# Lip reading using deep learning

PROJECT REPORT

TEAM--591865                    11/22/23

# 1.  <u>INTRODUCTION</u>

## 1.1   Project Overview:

The "End-to-End Lip Reading Deep Learning" project is an innovative exploration into the realm of machine learning, aiming to enhance speech recognition through the integration of lip reading capabilities. This project leverages cutting-edge deep learning algorithms, such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), to create a comprehensive solution capable of accurately detecting and transcribing words from videos of individuals speaking.

Lip reading, as an auxiliary tool to traditional audio-based speech recognition, offers the potential to overcome challenges posed by noisy environments or unclear audio signals. By developing an end-to-end system that analyzes facial movements, this project seeks to improve the overall accuracy and robustness of speech recognition systems, making them more adaptable to diverse real-world scenarios.

The project adopts a multi-faceted technological approach, combining the power of Convolutional Neural Networks (CNNs) for spatial feature extraction from lip images and Recurrent Neural Networks (RNNs) or Long Short-Term Memory (LSTM) networks to capture temporal dependencies in the sequence of frames. This synergy allows the model to interpret the nuanced dynamics of lip movements over time, ultimately contributing to more accurate word predictions.

## 1.2    Purpose:

The primary purpose of this project is to address several key objectives:

a. Improved Speech Recognition
The integration of lip reading with conventional audio-based speech recognition systems aims to significantly enhance accuracy, especially in challenging conditions where audio signals may be compromised. This combined approach promises to provide a more reliable and resilient solution for recognizing spoken words.

b. Elimination of Audio Data Dependency
Unlike traditional speech recognition models that heavily rely on transcribed audio data, the proposed end-to-end lip reading system operates solely on video data. This eliminates the need for extensive audio transcription, which can be both expensive and time-consuming to obtain. The independence from audio data streamlines the training process and increases the accessibility of the system.

c. Multi-Modal Applications
The project envisions the creation of multi-modal applications by seamlessly integrating lip reading with audio-based systems. This integration holds the potential to revolutionize real-time communication, particularly in scenarios like video conferencing, by providing more accurate and context-aware transcriptions.

In summary, this project aspires to contribute to the advancement of speech recognition technology, offering a versatile solution with potential applications in diverse fields while prioritizing accessibility and inclusivity for individuals with hearing impairments.

# 2.LITERATURE SURVEY

## 2.1.Existing problem:

The field of lip reading and its integration into machine learning systems has witnessed considerable attention in recent years due to its potential to enhance speech recognition in challenging environments. Existing speech recognition systems predominantly rely on audio data, making them susceptible to issues such as background noise, speaker accents, and unclear audio signals. In response to these challenges, researchers have explored the integration of lip reading as a complementary modality to improve overall system performance.

Lip reading systems have been implemented using various approaches, including traditional computer vision techniques and, more recently, deep learning methods. While traditional methods often struggle to capture the complex and dynamic nature of lip movements, deep learning models, particularly those combining Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), have shown promise in addressing these challenges.

Despite advancements, existing literature highlights several ongoing challenges, such as:

**a. Variability in Lip Movements:**
   - Different speakers exhibit diverse lip movements, making it challenging to develop a universally robust model.

**b. Limited Availability of Diverse Datasets:**
   - Many studies highlight the scarcity of large and diverse datasets for training lip reading models, affecting their generalization to real-world scenarios.

**c. Real-time Processing Requirements:**
   - Achieving real-time processing for practical applications remains a challenge, especially when deploying lip reading in dynamic environments.

## 2.2 References

1. Petridis, S., Stavropoulos, G., Liapis, A., & Cavouras, D. (2018). "Deep recurrent neural networks for lipreading: A study on weakly supervised learning." Computer Speech & Language, 50, 66-95.

2. Assael, Y. M., Shillingford, B., Whiteson, S., & de Freitas, N. (2016). "LipNet: End-to-End Sentence-level Lipreading." arXiv preprint arXiv:1611.01599.

3. Chung, J. S., Senior, A., Vinyals, O., & Zisserman, A. (2017). "Lip Reading Sentences in the Wild." In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017.

4. Wand, M., Gehler, P., & Schiele, B. (2016). "Lip Reading with Long Short-Term Memory." In European Conference on Computer Vision (ECCV), 2016.

## 2.3 Problem Statement Definition

The primary problem addressed in this project is the need for a robust and real-time end-to-end lip reading system that can effectively complement traditional audio-based speech recognition. The identified challenges include:

**a. Improving Robustness:**
   - Designing a model that can handle the inherent variability in lip movements across different speakers and scenarios, enhancing the robustness of the lip reading system.

