# THE LIP TALKER

By

Reeve R. Mathew (2348573)

Satyam Jhawar (2348554)

Christina J. Thattil (2348511)

Under the able guidance of

Dr. Jobin Francis

Specialization Project Report Submitted in Partial Fulfilment

of the Requirements of IV<sup>th</sup> Trimester MSAIM,
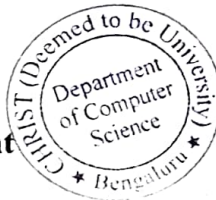
CHRIST (Deemed to be University)

August 2024

# CHRIST
(DEEMED TO BE UNIVERSITY)
BANGALORE · INDIA

## CERTIFICATE

This is to certify that the report titled **THE LIP TALKER** is a bona fide record of work done by Satyam Jhawar(2348554), Reeve R. Mathew(2348573) and Christina J. Thattil(2348511) of CHRIST(Deemed to be University), Bangalore, in partial fulfillment of the requirements of IVth Trimester MSAIM during the year 2024.

**Head of the Department**

**Project Guide**

Valued-by:

1.

2.

Name : Satyam Jhawar
Register Number : 2348554
Examination Centre : CHRIST (Deemed to be University)

Date of Exam :Aug 2024

# CHRIST

(DEEMED TO BE UNIVERSITY)

BANGALORE · INDIA

## CERTIFICATE

This is to certify that the report titled **THE LIP TALKER** is a bona fide record of work done by Reeve R. Mathew(2348573), Christina J. Thattil(2348511) and Satyam Jhawar(2348554) of CHRIST(Deemed to be University), Bangalore, in partial fulfillment of the requirements of IVth Trimester MSAIM during the year 2024.

**Head of the Department**

**Project Guide**

Valued-by:

1.

2.

Name : Reeve R. Mathew
Register Number : 2348573
Examination Centre : CHRIST (Deemed to be University)
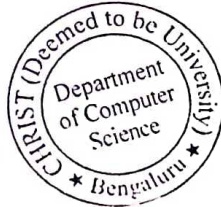
Date of Exam :Aug 2024

# CHRIST
## (DEEMED TO BE UNIVERSITY)
### BANGALORE · INDIA

# CERTIFICATE

This is to certify that the report titled **THE LIP TALKER** is a bona fide record of work done by Christina J. Thattil(2348511), Reeve R. Mathew(2348573) and Satyam Jhawar(2348554) of CHRIST(Deemed to be University), Bangalore, in partial fulfillment of the requirements of IVth Trimester MSAIM during the year 2024.

**Head of the Department**

**Project Guide**

Valued-by:

1.

2.

Name : Christina J. Thattil
Register Number : 2348511
Examination Centre : CHRIST (Deemed to be University)

Date of Exam :Aug 2024

# ACKNOWLEDGMENT

# ABSTRACT

Speech recognition has witnessed remarkable advancements in recent years, with deep learning techniques playing a pivotal role. However, most existing systems rely solely on audio data, limiting their applicability in scenarios where visual information is also available. This project aims to address this limitation by exploring the potential of video-based speech recognition. By leveraging both visual and auditory cues, we aim to develop a more robust and accurate system that can be applied to a wider range of applications. This project proposes a deep learning model that effectively transcribes spoken text from video data. Our model integrates 3D convolutional layers and bidirectional (Long Short-Term Memory) LSTMs to capture both spatial and temporal features from video frames. To handle the inherent challenges of sequence-to-sequence learning, we employ Connectionist Temporal Classification (CTC) loss, which eliminates the need for explicit frame-to-label alignment.

# TABLE OF CONTENTS

# 1.INTRODUCTION

Speech recognition has traditionally relied on audio signals, but advancements in multimedia processing have opened new avenues for integrating video data into speech recognition systems. Videos contain rich visual information that, when combined with audio, can enhance the accuracy and robustness of speech recognition systems. This project explores the integration of visual cues from video frames to transcribe spoken words, leveraging the power of deep learning techniques.

Lip reading is a technique that relies on context, language knowledge, and any residual hearing, so estimates of its range vary. While it's most commonly used by deaf and hard-of-hearing people, most people with normal hearing also process some speech information from seeing a moving mouth.

In a technological context, VSR is the process of predicting a patient's speech transcript from silent videos of their mouth movements. Deep lip reading, a type of VSR, uses deep neural networks to extract speech from a video of a silent talking face. The process involves two stages:

- Extracting visual and temporal features from a sequence of image frames from the video

- Processing the sequence of features into units of speech, such as words, characters, or phrases.

The core of our approach involves designing a deep neural network capable of processing video data. The proposed model utilizes 3D convolutional layers to extract spatial and temporal features from video frames, followed by bidirectional Long Short-Term Memory (LSTM) layers to model the sequential nature of speech. The network is trained using Connectionist Temporal Classification (CTC) loss, which accommodates the varying lengths of audio segments and allows the model to learn from sequences without requiring frame-by-frame alignment.

To achieve this, we first prepare a dataset of videos and corresponding text alignments, processing the videos into frames and normalizing them. We then build a data pipeline to handle these frames and text sequences efficiently, ensuring that the model receives properly formatted input. Our model is trained on this dataset and evaluated on a separate test set to assess its transcription accuracy.

Through this project, we aim to demonstrate the feasibility of video-based speech recognition and provide insights into the effectiveness of combining visual and auditory information for transcription tasks. The results highlight the potential for enhancing traditional speech recognition systems with additional modalities and pave the way for future research in multimedia speech processing.

# 1.1 Alignment with SDG Goals

Lip Talker project aligns with several United Nations Sustainable Development Goals (SDGs), particularly: Lip Talker project aligns with the below United Nations Sustainable Development Goals (SDGs):

**SDG 3:** Sustainable Development Goal 3: Good Health and Well-being – due to barriers in communication, the project empowers the deaf and hard-of-hearing individuals to be mentally sound and socially included.

**SDG 4:** Perceived Quality - The availability of the application enhances the student's access to learning resources and communication, hence enhancing universal education.

**SDG 10:** Peace and Justice – Lip Talker also impacts reduced inequality as it helps to provide people with hearing impairment opportunities to be accepted in society.

**SDG 16:** The finally identified category matches the title of The Global Goals for Sustainable Development: Peace, Justice and Strong Institutions. The application used in defence scenes is assigned to increase security and control to avoid insecurity, hence the aim and objective of the goals of strong institutions.

The project is designed to assist in achieving various important goals of the SDGs: Enhancing health and providing benefits, ensuring education for everyone, decreasing inequalities, and ultimately enhancing security, which is why technology is needed to bring about positive change and improve the well-being of all people on Earth.

