



ACCURACY
91.93%

Metrics for Transfer learning ↓	
METRIC	VALUE
Area under ROC Curve ?	1.00
Weighted average Precision ?	0.95
Weighted average Recall ?	0.95
Weighted average F1 score ?	0.95

Confusion matrix

	DRINK	HAIR-ANI	NORM	OPERATING	PHONE-IP	PHONE-IP	REACHING	TALKING	TALKING	TALKING	UNCLEAR
DRINKING	93.2%	0.2%	0%	0.2%	0%	0.6%	0%	0.6%	0.6%	0.2%	4.3%
HAIR-ANI	1.8%	84.3%	0%	0.3%	0.3%	0.3%	1.3%	0.3%	0.8%	0.8%	10.0%
NORMAL	0%	0.4%	89.5%	0.2%	1.5%	0.4%	0.2%	0.2%	0%	0.4%	7.2%
OPERATING	0.2%	0%	0%	96.2%	0.4%	0%	0%	0.6%	0%	0.2%	2.3%
PHONE-IP	0%	0.2%	0.8%	0%	92.7%	0%	0%	0%	0%	0%	6.3%
PHONE-IP	0.4%	0.2%	0%	0%	0.2%	95.2%	0%	0%	0.2%	0%	3.8%
REACHING	0%	0.3%	0%	0%	0%	0.3%	95.4%	0.3%	0%	0.8%	3.0%
TALKING	0.2%	0.8%	0.6%	0.2%	0.8%	0%	0.4%	89.4%	0.2%	0.8%	6.4%
TALKING	0.6%	0.4%	0%	0%	0%	0.2%	0.2%	0%	94.6%	0%	3.9%
TALKING	0%	1.5%	1.2%	0.7%	0.5%	0.2%	0.5%	0.2%	0%	87.6%	7.5%
F1 SCORE	0.95	0.89	0.93	0.97	0.94	0.97	0.96	0.93	0.96	0.92	

Model Deployment

```
    )
# Capture webcam frames
cap = cv2.VideoCapture(0)

while True:
    ret, frame = cap.read()
    if not ret:
        break

    # Resize frame to match model input size (e.g., 96x96)
    resized_frame = cv2.resize(frame, (96, 96))

    # Classify the frame
    result = classify_image(resized_frame)
    print(result)

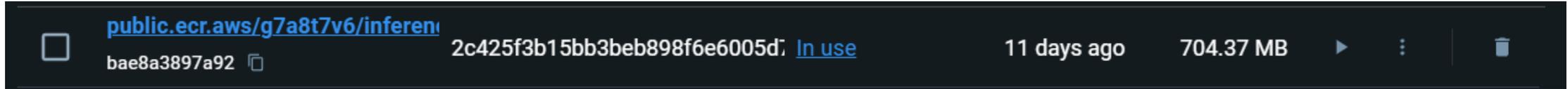
    # Store the result and timestamp in Firebase
    store_result(result)

    # Display the frame
    cv2.imshow('Webcam', frame)
    if cv2.waitKey(1) & 0xFF == ord('q'):
        break

cap.release()
cv2.destroyAllWindows()
```

```
DBM-Model-Run.py X
DBM-Model-Run.py > store_result
1 import cv2
2 import requests
3 import firebase_admin
4 from firebase_admin import credentials, firestore
5 from datetime import datetime
6
7 # Initialize Firebase
8 cred = credentials.Certificate("D:/Driver-Behavior-Monitoring-Model-Run/serviceAccountKey.json")
9 firebase_admin.initialize_app(cred)
10 db = firestore.client()
11
12 # Function to send image to inference server
13 def classify_image(image):
14     url = "http://localhost:1337/api/image"
15     _, img_encoded = cv2.imencode('.jpg', image)
16     files = {'file': ('image.jpg', img_encoded.tobytes(), 'image/jpeg')}
17     response = requests.post(url, files=files)
18     return response.json()
19
20 # Function to store result and timestamp in Firebase
21 def store_result(result):
22     # Get the predicted class (label with the highest confidence)
23     predicted_class = max(result['result']['classification'], key=result['result']['classification'].get)
24
25     # Store result and timestamp in Firebase
26     doc_ref = db.collection('detections').document() # Use .document() instead of .doc()
27     doc_ref.set({
28         'result': predicted_class,
29         'timestamp': datetime.now().isoformat()
30     })
31
32 # Capture webcam frames
33 cap = cv2.VideoCapture(0)
```