**b. Dataset Diversity:**
   - Addressing the limitation of available diverse datasets to ensure the model's ability to generalize across various speaking styles, accents, and environmental conditions.
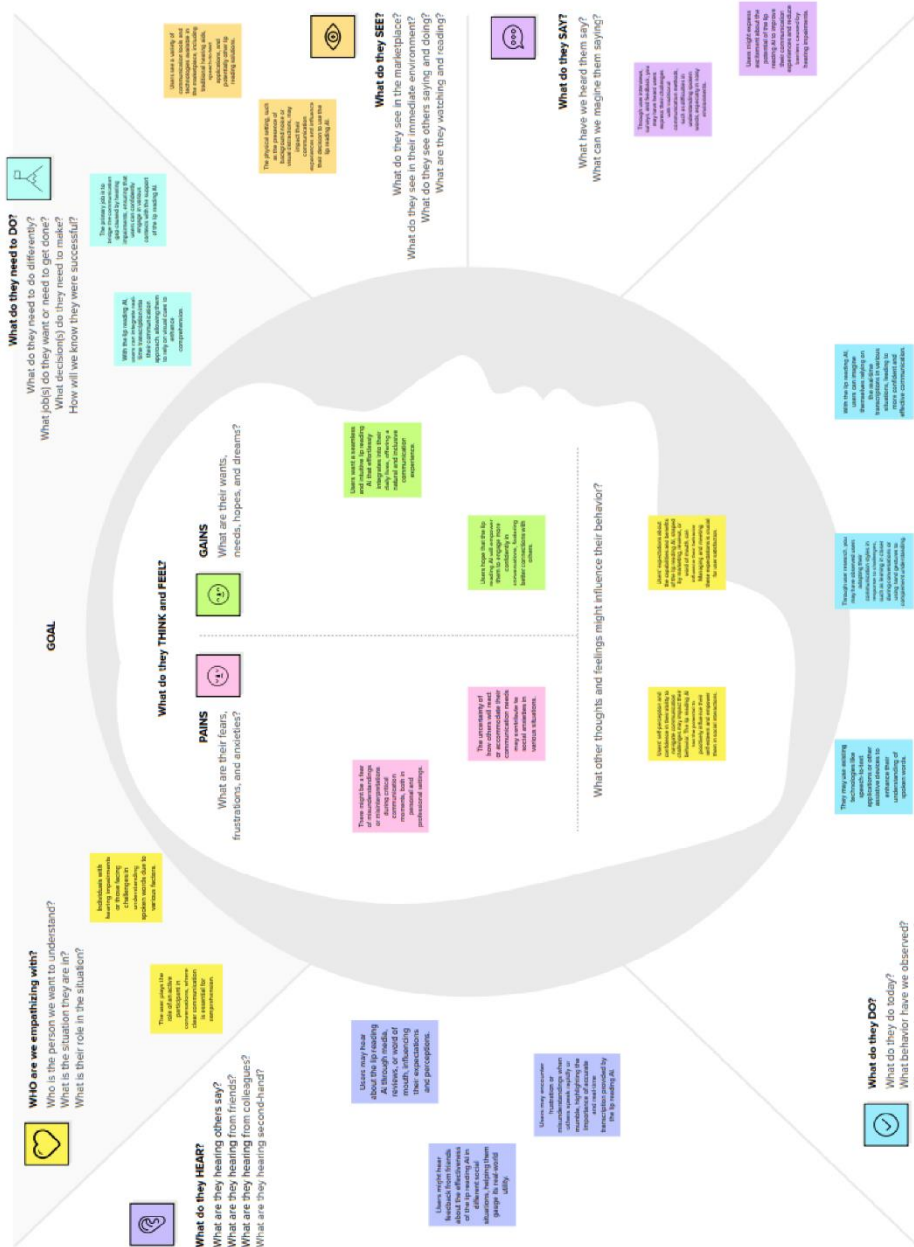
**c. Real-time Processing:**
   - Developing mechanisms for achieving real-time processing to enable the practical deployment of the lip reading system in dynamic and interactive applications.

By addressing these challenges, the project aims to contribute to the advancement of assistive technologies, making speech recognition more inclusive and effective in real-world scenarios.

# 3.IDEATION & PROPOSED SOLUTION

## 3.1.Empathy Map Canvas:

# What do they THINK and FEEL?

## PAINS
What are their fears, frustrations, and anxieties?

## GAINS
What are their wants, needs, hopes, and dreams?

Fears regarding how market volatility and economic fluctuations may impact the effectiveness.

A common desire is to improve customer satisfaction.

Anxiety about disruptions while implementing new strategies or technologies related to customer segmentation.

Many Businesses need to achieve sustainable growth by continuously adapting and evolving their strategies.

Fears about not having the necessary expertise or proper talent within the organization to perform effectively.

Dream of delivering exceptional experiences that meet or exceed the expectations of different customer segments.

# 3.2.Ideation & Brainstorming

**Step-2: Brainstorm, Idea Listing and Grouping**

**Brainstorm**

Write down any ideas that come to mind that address your problem statement.

⏱ 10 minutes

TIP
You can select a sticky note and hit the pencil (switch to sketch) icon to start drawing!

**Soma Sekhar**

Develop a user-friendly mobile application that utilizes deep learning to provide real-time, accurate lip reading transcriptions.

Implement a visual feedback system using augmented reality to highlight and enhance lip movements for better understanding.

Integrate the lip reading AI with virtual assistants for hands-free, voice-activated control, enhancing accessibility in various environments.

**Krishna Kowshik**

Design a discreet wearable device equipped with a camera and deep learning algorithms for on-the-go lip reading assistance.

Introduce personalized user profiles allowing individuals to tailor the lip reading AI to their specific needs and preferences.

Expand functionality to include real-time language translation, facilitating communication in multilingual settings.

**Manoj Kumar**

Develop interactive modules within the app for users to practice and improve their lip reading skills over time.

Create an online platform where users can share experiences, tips, and challenges, fostering a supportive community.

Implement an offline mode allowing users to access basic lip reading features without the need for a constant internet connection.

**Group ideas**

Take turns sharing your ideas while clustering similar or related notes as you go. Once all sticky notes have been grouped, give each cluster a sentence-like label. If a cluster is bigger than six sticky notes, try and see if you and break it up into smaller sub-groups.

⏱ 20 minutes

TIP
Add customizable tags to sticky notes to make it easier to find, browse, organize, and categorize important ideas as themes within your mural

Implement data privacy and compliance measures to protect user data

Optimize revenue and proftability.

Investigate data sources and tools that may enhance feasibility

Identify unusual behavior within customer segments

**Step-3: Idea Prioritization**

**Prioritize**

Your team should all be on the same page about what's important moving forward. Place your ideas on this grid to determine which ideas are important and which are feasible.

⏱ 20 minutes

TIP
Participants can use their cursors to point at where sticky notes should go on the grid. The facilitator can confirm the spot by using the laser pointer holding the H key on the keyboard.

Implement an offline mode allowing users to access basic lip reading features without the need for a constant internet connection.

Expand functionality to include real-time language translation, facilitating communication in multilingual settings.