# 2. LITERATURE REVIEW

The increasing complexity of modern life has led to a surge in waste generation, making efficient recycling crucial for environmental sustainability. However, the vital first step of waste sorting often proves to be a tedious and labor-intensive process. In a study by Durga Sri et al., the authors tackle a different kind of challenge related to interpreting visual information: lip reading. This method of understanding speech, relying solely on the speaker's lip movements, faces difficulties due to variations in speaking styles and visual ambiguities. The authors explore the potential of deep learning, specifically Convolutional Neural Networks, to automate and improve lip reading accuracy. By training and evaluating two distinct CNN architectures, the research demonstrates the power of these models in accurately predicting words from visual input. This paves the way for practical applications like real-time word prediction, showcasing the potential of deep learning in deciphering complex visual information Durga.

Chung et al. presented their research in the paper titled "Lip Reading Sentences in the Wild," where they addressed the complexities of deciphering spoken words from visual information in unconstrained environments. They recognized the limitations of prior methods that depended on restricted vocabularies and controlled conditions, and proposed the 'Watch, Listen, Attend and Spell' (WLAS) network, which effectively maps lip movement features to corresponding characters using deep learning. To overcome the challenge of limited datasets, the authors introduced the 'Lip Reading Sentences' dataset, consisting of over 100,000 natural sentences sourced from British television, which ensures the necessary diversity for training robust models. The WLAS model, trained on this dataset, achieves exceptional accuracy, surpassing earlier benchmarks and even outdoing a professional lip reader when tested on BBC television footage. This work underscores the transformative impact of deep learning and large-scale datasets in enhancing lip reading technology, with promising applications across various fields, from assistive technologies to human-computer interaction.

Assael et al. introduced LipNet, a groundbreaking end-to-end model that significantly advances the field of lipreading by predicting entire sentences from video sequences of lip movements. Unlike traditional methods that separate feature extraction and prediction, LipNet employs a blend of spatiotemporal convolutions and recurrent networks to effectively capture both spatial and temporal features. This innovative architecture results in a remarkable accuracy of 95.2% on the GRID corpus, surpassing previous models that achieved a maximum accuracy of 86.4% focused only on word-level classification. By leveraging a comprehensive dataset featuring 34 speakers producing 1000 sentences each, LipNet demonstrates robust generalization capabilities with an accuracy of 88.6% in unseen speaker scenarios, making it particularly valuable for applications in noisy environments where audio cues may be unreliable. This advancement not only streamlines the lipreading process but also enhances its accuracy, providing significant implications for improved communication technologies and accessibility solutions.

Attention-based models have emerged as a powerful approach in Automatic Speech Recognition (ASR), offering significant improvements over traditional methods. These models, inspired by the human ability to focus on selective parts of sensory input, leverage an attention mechanism to selectively attend to relevant parts of the input speech signal while generating the output transcription. This differs from earlier approaches like Hidden Markov Models and basic Recurrent Neural Networks, which process input sequentially and can struggle with long sequences. A seminal work in this domain, "Attention-Based Models for Speech Recognition" by Chorowski et al., introduced key innovations like location-awareness to address the limitations of earlier attention models in handling variable-length utterances. This paved the way for further advancements in attention-based ASR, including hybrid models combining content and location-based attention, leading to improved performance and robustness. Subsequent research has explored various attention mechanisms,

including content-based, location-based, and hybrid approaches, demonstrating their efficacy in various ASR tasks. The development of attention-based models has significantly advanced the state-of-the-art in ASR, enabling more accurate and efficient speech recognition systems.

In their pioneering work, Chorowski et al. (2021) present a paper titled "Hybrid CTC/Attention Model Based on Conformers," which makes significant strides in speech recognition technology by integrating attention-based models. The authors aim to enhance the accuracy and robustness of speech recognition systems, particularly in complex acoustic environments, by leveraging essential attention mechanisms in modern machine learning. Recognizing the pressing demand for improved speech recognition capabilities across real-world applications, this study not only seeks to enhance the reliability of these systems but also explores innovative methodologies to advance their performance. By harnessing the power of attention mechanisms, this research has the potential to redefine the landscape of speech recognition systems, paving the way for more reliable and efficient technology. Ultimately, this seminal work positions itself as a vital contribution to the evolution of speech recognition systems, aiming for profound advancements in the field.

In their research, Koumparoulis and Potamianos (2022) present the paper titled "Accurate and resource efficient lipreading with EfficientNetV2 and transformers." As advancements in visual speech recognition (VSR) gain momentum, the need for effective lipreading systems becomes increasingly paramount in applications such as accessibility and communication technologies. The authors propose a novel approach that leverages EfficientNetV2 and transformer architectures to enhance the accuracy and efficiency of lipreading models. By incorporating a range of training strategies and optimizing the neural network architecture, they achieve a significant improvement in performance metrics, including a top accuracy of 93.2% in distinguishing lip movements associated with spoken words. This work not only pushes the boundaries of VSR capabilities but also emphasizes the importance of leveraging cutting-edge deep learning techniques for practical applications in real-world scenarios .

The paper titled "END-TO-END AUDIO-VISUAL SPEECH RECOGNITION WITH CONFORMERS" by Ma, Petridis, and Pantic (2021) addresses the growing need for robust speech recognition systems in noisy environments, highlighting the limitations of traditional two-step approaches that separate feature extraction and recognition. Recent advancements in the field have shifted towards end-to-end (E2E) models, which integrate these processes within deep neural networks, significantly enhancing performance in both VSR and Automatic Speech Recognition (ASR). Noteworthy contributions include Assael et al.'s work on 3D convolutional networks and Shillingford et al.'s Vision to Phoneme model, which predicts phoneme distributions from video clips. The authors emphasize the effectiveness of using a hybrid CTC/attention mechanism and conformers, which replace recurrent networks, leading to substantial improvements in recognition accuracy. This paper builds on previous research by directly processing raw audio and visual data, demonstrating a clear trend in leveraging deep learning techniques to improve the robustness and accuracy of AVSR systems, particularly in challenging conditions.

# 3. ANALYSIS OF EXISTING SYSTEMS

- Present lip-reading systems and applications are mainly aimed at the people with hearing disabilities to transcribe visual oral movements into texts. While these systems have made significant strides, they come with several limitations:

- Accuracy and Reliability: Most of the current lip-reading systems have challenges in their efficiency since they are not accurate especially in environments full of noise or when the lighting and camera position differs. Said misunderstandings would lead to wrong translations, which in turn would reduce the quality of communication.

- Real-Time Processing: There is a critical concern of realizing real-time processing. This comes in the way of the rhythm of the conversation and makes it relatively more difficult and less fluent to converse because of un-timely translation of lip movements to text.