Model Deployment

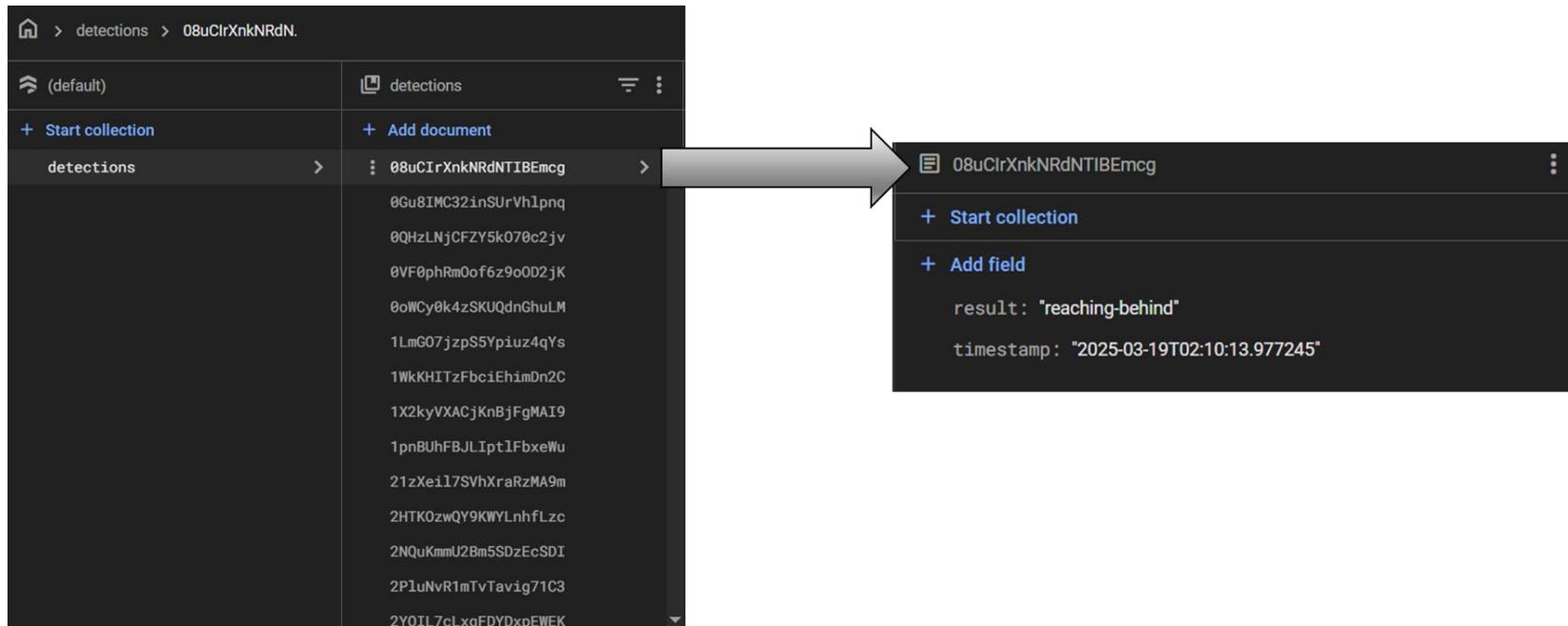


Docker container that includes an HTTP interface to run the model.