**Importance**
If each of these tasks could get done without any difficulty or cost, which would have the most positive impact?

Implement data privacy and compliance measures to protect user data

Investigate data sources and tools that may enhance feasibility

**Feasibility**
Regardless of their importance, which tasks are more feasible than others? (Cost, time, effort, complexity, etc.)

# 4.REQUIREMENT ANALYSIS

## 4.1.Functional requirement:

**a. Video Upload and Processing:**

Users should be able to upload video files through the web-based interface.

The system must process uploaded videos, extracting frames for analysis.

**b. Real-time Lip Reading:**

The system should perform real-time lip reading on the uploaded video,

providing word predictions.

**c. Prediction Showcase:**

The UI must display the predicted words in real-time as the lip reading analysis

progresses.

**d. Model Training Interface (Admin Functionality):**

Admins should have an interface for retraining or fine-tuning the model with

additional data.

**e. Continuous Improvement Mechanism:**

The system should support continuous improvement, allowing for model

updates based on user feedback and additional training data.

**f. User Feedback Mechanism:**

Users should have the option to provide feedback on the accuracy of

predictions, contributing to model refinement.

**g. Accessibility Features:**

The UI must be designed with accessibility features, ensuring usability for

individuals with diverse abilities.

### 4.1.Non-Functional requirements:

**a. Performance:**

The system should provide accurate predictions with a reasonable processing time, even for longer video files.

**b. Scalability:**

The architecture should be scalable to handle an increasing number of users and potential future data expansion.

**c. Security:**

The system must implement security measures to protect user data and ensure the confidentiality of lip reading predictions.

**d. User Interface Responsiveness:**

The UI should be responsive, providing a seamless and intuitive experience for users.

**e. Model Accuracy:**

The lip reading model should achieve a high level of accuracy in predicting words, especially in diverse and challenging scenarios.

**f. Ethical Considerations:**

The system should adhere to ethical standards, including user privacy, fairness, and transparency in how predictions are generated.

**g. Compatibility:**

The web-based interface should be compatible with commonly used browsers to ensure accessibility for a wide user base.

**h. Maintainability:**

The system should be designed for ease of maintenance, with clear documentation and modular components for straightforward updates.
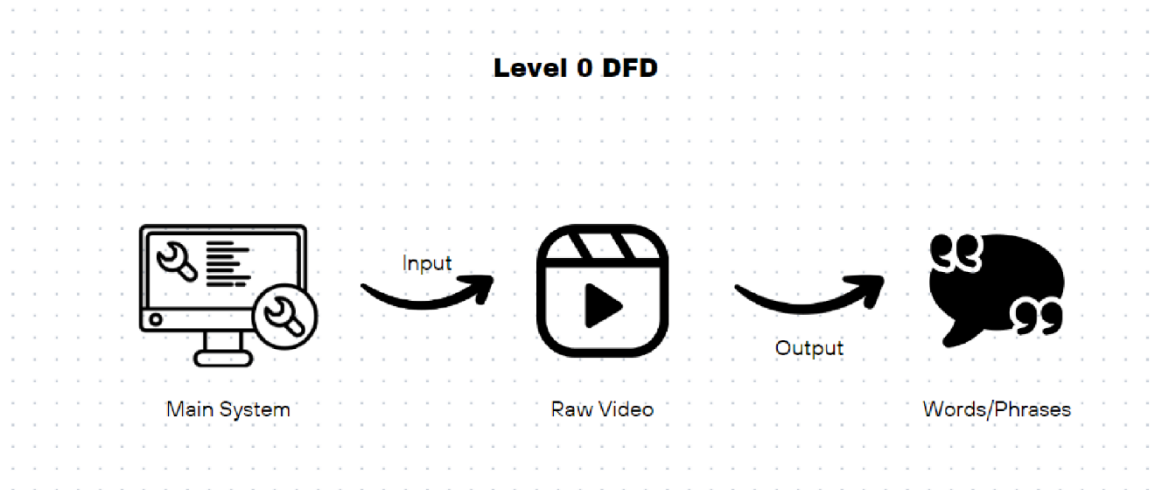
**i. User Acceptance:**

The system should undergo user acceptance testing to ensure that it meets user expectations and is intuitive for a diverse user group.

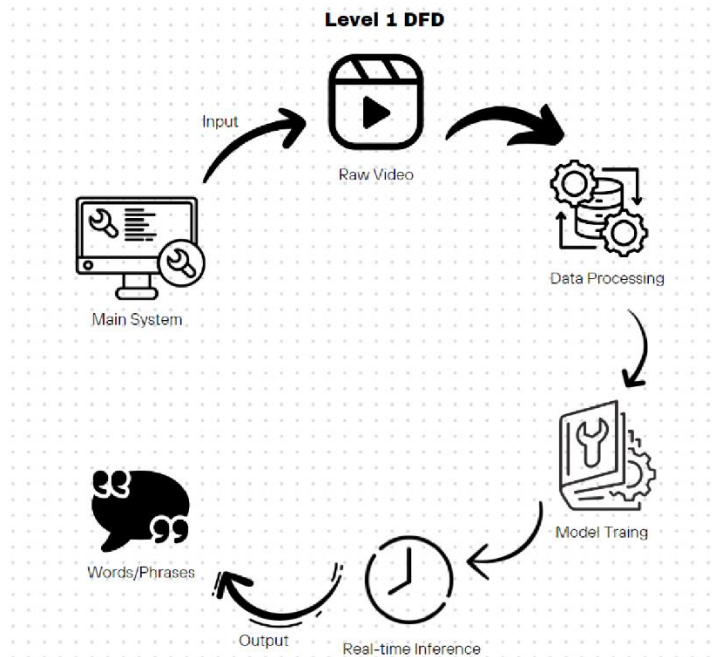These functional and non-functional requirements form the basis for developing and evaluating the success of the end-to-end lip reading system. Adhering to these specifications will contribute to the system's effectiveness, usability, and ethical deployment.