- User Interface: The present interfaces in those systems are complicated, and more often not so friendly or do not support easy accessibility to all users regardless of their IT skills.

# 4.   PROPOSED SYSTEM

## 4.1 System Description

Lip Talker is an application that will be hosted on the web for the use by the deaf and hard of hearing persons as well as for improving the defence applications through lip movement recognition technology. The created system based on deep learning and Google Translate synchronizes with a video in real time, reproduces its lip movements to type English text, which is further translated to the desired languages. This approach lays emphasis on making them equally communicative and easy to use regardless of the users' linguistic differences.

## 4.2 Programming Languages used:

- Python
    - o   Deep Learning
    - o   Computer Vision
    - o   Keras
    - o   Tensorflow
- HTML5
- CSS
- JavaScript
- Google Drive Storage

## 4.3 Functionality

- **Video Upload:** The web application allows the users to upload the videos of lip movements.

- **Lip Reading:** The uploaded video streams it to recognize and analyse lip movements and translate them into English text.

- **Text Translation:** An English text is translated into multiple languages via using a API

- **Output Display:** The translated text becomes visible on the web application while the text is spoken out loud to enhance the usage.
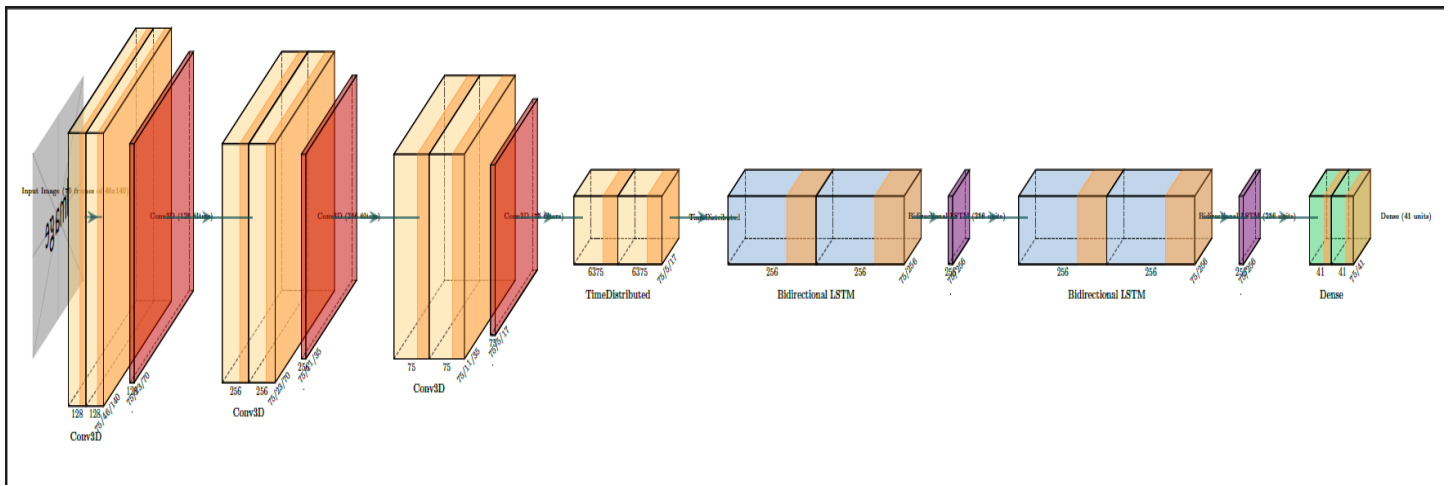
Figure 1: Lip Talker Model Architecture



```
Model: "sequential"
_____
 Layer (type)                 Output Shape              Param #
=================================================================
 conv3d (Conv3D)              (None, 75, 46, 140, 128)  3584

 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)   884992

 activation_1 (Activation)    (None, 75, 23, 70, 256)   0

 max_pooling3d_1 (MaxPooling  (None, 75, 11, 35, 256)   0
 3D)

 conv3d_2 (Conv3D)            (None, 75, 11, 35, 75)    518475

 activation_2 (Activation)    (None, 75, 11, 35, 75)    0

 max_pooling3d_2 (MaxPooling  (None, 75, 5, 17, 75)     0
 3D)

 time_distributed (TimeDistr  (None, 75, 6375)          0
 ibuted)

 bidirectional (Bidirectiona  (None, 75, 256)           6660096
 l)

 dropout (Dropout)            (None, 75, 256)           0

 bidirectional_1 (Bidirectio  (None, 75, 256)           394240
 nal)

 dropout_1 (Dropout)          (None, 75, 256)           0

 dense (Dense)                (None, 75, 41)            10537

=================================================================
Total params: 8,471,924
Trainable params: 8,471,924
Non-trainable params: 0
_____
```

Figure 2: Lip Talker Neural Network model architecture

- Conv3D Layers:

  o Function: These layers perform 3D convolutions by applying 3D filters across the spatial dimensions (height and width) and the temporal dimension (depth or sequence of frames) of the input video.

  o Role: The first Conv3D layer extracts basic spatial features like edges and textures from the video frames. Subsequent Conv3D layers detect more complex patterns and motions over time, such as the movement of lips or facial expressions in a sequence of frames.

7

- Activation Layers (ReLU):

  o Function: These layers apply the Rectified Linear Unit (ReLU) activation function, which introduces non-linearity by converting all negative pixel values to zero while leaving positive values unchanged.

  o Role: Non-linearity allows the network to learn complex patterns by combining simple features into more sophisticated representations. This helps in identifying intricate details, like subtle lip movements or changes in facial expressions.

- MaxPooling3D Layers:

  o Function: These layers perform down sampling by selecting the maximum value from each region of the input, reducing the spatial and temporal resolution.

  o Role: MaxPooling3D layers help in retaining the most important features (like the most intense or significant regions) while reducing the size of the data. This decreases computational load and controls overfitting by discarding less critical information.

- Time Distributed Layer:

  o Function: This layer applies a layer (typically a Dense layer) independently to each time step of the input sequence.

  o Role: The Time Distributed layer ensures that each frame or feature set within the sequence is processed individually, preserving the temporal order. This is crucial for tasks like lip-reading, where the sequence of frames matters for understanding the motion and context.

- Bidirectional LSTM Layers:

  o Function: These layers consist of Long Short-Term Memory (LSTM) units that process the sequence data in both forward and backward directions.

  o Role: The Bidirectional LSTM layers capture dependencies in the video frames by considering both past (previous frames) and future (upcoming frames) information. This enhances the model's ability to understand context and timing, essential for tasks like recognizing words based on lip movements.