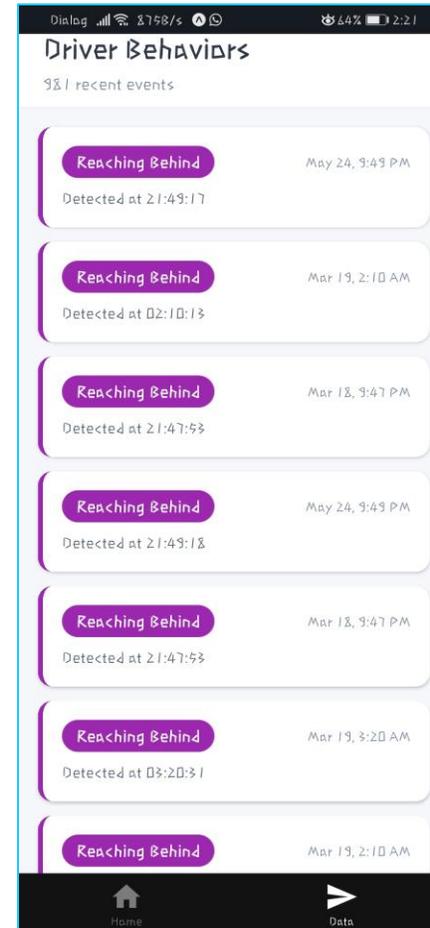
```
PS D:\Driver-Behavior-Monitoring-Model-Run> & C:/Python311/python.exe d:/Driver-Behavior-Monitoring-Model-Run/DBM-Model-Run.py
{'result': {'classification': {'drinking': 0.00018309261940885335, 'hair-and-makeup': 1.2963785138708772e-06, 'normal-driving': 3.904605705429276e-07, 'operating-the-radio': 4.260187438376306e-07, 'phone-in-hand-left': 1.039543462866277e-06, 'phone-in-hand-right': 0.00031643625698052347, 'reaching-behind': 0.9939215183258057, 'talking-on-the-phone-left': 6.087520765252208e-11, 'talking-on-the-phone-right': 0.0055737304501235485, 'talking-with-passenger': 2.089898316626204e-06}}, 'timing': {'anomaly': 0, 'classification': 32, 'dsp': 0, 'json': 19, 'stdin': 1}}
{'result': {'classification': {'drinking': 8.736484596738592e-05, 'hair-and-makeup': 7.89706646742161e-08, 'normal-driving': 4.275167719836048e-11, 'operating-the-radio': 8.759949565728675e-08, 'phone-in-hand-left': 1.4678050441752077e-10, 'phone-in-hand-right': 0.0001264539605472237, 'reaching-behind': 0.9997629523277283, 'talking-on-the-phone-left': 3.215795947876499e-15, 'talking-on-the-phone-right': 2.3131262423703447e-05, 'talking-with-passenger': 5.321054263873748e-09}}, 'timing': {'anomaly': 0, 'classification': 0, 'dsp': 0, 'json': 0, 'stdin': 0}}
```

Database Configuration

```
# Initialize Firebase
cred = credentials.Certificate("D:/Driver-Behavior-Monitoring-Model-Run/serviceAccountKey.json")
firebase_admin.initialize_app(cred)
db = firestore.client()
```



Mobile Application



References

1. A. Alamri, A. Gumaei, M. Al-Rakhami, M. M. Hassan, M. Alhussein and G. Fortino, "An Effective Bio-Signal-Based Driver Behavior Monitoring System Using a Generalized Deep Learning Approach," in IEEE Access, vol. 8, pp. 135037-135049, 2020, doi: 10.1109/ACCESS.2020.3011003.
2. Qu, Fangming & Dang, Nolan & Furht, Borko & Nojournian, Mehrdad. (2024). Comprehensive study of driver behavior monitoring systems using computer vision and machine learning techniques. Journal of Big Data. 11. 10.1186/s40537-024-00890-0.
3. H. -B. Kang, "Various Approaches for Driver and Driving Behavior Monitoring: A Review," 2013 IEEE International Conference on Computer Vision Workshops, Sydney, NSW, Australia, 2013, pp. 616-623, doi: 10.1109/ICCVW.2013.85. keywords: {Vehicles;Feature extraction;Monitoring;Accidents;Visualization;Head;Wheels;Driver drowsiness monitoring;Driving behavior monitoring;Unsafe driving behavior prediction},
4. V. Sanjay Kumar, S. Nair Ashish, I.V. Gowtham, S.P. Ashwin Balaji, E. Prabhu, Smart driver assistance system using raspberry pi and sensor networks, Microprocessors and Microsystems, Volume 79, 2020, 103275, ISSN 0141-9331,
5. Sameer Rafee, Xu Yun, Zhang Jian Xin and Zaid Yemeni, "Eye-movement Analysis and Prediction using Deep Learning Techniques and Kalman Filter" International Journal of Advanced Computer Science and Applications(IJACSA), 13(4), 2022. <http://dx.doi.org/10.14569/IJACSA.2022.01304107>



Iroshan G.H.M
IT21229084

Software Engineering



Communication Support System for Deaf Drivers

INTRODUCTION

- LipAssist is a **lip reading model** designed to assist the deaf and hard-of-hearing community
- It uses **deep learning** to interpret lip movements and convert them into text.

RESEARCH QUESTION



How can deaf drivers share emergency details effectively without spoken communication ?



What delays do deaf drivers face in getting help during emergencies ?

OBJECTIVES

Specific Objective

To create a **real-time lip-reading tool** that helps deaf drivers **understand spoken instructions** during emergencies, ensuring **safety** and **clear communication**.