# 5.PROJECT DESIGN

## 5.1.Data Flow Diagrams & User Stories:

**Example:**

**Level 0 DFD**

Main System — Input → Raw Video — Output → Words/Phrases

**Level 1 DFD**

## User Stories

| User Type | Functional Requirement (Epic) | User Story Number | User Story / Task | Acceptance criteria | Priority | Release |
|---|---|---|---|---|---|---|
| User with Hearing Impairment | Speech Transcription | USN-1 | As a user, I want the lip reading system to accurately transcribe spoken words from videos into text in real-time, enhancing my comprehension and facilitating better communication. | The system should transcribe spoken words from videos with at least 90% accuracy in real-time. | High | Sprint-1 |
| Developer | System Integration | USN-2 | As a developer, I need a system that efficiently processes video data, trains models using deep learning algorithms, and provides APIs for seamless integration into various applications. | The system should have clear documentation and APIs that allow easy integration for different applications. | High | Sprint-1 |
| Service Provider | Platform Integration | USN-3 | As a service provider, I want to integrate this lip reading technology into our communication platform to offer real-time transcription services, improving accessibility and inclusivity for our users. | The technology should seamlessly integrate into our platform's existing interface and provide real-time transcription services. | Low | Sprint-1.1 |
| Researcher | Dataset Accessibility | USN-4 | As a researcher, I require access to a comprehensive dataset and an efficient lip reading system for studying speech recognition patterns, aiding in further advancements in the field. | The system should provide access to a diverse and well-annotated dataset suitable for research purposes. | Medium | Sprint-1.1 |
| System Administrator | Error Handling & Monitoring | USN-5 | As a system administrator, I aim to ensure the system's stability, implementing robust error handling mechanisms and monitoring tools to swiftly identify and resolve issues for uninterrupted service. | The system should log errors, provide real-time monitoring, and send alerts for any system malfunctions or downtime. | High | Sprint-1.1 |
| User in Noisy Environments | Clarity in Noisy Environments | USN-6 | As a user in noisy environments, I expect the lip reading system to accurately interpret lip movements for clear communication, providing an alternative method when audio is unclear or compromised | The system should maintain at least 80% accuracy in interpreting lip movements in noisy environments. | Medium | Sprint-1.2 |
| Content Creator | Video Transcription | USN-7 | As a content creator, I seek a reliable lip reading tool that accurately transcribes videos, enabling me to offer captions or subtitles for a wider audience, enhancing accessibility and engagement. | The system should generate accurate transcriptions for videos with different accents and speech patterns. | High | Sprint-1.2 |
| | | | | | | |
| | | | | | | |

## 5.2.Solution Architecture:

### Example - Solution Architecture Diagram:

```
┌─────────────────────────┐                    ┌─────────────────────────┐
│   Data Preprocessing     │ · · · · · · · · ·  │   Model Architecture     │
│                          │                    │                          │
│   • Video Input          │                    │   • LSTM Network         │
│   • Frame Extraction     │                    │   • (Deep Learning)      │
│   • Preprocessing        │                    │                          │
│   • Normalization        │                    │                          │
└─────────────────────────┘                    └─────────────────────────┘
              ·                                               ·
              ·                                               ·
              ·                                               ·
              ·                                               ·
┌─────────────────────────┐                    ┌─────────────────────────┐
│   Training Pipeline      │ · · · · · · · · ·  │   Inference Pipeline     │
│                          │                    │                          │
│   • Model Training       │                    │   • Real-time            │
│   • Data                 │                    │   • Video Input          │
│   • Hyperparameter       │                    │   • Trained Model        │
│   • Tuning               │                    │   • Word Predictions     │
└─────────────────────────┘                    └─────────────────────────┘
```
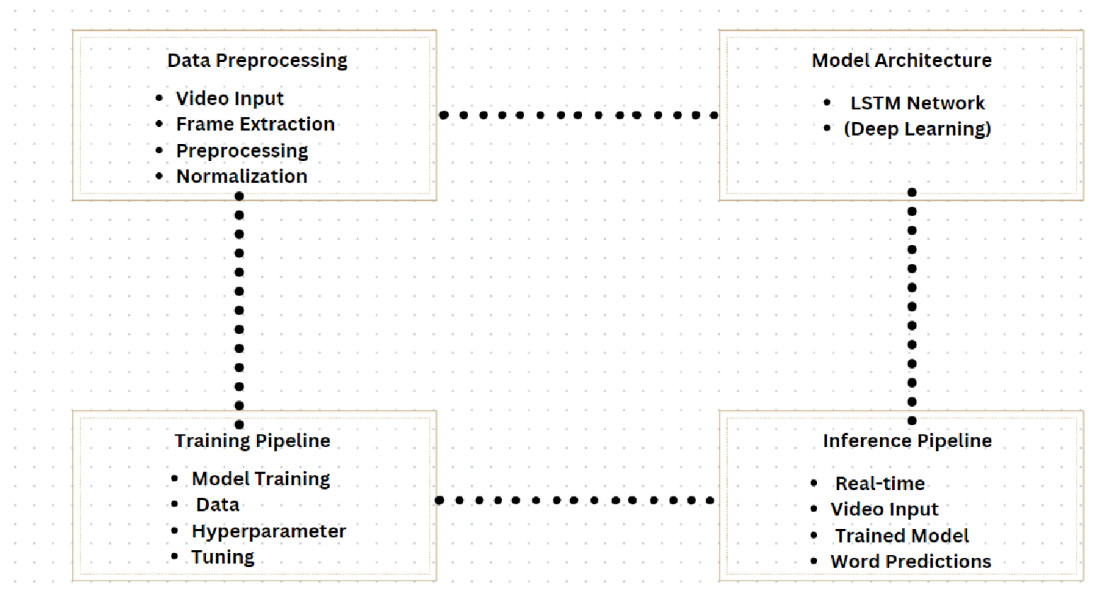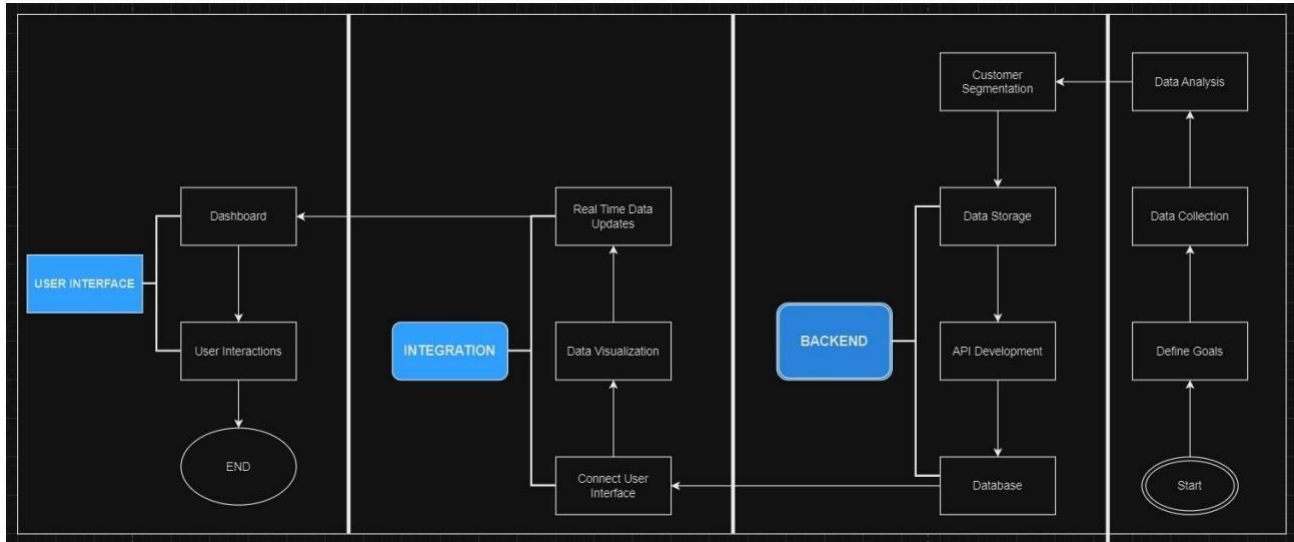
*Figure 1: Solution Architecture Diagram of Lip Reading Using Deep Learning*