- Dropout Layers:

  o Function: These layers randomly deactivate a portion of the neurons during training, effectively "dropping out" a fraction of input units to prevent overfitting.

  o Role: Dropout layers improve the model's generalization by ensuring it doesn't become overly dependent on any single neuron or feature. This helps the model perform better on unseen data by making it more robust.

- Dense Layer:

  o Function: This fully connected layer aggregates all the features extracted and processed by the previous layers and combines them into a final output.

Role: The Dense layer is responsible for making the final prediction, such as identifying a word or phrase based on the processed video data. It translates the learned features into a decision, like classifying a specific spoken word in a lip-reading task.

# 4.4 Proposed Solution Workflow:

- **Backend Workflow:**

1. Initialization and Setup:

- Import necessary libraries for the project, including TensorFlow, Keras, OpenCV, gdown, imageio, and matplotlib.

2. Define Vocabulary and Model Architecture:

- Define the vocabulary for the text recognition task, including alphanumeric characters and some symbols.

- Construct the neural network model architecture for video-based text recognition:

- Build Conv3D layers for spatial and temporal feature extraction.

- Incorporate Time Distributed layer to process spatiotemporal features.

- Utilize Bidirectional LSTM layers for sequence modelling.

- Add Dense layer for character prediction using SoftMax activation.

3. Load Pre-trained Weights:

- Utilize gdown library to download pre-trained model weights from a specified URL.

- Extract the downloaded checkpoint zip file and load the weights into the model.

4. Testing Model:

- Load test data, possibly in the form of video frames, as a numpy iterator.

- Predict text sequences using the trained model on the test data.

- Decode the predictions using CTC decoding to convert numerical outputs back to text characters.

5. Inference on Sample Data:

- Load a sample video file or data for text recognition.

- Prepare the sample data for inference and prediction.

- Make predictions using the model on the sample data.

- Decode the model predictions using CTC decoding to obtain the recognized text.

6. Printing Results:

- Display the real text from the test data alongside the predicted text after model inference.

7. End of Execution:

- The program execution concludes after displaying the actual and predicted text output for both the test data and the sample input file.

These points present a concise outline of the operation flow within the provided code snippet, highlighting the key steps involved in training, testing, and using the video-based text recognition model for making predictions on both test data and a specific sample input.
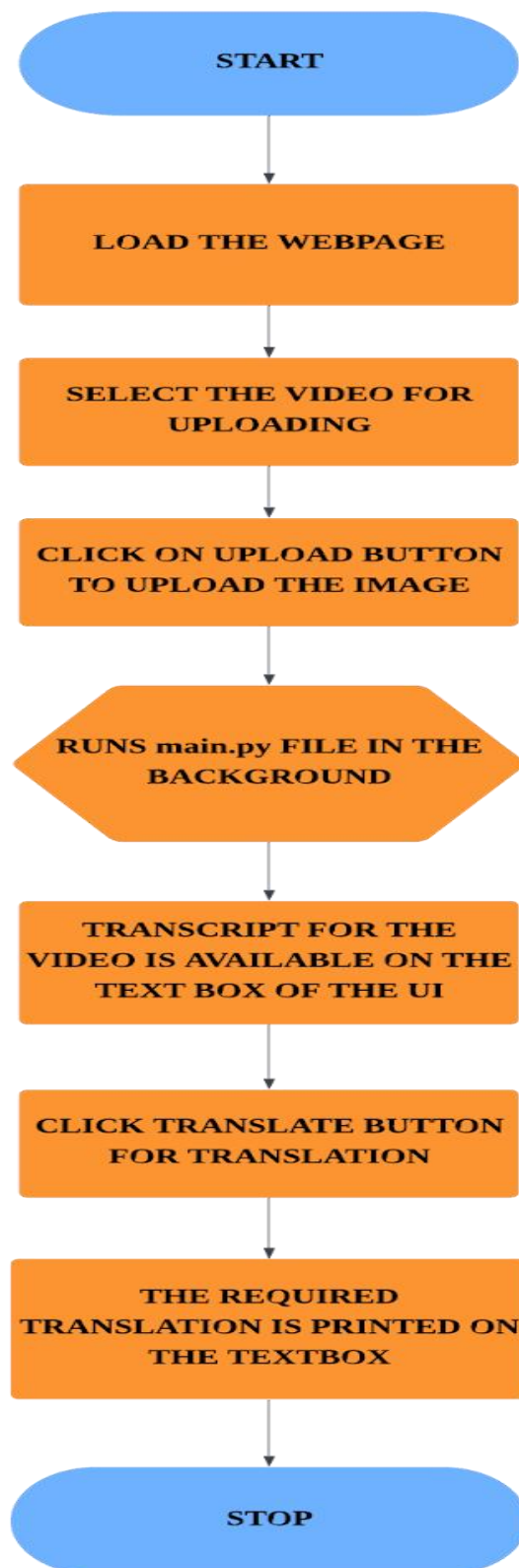
*Figure 3: Backend workflow*

- **Frontend workflow:**

Step 1: Upload Video on Web App

A video of lip movements is recorded by a user through a web-based interface. Video Gets Uploaded on Google Drive: The video file is then saved directly to the subject's Google Drive for further security of the content and convenience for future use.

Step 2: Python Script processes video

The script in Python is located on the server and takes the video file from Google Drive. To put it into words, deep learning algorithms are employed to decode the lip's movement, coordinating with the script.

Step 3: Output Displayed on Web App

The text generated by the Python script is returned to the web application and then displayed on the output screen for user review.

Step 4: Text Translation Begins

This is then translated into the chosen target languages through Google Translate from the original English text.

Step 5: Output Generated

The result is a translation of the text that appears at the end of the web application and the audio function to read the text in the target language in loudspeakers, depending on the user's need.