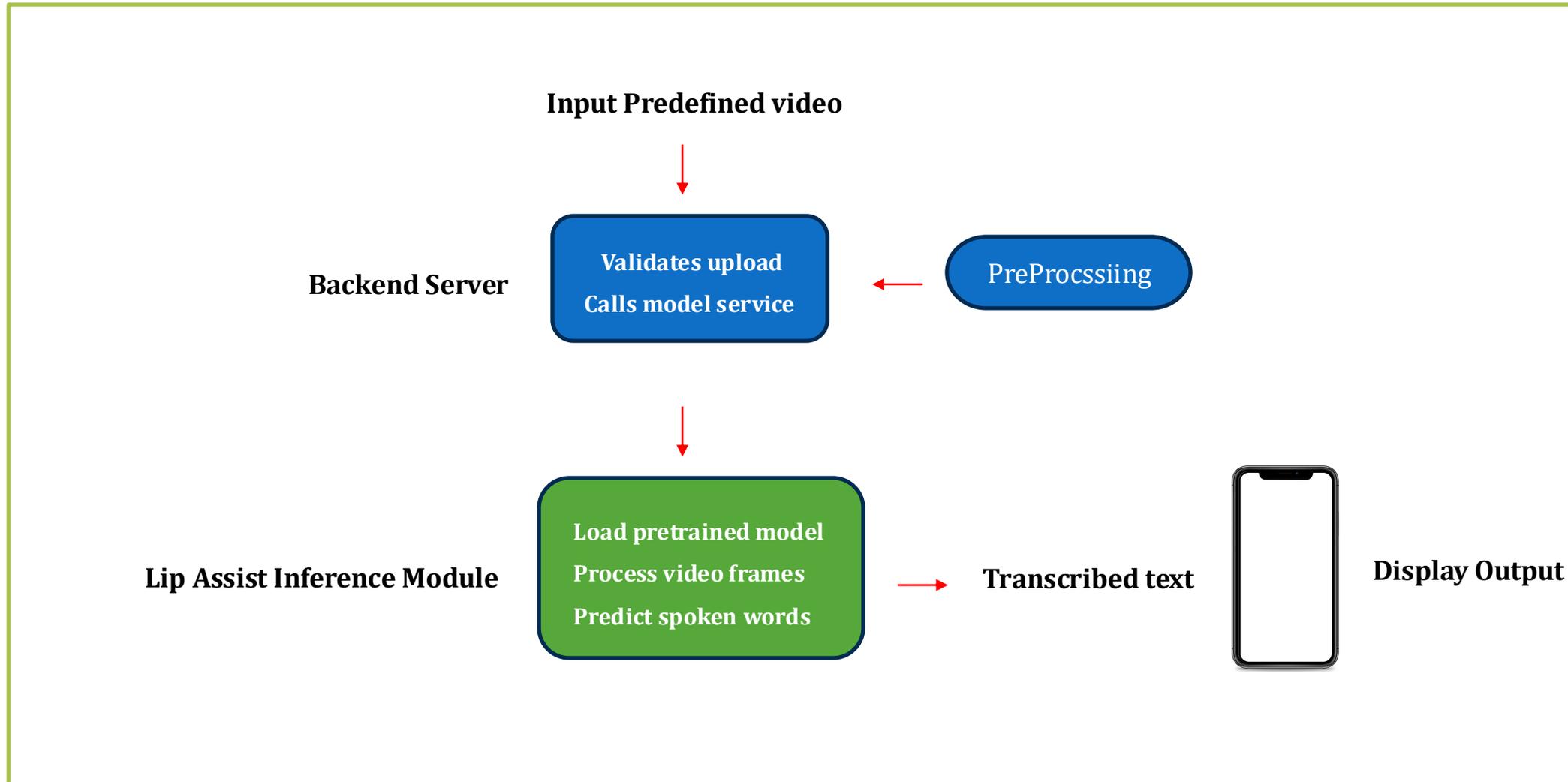
Sub Objective

Communication Facilitation

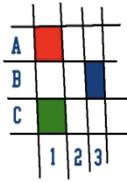
lip reading model

Emergency scenarios

Offline Capability



Data collection



The GRID audiovisual sentence corpus

[What is GRID?](#) | [Examples](#) | [Downloading](#) | [Documentation](#) | [Credits](#)