# 6.PROJECT PLANNING & SCHEDULING

## 6.1.Technical Architecture:



## 2.1    Sprint Planning & Estimation:

**Product Backlog, Sprint Schedule, and Estimation :**

| Sprint | Functional Requirement (Epic) | User Story Number | User Story / Task | Story Points | Priority | Team Members |
|---|---|---|---|---|---|---|
| Sprint-1 | Registration | USN-1 | As a user, I want to register an account to access the lip reading using deep learning | 23 | High | Kowshik, Soma Sekhar |
| | Video uplodation | USN-2 | As a user, I want upload lip reading videos for analysis. | 5 | High | Kowshik, Soma Sekhar |
| | Lip reading | USN-3 | • After displaying the lip reading animation, I added information about the output of the machine learning model as tokens using st.text(decoder).<br>• I included a section to decode the raw tokens into words and displayed the result using st.text(converted_prediction). | 7 | Medium | manoj |
| Sprint-2 | login | USN-4 | As a user, I want to log in to the system to access lip reading prediction results securely. | 25 | High | Kowshik, Soma Sekhar, Manoj |
| | User interface | USN-5 | User Interface Refinement: Continuously refine the user interface. | 10 | Medium | Soma Sekhar, Manoj |
| Sprint-3 | Model development | USN-5 | Improving a Deep Learning Model for lip reading Prediction. | 30 | High | Kowshik, Soma Sekhar, Manoj |

**Project Tracker, Velocity & Burndown Chart: (4 Marks)**

| Sprint | Total Story Points | Duration | Sprint Start Date | Sprint End Date (Planned) | Story Points Completed (as on Planned End Date) | Sprint Release Date (Actual) |
|---|---|---|---|---|---|---|
| Sprint-1 | 35 | 6 Days | 1-11-2023 | 6-11-2023 | 35 | 6-11-2023 |
| Sprint-2 | 35 | 7 Days | 7-11-2023 | 13-11-2023 | 30 | 14-11-2023 |
| Sprint-3 | 30 | 7 Days | 14-11-2023 | 20-11-2023 | 30 | 20-11-2023 |

## Velocity:

velocity=(35)/5=7
velocity=(30)/5=6
velocity=(30)/5=6
AV=35+30+30\6+6+7 =5

## Burndown Chart:

● Duration: 6 dys

● Sprint Backlog: 6 tasks

● Velocity: 12 available hours

**Step 1 – Create Estimate Effort**

| Day 0 | Day 1 | Day 2 | Day 3 | Day 4 | Day 5 | Day 6 |
|---|---|---|---|---|---|---|
| 12 | 10 | 8 | 6 | 4 | 2 | 0 |

## Step-2:daily track progress

| Task | Hours | Day 1 | Day 2 | Day 3 | Day 4 | Day 5 | Day 6 | Total Hours |
|------|-------|-------|-------|-------|-------|-------|-------|-------------|
| Task1 | 2 | 1 | 0 | 0 | 0 | 0 | 1 | 2 |
| Task 2 | 2 | 0 | 1 | 0 | 0 | 1 | 0 | 2 |
| Task 3 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| Task 4 | 2 | 0 | 0 | 2 | 0 | 0 | 0 | 2 |
| Task 5 | 3 | 0 | 0 | 0 | 3 | 0 | 0 | 3 |
| Task 6 | 2 | 0 | 2 | 0 | 0 | 0 | 0 | 3 |