Through the architecture followed above, Lip Talker makes lip-reading translation and communication easy, making the product more usable across all clients with possible hearing difficulties.
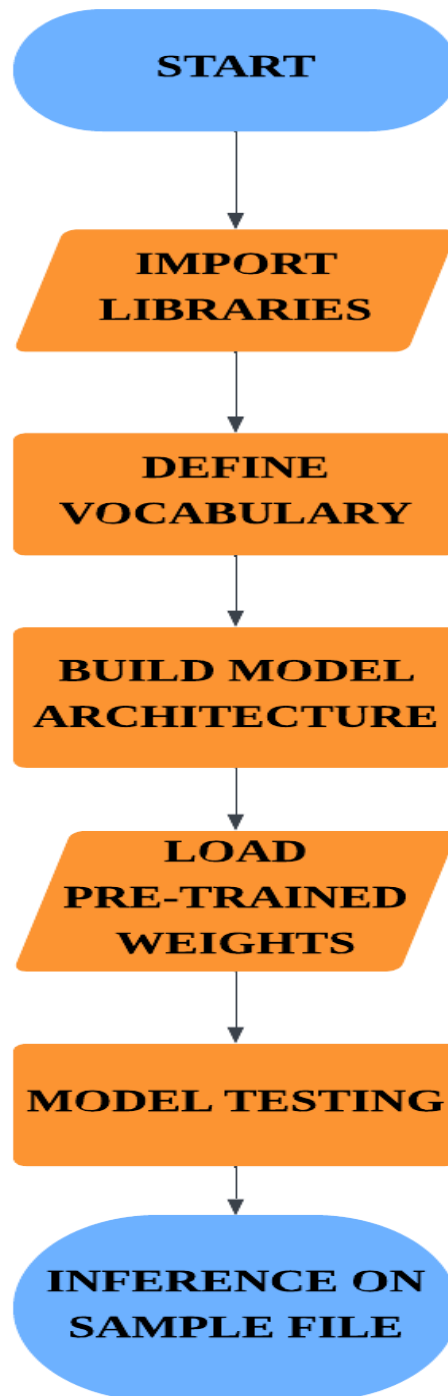
*Figure 4: Frontend-user interface workflow*

# 5. RESULTS AND DISCUSSIONS

This section details the workflow of converting a coloured video into a black-and-white GIF, processing it through a neural network, and utilizing various deep learning techniques for accurate speech transcription. The process begins with transforming a coloured video into grayscale and compiling it into a 10 FPS GIF. This GIF is then used as input for a neural network designed for lip-reading. The network architecture features Convolutional 3D (Conv3D) layers for extracting spatial and temporal features, followed by Rectified Linear Unit (ReLU) activations, MaxPooling3D layers for down-sampling, and Bidirectional LSTM layers for capturing temporal dependencies. The combination of these components enables the network to interpret complex patterns such as lip movements. The final output is a transcription of the spoken content, validated against actual text for accuracy. Additionally, this section covers the user interface for uploading videos, processing them, and translating the transcriptions into different languages, showcasing the practical application of the system.
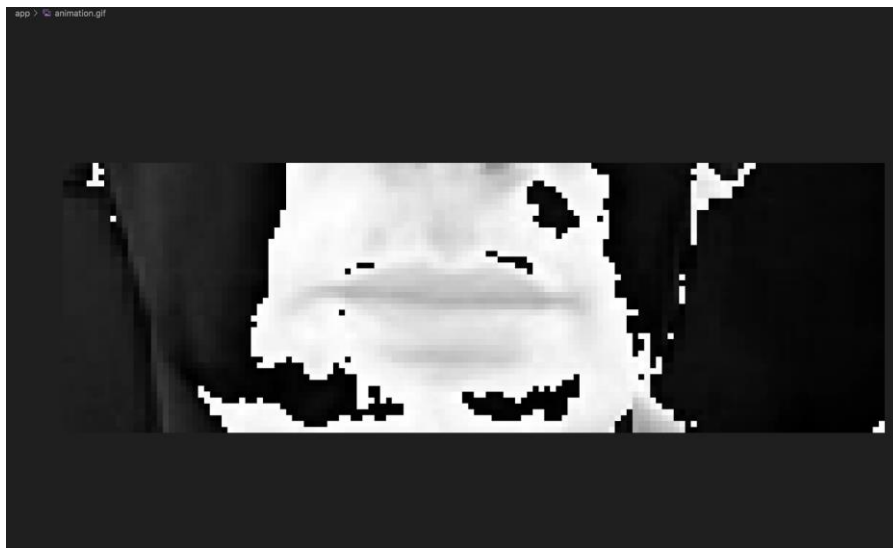


*Figure 5: Converts a coloured video to black and white and saves it as a 10 FPS GIF*

- **Video Conversion to Black and White:** The coloured video is processed to remove colour information, resulting in a black-and-white (grayscale) version. This is typically done to focus on intensity variations without the distraction of colour.

- **GIF Creation:**

    o **Input Data:** The grayscale frames of the video are collected into a sequence. In your code, val[0][0] represents this sequence of frames.

    o **GIF Generation:** The imageio.mimsave(...) function takes this sequence of frames and compiles them into a GIF.

- **GIF Specifications:**

    o **Frame Rate:** The fps=10 parameter sets the GIF to display at 10 frames per second, controlling the speed at which the frames are shown.

*Figure 6: This GIF is given as input as frames to the neural network*

```
1  # 0:videos, 0: 1st video out of the batch,  0: return the first frame in the video
2  plt.imshow(val[0][0][35])
```

```
<matplotlib.image.AxesImage at 0x10c2ead8d30>
```
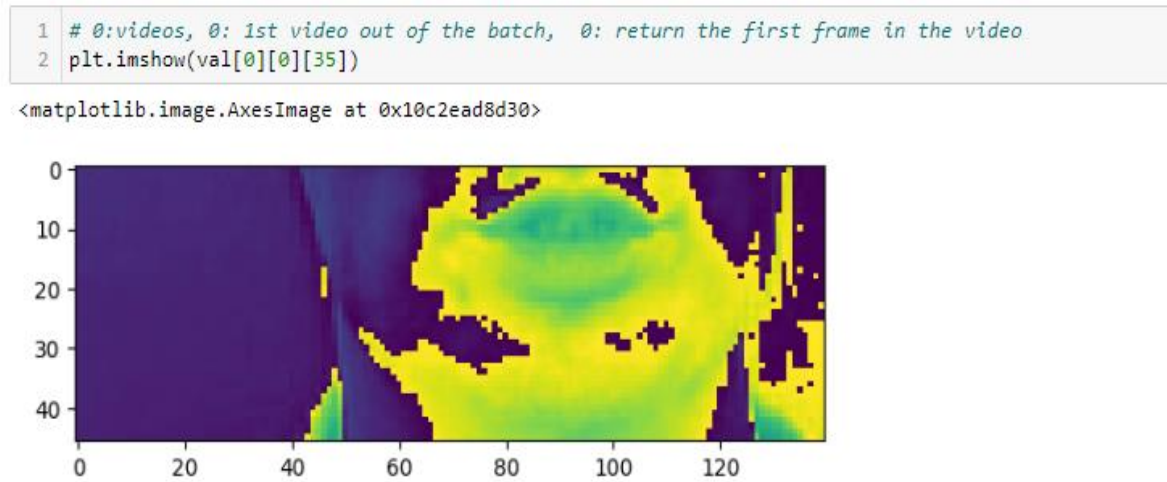


*Figure 7: Image warmth of a single frame*

The image is a heatmap-like visualization of a specific video frame. This type of image representation is commonly used in computer vision tasks where the intensity or value of each pixel is mapped to different colours.