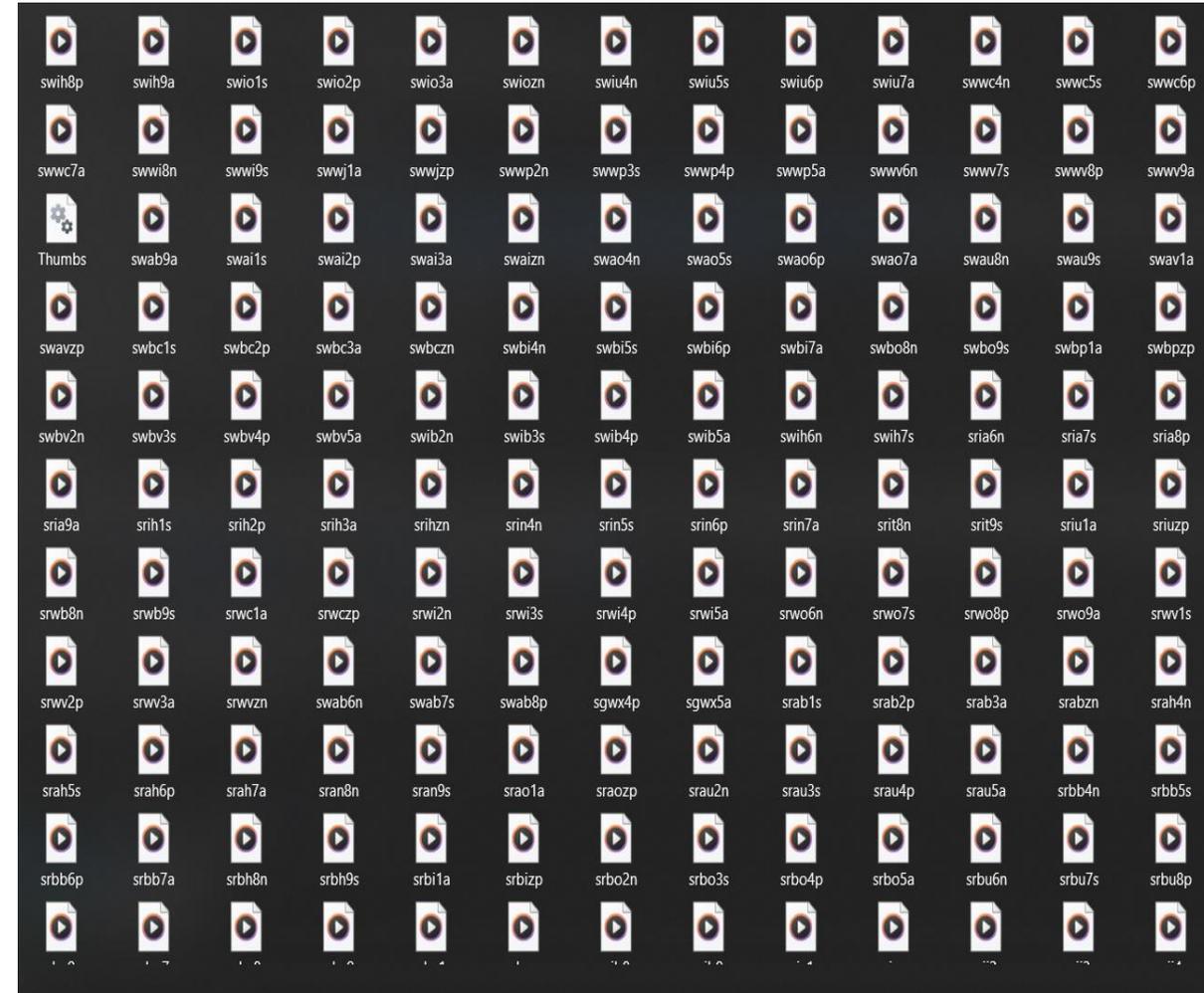
What is GRID?

GRID is a large multitalker audiovisual sentence corpus to support joint computational-behavioral studies in speech perception. In brief, the corpus consists of high-quality audio and video (facial) recordings. The corpus, together with transcriptions, is freely available for research use. GRID is described in more detail in this [paper](#).

Examples

talker	audio only	video (normal)	video (high)	transcriptions
male	download	download	download	download
female	download	download	download	download