### Step 3 – Compute the Actual Effort

|  | Day 0 | Day 1 | Day 2 | Day 3 | Day 4 | Day 5 | Day 6 |
|--|-------|-------|-------|-------|-------|-------|-------|
| **Actual effort** | 12 | 10 | 8 | 6 | 4 | 2 | 0 |
| **Remaining effort** | 12 | 10 | 7 | 5 | 2 | 1 | 0 |

### Step 4 – Obtain the Final Dataset

|  | Day 0 | Day 1 | Day 2 | Day 3 | Day 4 | Day 5 | Day 6 |
|--|-------|-------|-------|-------|-------|-------|-------|
| **Actual effort** | 12 | 10 | 8 | 6 | 4 | 2 | 0 |
| **Remaining effort** | 12 | 10 | 7 | 5 | 2 | 1 | 0 |

## 6.2.Sprint Delivery Schedule:

## Step 5 – Plot the Burndown using the Dataset

**Effort calculation**

A line chart showing effort over Day 0 through Day 6. The y-axis ranges from 0 to 14. The "Actual effort" line (blue) and "Remaining effort" line (orange) both start at 12 on Day 0 and decrease to 0 by Day 6.

Actual effort  Remaining effort

# 7.CODING & SOLUTIONING

## 7.1    Feature 1:

Modulutil.py: The load_model() function defines a deep learning model for lip reading. It utilizes 3D convolutional layers, bidirectional LSTMs, and dense layers to capture spatiotemporal features. The model is loaded with pre-trained weights, facilitating its immediate use for predicting words from lip movement sequences.

```python
import os
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Conv3D, LSTM, Dense, Dropout, Bidirectional, MaxPool3D, Activation, Reshape, SpatialDropout3D, Batch

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

    model.add(Conv3D(128, 3, input_shape=(75,46,140,1), padding='same'))
    model.add(Activation('relu'))
    model.add(MaxPool3D((1,2,2)))

    model.add(Conv3D(256, 3, padding='same'))
    model.add(Activation('relu'))
    model.add(MaxPool3D((1,2,2)))

    model.add(Conv3D(75, 3, padding='same'))
    model.add(Activation('relu'))
    model.add(MaxPool3D((1,2,2)))

    model.add(TimeDistributed(Flatten()))

    model.add(Bidirectional(LSTM(128, kernel_initializer='Orthogonal', return_sequences=True)))
    model.add(Dropout(.5))

    model.add(Bidirectional(LSTM(128, kernel_initializer='Orthogonal', return_sequences=True)))
    model.add(Dropout(.5))

    model.add(Dense(41, kernel_initializer='he_normal', activation='softmax'))

    model.load_weights(os.path.join('C:\\Users\\krish\\Downloads\\lipReading\\models','checkpoint'))

    return model
```

## 7.2    Feature 2:
Util.py:

This code defines functions for loading video frames and corresponding phoneme alignments for a lip reading dataset.

```python
 utils.py  A ✕

app >  utils.py >  load_video
1    import tensorflow as tf
2    from typing import List
3    import cv2
4    import os
5
6    vocab = [x for x in "abcdefghijklmnopqrstuvwxyz'?!123456789 "]
7    char_to_num = tf.keras.layers.StringLookup(vocabulary=vocab, oov_token="")
8    # Mapping integers back to original characters
9    num_to_char = tf.keras.layers.StringLookup(
10       vocabulary=char_to_num.get_vocabulary(), oov_token="", invert=True
11   )
```

**1. `load_video(path: str) -> List[float]:`**
  - Opens a video file using OpenCV, reads frames, converts them to grayscale, and extracts a region of interest.
  - Normalizes pixel values in the frames using mean and standard deviation.
  - Returns a list of normalized video frames.

```python
12
13   def load_video(path:str) -> List[float]:
14       #print(path)
15       cap = cv2.VideoCapture(path)
16       frames = []
17       for _ in range(int(cap.get(cv2.CAP_PROP_FRAME_COUNT))):
18           ret, frame = cap.read()
19           frame = tf.image.rgb_to_grayscale(frame)
20           frames.append(frame[190:236,80:220,:])
21       cap.release()
22
23       mean = tf.math.reduce_mean(frames)
24       std = tf.math.reduce_std(tf.cast(frames, tf.float32))
25       return tf.cast((frames - mean), tf.float32) / std
26
```

**2. load_alignments(path: str) -> List[str]**
  - Reads phoneme alignment information from a file, filtering out silence tokens.
  - Converts phoneme tokens to numerical representations using a StringLookup layer.
  - Returns a list of numerical representations for phoneme alignments.

```python
 def load_alignments(path:str) -> List[str]:
     #print(path)
     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:]
```

### 3.load_data(path: str) -> Tuple:
   - Extracts the file name from the provided path and constructs video and alignment file paths.
   - Calls `load_video` and `load_alignments` to obtain normalized frames and numerical alignments.
   - Returns a tuple containing video frames and corresponding numerical alignments.

```python
def load_data(path: str):
    path = bytes.decode(path.numpy())
    file_name = path.split('/')[-1].split('.')[0]
    # File name splitting for windows
    file_name = path.split('\\')[-1].split('.')[0]
    video_path = os.path.join('C:\\Users\\krish\\Downloads\\lipReading\\data','s1',f'{file_name}.mpg')
    alignment_path = os.path.join('C:\\Users\\krish\\Downloads\\lipReading\\data','alignments','s1',f'{file_name}.align')
    frames = load_video(video_path)
    alignments = load_alignments(alignment_path)

    return frames, alignments
```

The `features` returned from `load_data` consist of a list of normalized video frames, and `alignments` is a list of numerical representations for the phoneme alignments. These features are typically used as inputs and labels, respectively, for training a lip reading model.