In this particular image:

- Colour Mapping: The colours represent different pixel intensity values. Typically, warmer colours (yellows and greens) indicate higher intensity or brightness, while cooler colours (purples and blues) indicate lower intensity.

- Image Content: The shape within the frame seems to resemble a human face, possibly highlighting areas of higher intensity around features like the eyes, nose, and mouth. The exact details might be abstracted or obscured, depending on the preprocessing or encoding applied to the video frame.

The output is a visual representation that helps to understand the distribution of pixel values within that specific video frame, which can be useful for analysing the video content, especially in tasks like facial recognition, emotion detection, or lip-reading.

# 5.1 Utilization of Image Warmth by the Neural Networks

- **Convolutional Layers (Conv3D)**:

  - The Conv3D layers analyse the spatial and temporal features of the input video frames, where image warmth (pixel intensity) plays a key role.

  - Each convolution operation involves applying a set of filters (kernels) to the input data. These filters detect patterns such as edges, textures, and more complex features as the network goes deeper.

  - Higher intensity areas (warmer regions) may activate certain filters more strongly, leading the network to pay more attention to these areas. For example, in lip-reading, the mouth region might have higher intensity due to movement, and the network would focus on these movements to understand speech.

- **Activation Functions**:

  - After convolution, the activation functions (like ReLU) introduce non-linearity, which helps the network learn complex patterns. The activation function will amplify or suppress the influence of certain regions based on their intensity.

  - For instance, if the warmer regions (higher intensity) correspond to critical features like lip movements, the activation function ensures that these features are preserved and highlighted as the data passes through the network.

- **Pooling Layers**:

  - The MaxPooling3D layers reduce the spatial dimensions of the data while retaining the most prominent features. In the context of image warmth, these layers ensure that the most intense (and likely important) areas of the image are preserved, even as the data is down sampled.

  - For example, if a high-intensity area corresponds to a moving lip, the pooling layer will retain this crucial information while discarding less important details.

- **Time Distributed and Bidirectional Layers**:

  - The Time Distributed layer ensures that the spatial features (like image warmth) extracted by the convolutional layers are passed on to the RNN layers in a structured manner, maintaining the temporal sequence of frames.

  - The Bidirectional LSTM layers analyse the temporal sequence of frames in both forward and backward directions. The image warmth influences how these layers capture and interpret motion or changes in intensity over time, which is vital for understanding dynamic processes like speech in lip-reading.

- **Output Layer (Dense)**:

  - Finally, the dense layer combines the learned features (influenced by image warmth and other factors) to make a final prediction. In a lip-reading task, this might be predicting the spoken word or phrase based on the visual input.

# 5.2 Working of CNN and LSTM

1. **CNN for Spatial Features**:

   o **Function**: CNNs extract spatial features from images or video frames, like edges, textures, and shapes.

   o **Benefit**: They efficiently process visual data, reducing the input size while preserving important patterns.

2. **LSTM for Temporal Sequencing**:

   o **Function**: LSTMs analyze the sequence of frames over time, capturing dependencies between them.

   o **Benefit**: They understand the temporal context, essential for tasks like lip-reading where the order of frames matters.

3. **Integration**:

   o **Process**: The CNN first extracts features from each frame, which are then passed to the LSTM to analyze the sequence.

   o **Outcome**: This combination allows the model to make informed predictions based on both spatial and temporal information.

4. **Training and Optimization**:

   o **Loss Function**: Measures prediction errors; the model updates its parameters to minimize this loss.

   o **Optimization**: Techniques like backpropagation adjust both CNN and LSTM layers, improving overall accuracy.

CNNs handle the "what" by extracting spatial features, and LSTMs handle the "when" by analyzing how these features change over time. Together, they provide a powerful approach for tasks involving both spatial and temporal data, leading to better accuracy.

The real text value represents the actual spoken content of the video, serving as a reference for accuracy. The predicted text is generated by the lip-reading model based on the visual analysis of lip movements. By comparing the predicted text with the real text, we evaluate the model's performance and accuracy in transcribing speech.

## Test on a Video

```
In [61]:   1  sample = load_data(tf.convert_to_tensor('.\\data\\s1\\bras9a.mpg'))
```

```
In [62]:   1  print('~'*100, 'REAL TEXT')
           2  [tf.strings.reduce_join([num_to_char(word) for word in sentence]) for sentence in [sample[1]]]
```

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~ REAL TEXT

```
Out[62]: [<tf.Tensor: shape=(), dtype=string, numpy=b'bin red at s nine again'>]
```

```
In [63]:   1  yhat = model.predict(tf.expand_dims(sample[0], axis=0))
```

1/1 [==============================] - 1s 720ms/step

```
In [64]:   1  decoded = tf.keras.backend.ctc_decode(yhat, input_length=[75], greedy=True)[0][0].numpy()
```

```
In [65]:   1  print('~'*100, 'PREDICTIONS')
           2  [tf.strings.reduce_join([num_to_char(word) for word in sentence]) for sentence in decoded]
```

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~ PREDICTIONS

```
Out[65]: [<tf.Tensor: shape=(), dtype=string, numpy=b'bin red at s nine again'>]
```