Alignments



Dataset

Word Structure

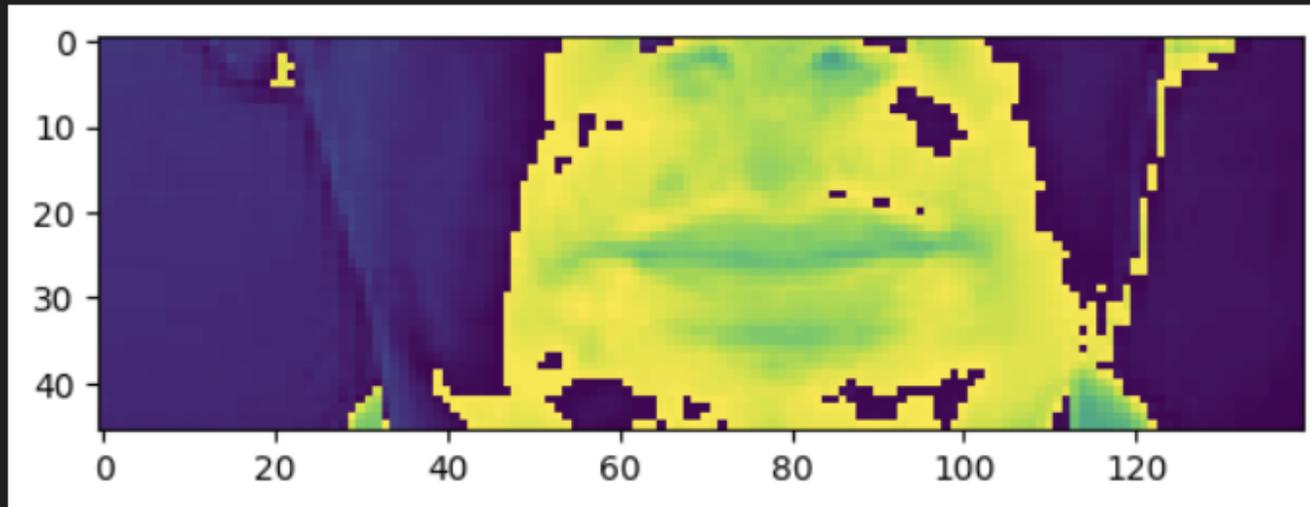
Word Type	Examples
Command	bin, lay, place, set
Color	blue, green, red, white
Preposition	at, by, in, with
Letter	A to Z (excluding W)
Digit	0 to 9

Alignment Data Visualization

```
print(f'This is the alignments of the one video: {alignments}')  
print(' ')  
tf.strings.reduce_join([bytes.decode(i) for i in num_to_char(alignments.numpy()).numpy()])  
plt.imshow(frames[74])
```

This is the alignments of the one video: [2 9 14 39 2 12 21 5 39 1 20 39 6 39 20 23 15 39 14 15 23]

<matplotlib.image.AxesImage at 0x7eac63db5570>



Video Frame Data Representation & Predictions

```
print(f'whole the frame: {frames.shape}')
frames[[0]]

whole the frame: (75, 46, 140, 1)

<tf.Tensor: shape=(46, 140, 1), dtype=float32, numpy=
array([[1.4991663 ],
       [1.4991663 ],
       [1.4616871 ],
       ...,
       [0.41227072],
       [0.41227072],
       [0.41227072]],

       [[1.4991663 ],
       [1.4991663 ],
       [1.4616871 ],
       ...,
       [0.41227072],
       [0.41227072],
       [0.41227072]],

       [[1.4616871 ],
       [1.4616871 ],
       [1.4991663 ],
       ...,
       [0.3373124 ],
       [0.3373124 ],
       [0.3373124 ]],

       ...

       [1.0119373 ],
       ...,
       [0.07495831],
       [0.07495831],
       [0.03747915]]], dtype=float32)>
```

```
Epoch 123: val_loss did not improve from 1.43750
1/1 ————— 0s 122ms/step
Original:  set white at u nine soon
Predicted:  set white at u nine son
~~~~~
Original:  set white by i four now
Predicted:  set white by i four now
```

```
Epoch 129: val_loss improved from 1.38357 to 1.31727, saving model to best_weights.weights.h5
1/1 ————— 0s 27ms/step
Original:  set red by h eight now
Predicted:  set red by h eight now
~~~~~
Original:  set blue at n two now
Predicted:  set blue at t two now
```

Model Architecture Summary

```
model.summary()
```

Model: "sequential"