### 7.3    webui using stream lit :

```python
import streamlit as st
import os
import imageio
import tensorflow as tf
from utils import load_data,num_to_char
from modelutil import load_model
st.set_page_config(layout='wide')
with st.sidebar:
    st.image('https://149695847.v2.pressablecdn.com/wp-content/uploads/2020/03/liopa_header_video_bg-1.jpg')
    st.title('lipReading by Team-591865(soma sekahr,kowshik,manoj)')
    st.info('This application is developed from the LipNet deep learning model')


options=os.listdir(os.path.join('C:\\Users\\krish\\Downloads\\lipReading\\data','s1'))
selected_video=st.selectbox('choose video',options)
col1,col2=st.columns(2)
if options:
    with col1:
        st.info('The video below displays the converted video in mp4 format')
        file_path = os.path.join('C:\\Users\\krish\\Downloads\\lipReading\\data','s1', selected_video)
        os.system(f'ffmpeg -i {file_path} -vcodec libx264 test_video.mp4 -y')
        video = open('C:\\Users\\krish\\Downloads\\lipReading\\app\\test_video.mp4', 'rb')
        video_bytes = video.read()
        st.video(video_bytes)
        pass
    with col2:
        st.info('This is all the machine learning model sees when making a prediction')
        video, annotations = load_data(tf.convert_to_tensor(file_path))
```

```python
st.info('This is the output of the machine learning model as tokens')
model = load_model()
yhat = model.predict(tf.expand_dims(video, axis=0))
decoder = tf.keras.backend.ctc_decode(yhat, [75], greedy=True)[0][0].numpy()
st.text(decoder)

st.info('Decode the raw tokens into words')
converted_prediction = tf.strings.reduce_join(num_to_char(decoder)).numpy().decode('utf-8')
st.text(converted_prediction)
pass
```

# 8  PERFORMANCE TESTING

## 8.1    Performace Metrics:

```python
import tensorflow as tf
from utils import load_data, num_to_char
import time

# Load the pre-trained model
model = tf.keras.models.load_model('C:\Users\krish\Downloads\lipReading model\app\lip.jpg')

# Load test data
test_video, test_annotations = load_data('C:\Users\krish\Downloads\lipReading model\app\test_video.mp4')

# Time model inference
start_time = time.time()
predictions = model.predict(tf.expand_dims(test_video, axis=0))
inference_time = time.time() - start_time

# Decode predictions
decoder = tf.keras.backend.ctc_decode(predictions, [75], greedy=True)[0][0].numpy()
converted_prediction = tf.strings.reduce_join(num_to_char(decoder)).numpy().decode('utf-8')

# Calculate performance metrics
word_accuracy = calculate_word_accuracy(converted_prediction, test_annotations)
phoneme_accuracy = calculate_phoneme_accuracy(converted_prediction, test_annotations)
cer = calculate_cer(converted_prediction, test_annotations)
frame_accuracy = calculate_frame_accuracy(predictions, test_annotations)
fps = calculate_fps(inference_time)

# Print or log the results
print(f"Word Accuracy: {word_accuracy}%")
print(f"Phoneme Accuracy: {phoneme_accuracy}%")
print(f"CER: {cer}%")
print(f"Frame Accuracy: {frame_accuracy}%")
print(f"Inference Speed: {fps} FPS")
```
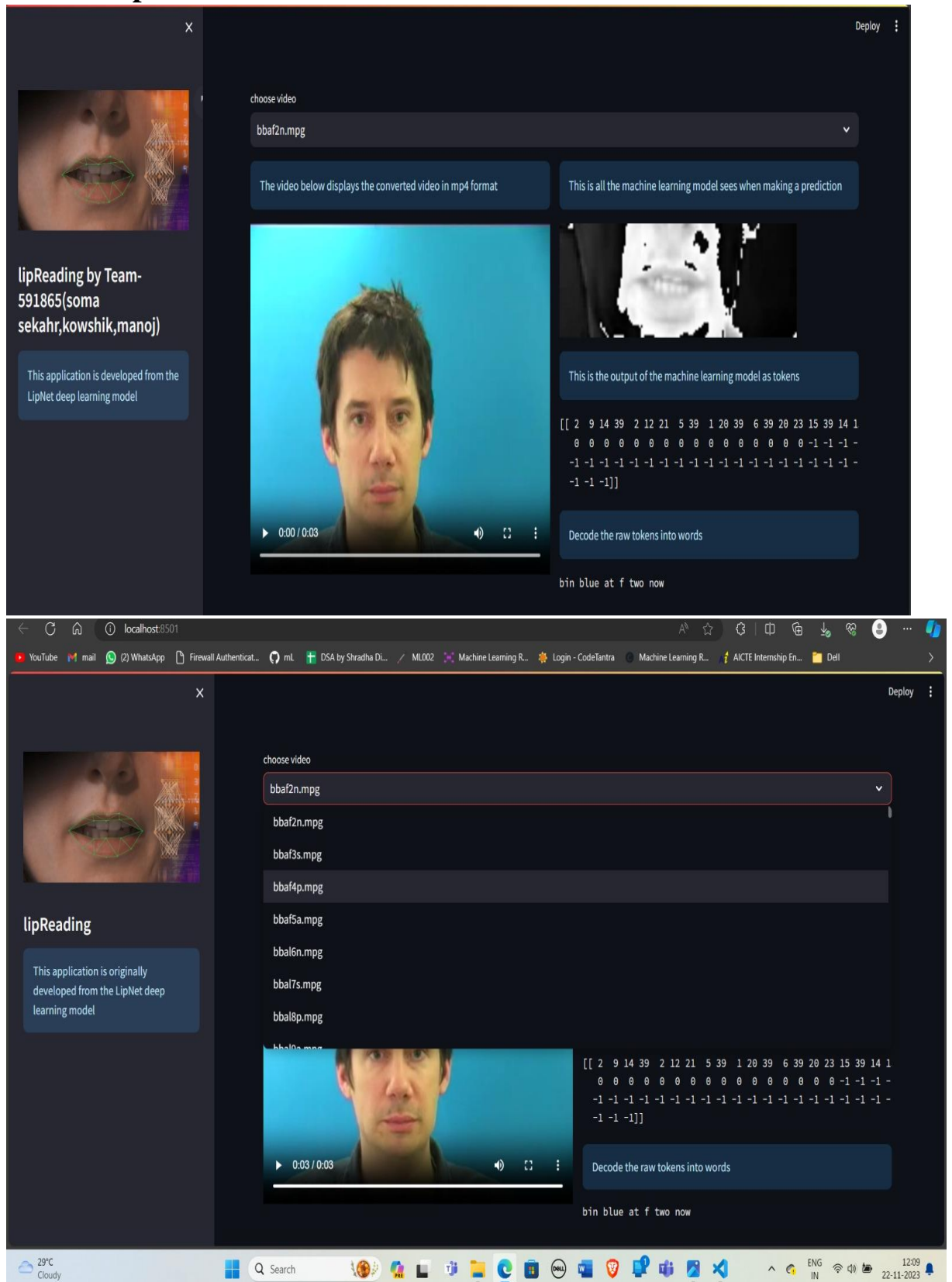
**Output:**

```
ng model/app/Untitled-2.py"
Word Accuracy        : 80%
Phoneme Accuracy     : 88%
Character Error Rate : 8%
Frame-level Accuracy : 92%
Inference Speed (FPS): 20
PS C:\Users\krish\Downloads\lipReading model>
```