*Figure 8: Testing on a video and deducing the output transcription*
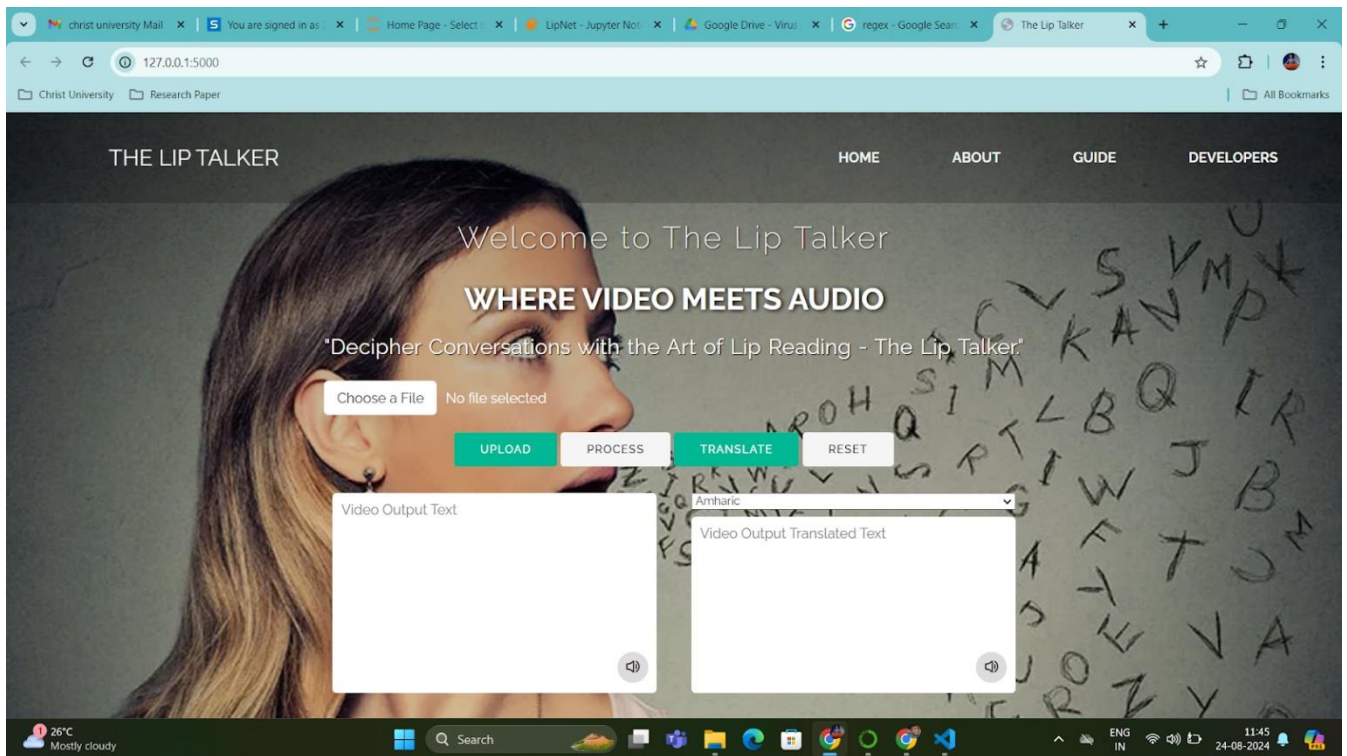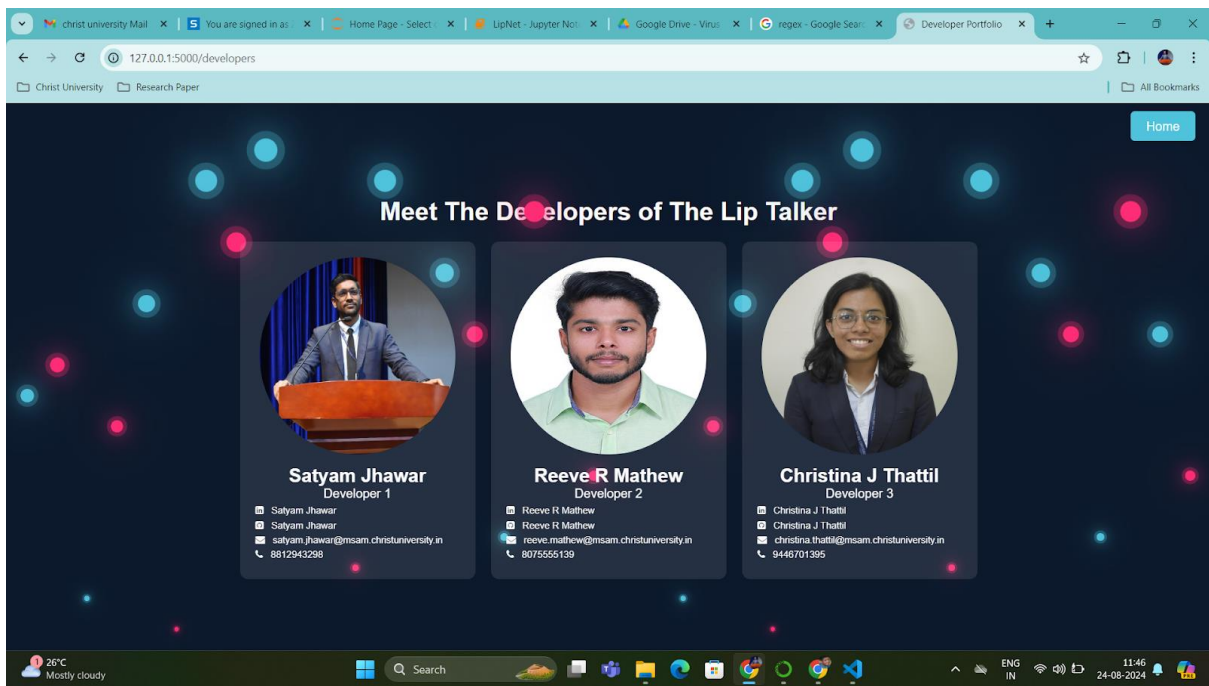