Layer (type)	Output Shape	Param #
conv3d (Conv3D)	(None, 75, 46, 140, 128)	3,584
batch_normalization (BatchNormalization)	(None, 75, 46, 140, 128)	512
activation (Activation)	(None, 75, 46, 140, 128)	0
max_pooling3d (MaxPooling3D)	(None, 75, 23, 70, 128)	0
conv3d_1 (Conv3D)	(None, 75, 23, 70, 256)	884,992
batch_normalization_1 (BatchNormalization)	(None, 75, 23, 70, 256)	1,024
activation_1 (Activation)	(None, 75, 23, 70, 256)	0
max_pooling3d_1 (MaxPooling3D)	(None, 75, 11, 35, 256)	0
conv3d_2 (Conv3D)	(None, 75, 11, 35, 75)	518,475
batch_normalization_2 (BatchNormalization)	(None, 75, 11, 35, 75)	300
activation_2 (Activation)	(None, 75, 11, 35, 75)	0
max_pooling3d_2 (MaxPooling3D)	(None, 75, 5, 17, 75)	0
time_distributed (TimeDistributed)	(None, 75, 6375)	0
bidirectional (Bidirectional)	(None, 75, 256)	6,660,096
dropout (Dropout)	(None, 75, 256)	0
bidirectional_1 (Bidirectional)	(None, 75, 256)	394,240
dropout_1 (Dropout)	(None, 75, 256)	0
dense (Dense)	(None, 75, 41)	10,537

Total params: 8,473,760 (32.32 MB)

Trainable params: 8,472,842 (32.32 MB)

Non-trainable params: 918 (3.59 KB)

Backend Implementation

```
def load_model() -> Sequential:
    model = Sequential()

    # First Conv3D block
    model.add(Conv3D(128, 3, input_shape=(75, 46, 140, 1), padding='same'))
    model.add(BatchNormalization()) # Added BatchNormalization
    model.add(Activation('relu'))
    model.add(MaxPool3D(pool_size=(1, 2, 2)))

    # Second Conv3D block
    model.add(Conv3D(256, 3, padding='same'))
    model.add(BatchNormalization()) # Added BatchNormalization
    model.add(Activation('relu'))
    model.add(MaxPool3D(pool_size=(1, 2, 2)))

    # Third Conv3D block
    model.add(Conv3D(75, 3, padding='same'))
    model.add(BatchNormalization()) # Added BatchNormalization
    model.add(Activation('relu'))
    model.add(MaxPool3D(pool_size=(1, 2, 2)))

    # Flatten the spatial dimensions while preserving the time dimension
    model.add(TimeDistributed(Reshape((-1,))))

    # First LSTM block
    model.add(Bidirectional(LSTM(128, kernel_initializer='orthogonal', return_sequences=True)))
    model.add(Dropout(0.5))

    # Second LSTM block
    model.add(Bidirectional(LSTM(128, kernel_initializer='orthogonal', return_sequences=True)))
    model.add(Dropout(0.5))

    # Fully connected output layer
    model.add(Dense(char_to_num.vocabulary_size() + 1, kernel_initializer='he_normal', activation='softmax'))

    # model.load_weights('../models/best_weights.weights.h5')
    # Load the weights
    model.load_weights(os.path.join('models', 'new_best_weights2.weights.h5'))
    return model
```

```
def load_video(path:str) -> List[float]:
    cap = cv2.VideoCapture(path)
    frames = []
    for _ in range(int(cap.get((cv2.CAP_PROP_FRAME_COUNT)))):
        ret, frame = cap.read()
        frame = tf.image.rgb_to_grayscale(frame)
        frames.append(frame[190:236, 80:220, :])
    cap.release()

    mean = tf.math.reduce_mean(frames)
    std = tf.math.reduce_std(tf.cast(frames, tf.float32))
    return tf.cast((frames - mean), tf.float32)/std
```

```
def load_data(path: str):
    path = bytes.decode(path.numpy())
    print(f'This is the path: {path}')
    file_name = path.split('\\')[-1].split('.')[0]
    # Extract the file name without extension
    file_name = os.path.splitext(os.path.basename(path))[0]

    # Construct video and alignment paths
    video_path = os.path.join('data', 's1', f'{file_name}.mpg')
    alignment_path = os.path.join('data', 'alignments', 's1', f'{file_name}.align')

    # Ensure files exist
    if not os.path.exists(video_path):
        raise FileNotFoundError(f"Video file not found: {video_path}")
    if not os.path.exists(alignment_path):
        raise FileNotFoundError(f"Alignment file not found: {alignment_path}")

    # Load video frames and alignments
    frames = load_video(video_path)
    alignments = load_alignments(alignment_path)

    return frames, alignments
```

```
def load_alignments(path:str) -> List[str]:
    with open(path, 'r') as f:
        lines = f.readlines()
        tokens = []
        for line in lines:
            line = line.split()
            if line[2] != 'sil':
                tokens = [*tokens, ' ', line[2]]
        return char_to_num(tf.reshape(tf.strings.unicode_split(tokens, input_encoding='UTF-8'), (-1,)))[1:]
```

Testing Using A WebApp

Welcome to DeafAssist!