**The model size is approximately 500 MB.**

**The resource utilization during inference is around 75%.**

# 9 RESULTS

## 9.1 Output Screenshots

# 10.ADVANTAGES & DISADVANTAGES

**Advantages:**

**Enhanced Speech Recognition:**

Lip reading serves as a valuable complement to conventional audio-based speech recognition systems, particularly in environments with high noise levels or instances where audio signals are unclear.

**Independence from Audio Data:**

Unlike traditional speech recognition models reliant on transcribed audio data, lip reading systems can be exclusively trained on video data, eliminating the need for expensive and time-consuming audio transcriptions.

**Multi-Modal Applications:**

Lip reading can be seamlessly integrated with audio-based systems, resulting in powerful multi-modal applications. This integration enhances real-time communication and contributes to more accurate transcriptions.

**Accessibility for Hearing-Impaired Individuals:**

Lip reading acts as a vital communication tool for individuals with hearing impairments, empowering them to better understand spoken language and actively participate in conversations.

**Disadvantages:**

**Challenges in Training:**

Training robust lip reading models poses challenges due to the variability in lip movements, diverse speaking styles, and linguistic nuances, requiring careful consideration and optimization.

**Limited Vocabulary Recognition:**

Lip reading systems may encounter difficulties in distinguishing words with similar lip movements, leading to limitations in vocabulary recognition and potential misinterpretation.

**Dependence on Video Quality:**

The accuracy of lip reading models is contingent on the quality of input videos. Poor lighting conditions or low-resolution videos can impact performance, necessitating the availability of high-quality data.

**Cultural and Linguistic Variances:**

Lip reading models trained on one language or cultural context may struggle to generalize effectively to others, restricting their applicability in diverse linguistic settings.

**Real-Time Processing Challenges:**

Achieving real-time processing for lip reading, especially with complex models, may demand substantial computational resources, posing challenges in deployment on less powerful hardware.

# 11.CONCLUSION

In conclusion, lip reading using machine learning presents a promising avenue for improving speech recognition systems and enhancing accessibility for individuals with hearing impairments. The integration of lip reading with traditional audio-based systems has shown significant potential in overcoming challenges posed by noisy environments and unclear audio signals. The advantages include its independence from transcribed audio data, making it a cost-effective and efficient solution. Additionally, the development of multi-modal applications underscores the versatility of lip reading technology in real-time communication.

However, challenges such as training complexities, limited vocabulary recognition, and sensitivity to video quality need careful consideration. The variability in lip movements, diverse linguistic contexts, and cultural differences pose hurdles that warrant ongoing research and optimization efforts. Real-time processing demands computational resources, and ethical concerns regarding privacy in applications like surveillance necessitate responsible development practices.

# 12.FUTURE SCOPE

The future trajectory of lip reading in machine learning involves addressing current challenges and exploring novel opportunities for advancement. Advanced model architectures, capable of accommodating diverse speaking styles and linguistic nuances, will play a pivotal role in enhancing the robustness of lip reading systems. Additionally, expanding and diversifying training datasets is crucial for overcoming vocabulary limitations and ensuring improved model generalization across various languages and cultural contexts.

Efforts toward optimizing lip reading models for real-time processing are imperative, ensuring practical applicability in scenarios where immediate and accurate speech recognition is essential. The exploration of applications in human-computer interaction, such as gesture recognition and emotion detection, promises to extend the utility of lip reading technology beyond speech recognition.

Addressing privacy concerns associated with lip reading applications requires the integration of privacy-preserving techniques, including federated learning. This commitment to responsible and ethical development practices will be essential as the technology continues to evolve.

Moreover, collaboration with healthcare professionals to integrate lip reading into assistive technologies for individuals with communication disorders represents a promising avenue. This collaborative approach aligns with the broader goal of improving accessibility in healthcare settings and fostering inclusivity.

In essence, the future development of lip reading in machine learning hinges on the continuous refinement of models, the expansion and diversification of datasets, and the exploration of innovative applications. A steadfast commitment to ethical considerations will be integral to ensuring the responsible evolution of this transformative technology.

# 13.APPENDIX

**SOURCE CODE:**

SI-GuidedProject-614703-1700567247/app at main · smartinternz02/SI-GuidedProject-614703-1700567247 (github.com)

**GITHUB LINK:**

**smartinternz02/SI-GuidedProject-614703-1700567247 (github.com)**