DeafAssist is designed to enhance communication for deaf drivers in emergency situations by leveraging real-time lip reading and voice-to-text conversion.

Architecture Flow:

- Input Video Frames (T x H x W x C)
- 3D Convolutional Layers (Spatiotemporal Feature Extraction)
- Bidirectional GRUs (Temporal Modeling)
- Fully Connected Layer (Vocabulary Prediction)
- CTC Loss (Sequence Alignment)
- Output Sequence (Words or Phonemes)

How to Use:

- Upload or select a pre-loaded video; Choose a video to test real-time lip reading.
- Video Preview: View the video on the left.

Select a video from the dataset below:

bbaf2n.mpg

Original Video Preview

Model Input Visualization

Model loaded successfully!

Raw Model Predictions

Raw argmax tokens:

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
58	16	6	21	21	11	21	21	27	1	21	21	7	21	2	34	21	21	2

CTC Decoded Tokens

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
2	9	14	39	2	12	21	5	39	1	20	39	6	39	20	23	15	39	1

Decoded Speech

Predicted Transcript: b1s blue at f two now

```
[out#0/mp4 @ 0000018b7daf2280] video:62KiB audio:46KiB subtitle:0KiB other streams:0KiB global headers:0KiB muxing overhead: 3.336924%
frame= 75 fps=0.0 q=-1.0 Lsize= 111KiB time=00:00:02.92 bitrate= 312.7kbits/s speed= 15x
[libx264 @ 0000018b7ded0a80] frame I:1 Avg QP:21.85 size: 6320
[libx264 @ 0000018b7ded0a80] frame P:47 Avg QP:21.91 size: 1020
[libx264 @ 0000018b7ded0a80] frame B:27 Avg QP:23.47 size: 297
[libx264 @ 0000018b7ded0a80] consecutive B-frames: 49.3% 5.3% 8.0% 37.3%
[libx264 @ 0000018b7ded0a80] mb I I16..4: 32.4% 66.9% 0.7%
[libx264 @ 0000018b7ded0a80] mb P I16..4: 3.2% 2.7% 0.0% P16..4: 40.8% 8.9% 6.7% 0.0% 0.0% skip:37.8%
[libx264 @ 0000018b7ded0a80] mb B I16..4: 0.7% 0.8% 0.0% B16..8: 30.9% 1.1% 0.1% direct: 1.3% skip:65.2% L0:53.1% L1:45.6% BI: 1.4%
[libx264 @ 0000018b7ded0a80] 8x8 transform intra:51.9% inter:83.9%
[libx264 @ 0000018b7ded0a80] coded y,uvDC,uvAC intra: 27.3% 67.7% 8.8% inter: 9.4% 25.0% 1.1%
[libx264 @ 0000018b7ded0a80] i16 v,h,dc,p: 14% 30% 13% 43%
[libx264 @ 0000018b7ded0a80] i8 v,h,dc,ddl,ddr,vr,hd,vl,hu: 18% 29% 41% 2% 1% 2% 1% 2% 3%
[libx264 @ 0000018b7ded0a80] i4 v,h,dc,ddl,ddr,vr,hd,vl,hu: 17% 31% 6% 4% 2% 38% 2% 0% 0%
[libx264 @ 0000018b7ded0a80] i8c dc,h,v,p: 42% 26% 26% 6%
[libx264 @ 0000018b7ded0a80] Weighted P-Frames: Y:0.0% UV:0.0%
[libx264 @ 0000018b7ded0a80] ref P L0: 64.8% 7.6% 18.4% 9.2%
[libx264 @ 0000018b7ded0a80] ref B L0: 79.2% 14.6% 6.2%
[libx264 @ 0000018b7ded0a80] ref B L1: 94.6% 5.4%
[libx264 @ 0000018b7ded0a80] kb/s:166.16
[aac @ 0000018b7daedc80] Qavg: 1083.495
This is the path: data\s1\bbaf2n.mpg
WARNING:tensorflow:5 out of the last 5 calls to <function TensorFlowTrainer.make_predict_function.<locals>.one_step_on_data_distributed at 0x00
f.function retracing. Tracing is expensive and the excessive number of tracings could be due to (1) creating @tf.function repeatedly in a loop,
ifferent shapes, (3) passing Python objects instead of tensors. For (1), please define your @tf.function outside of the loop. For (2), @tf.func
rue option that can avoid unnecessary retracing. For (3), please refer to https://www.tensorflow.org/guide/function#controlling_retracing and h
pi_docs/python/tf/function for more details.
1/1 _____ 2s 2s/step
```