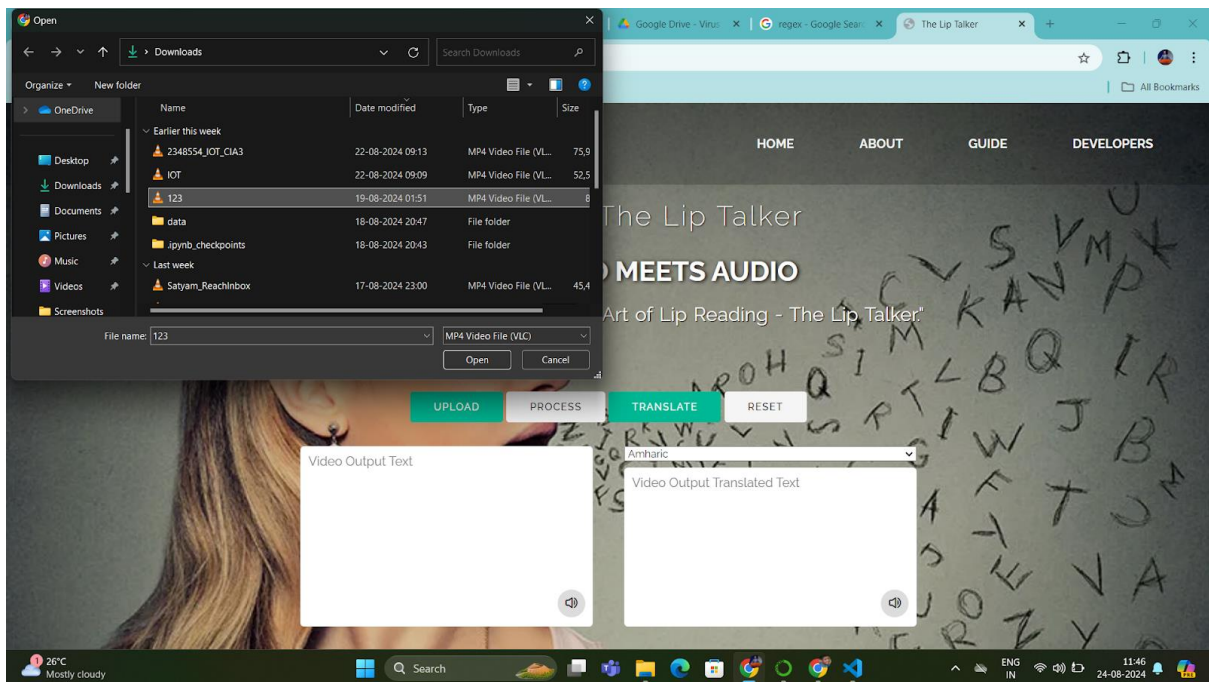*Figure 9: User Interface - Home page*

*Figure 10: Developers page*



*Figure 11: Clicking the UPLOAD button and uploading a video for transcription*
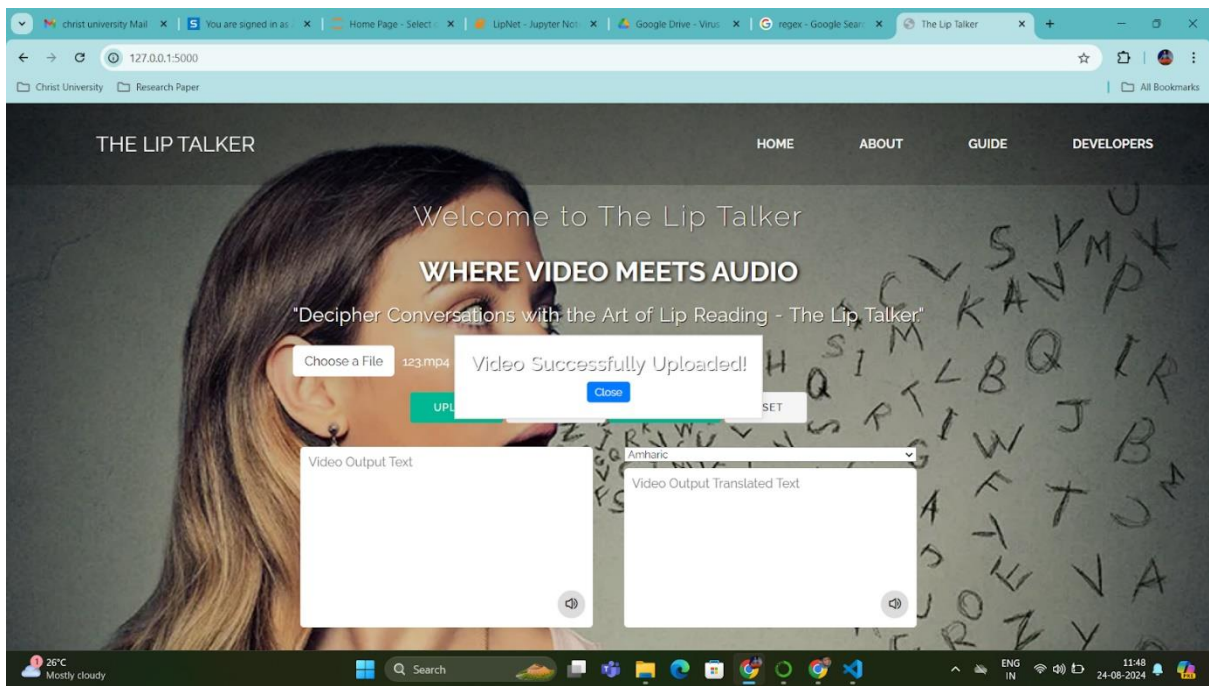
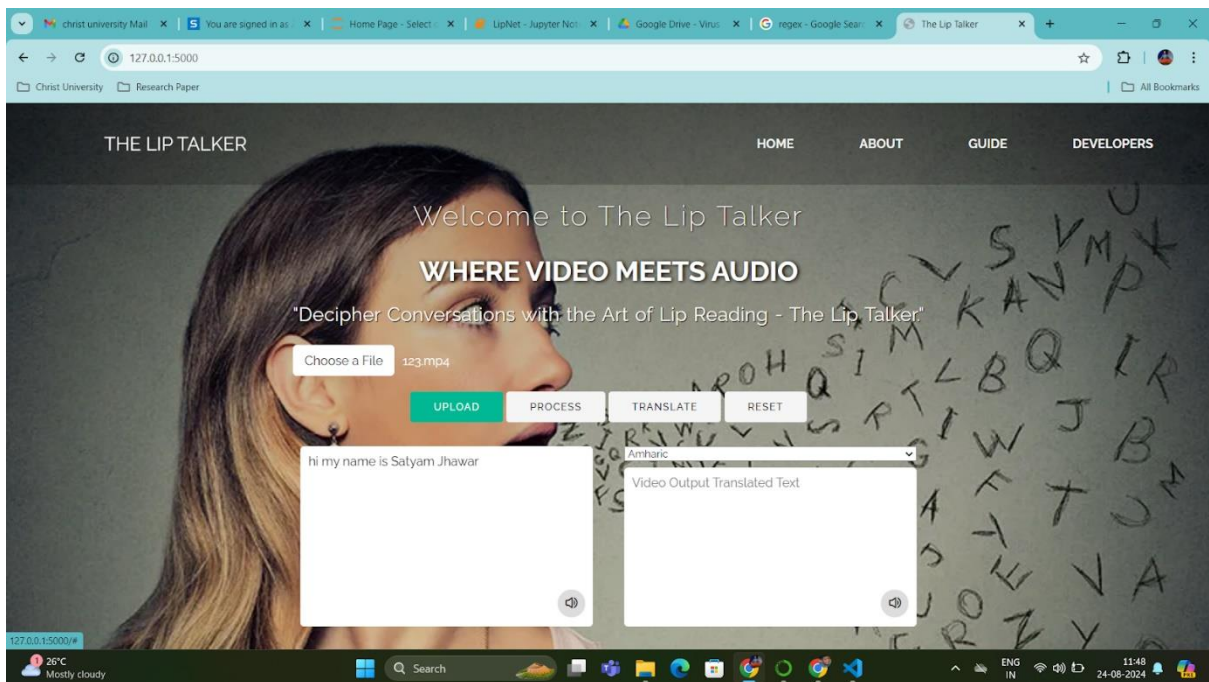*Figure 12: Video Successfully Uploaded dialog box after successful upload of the video*



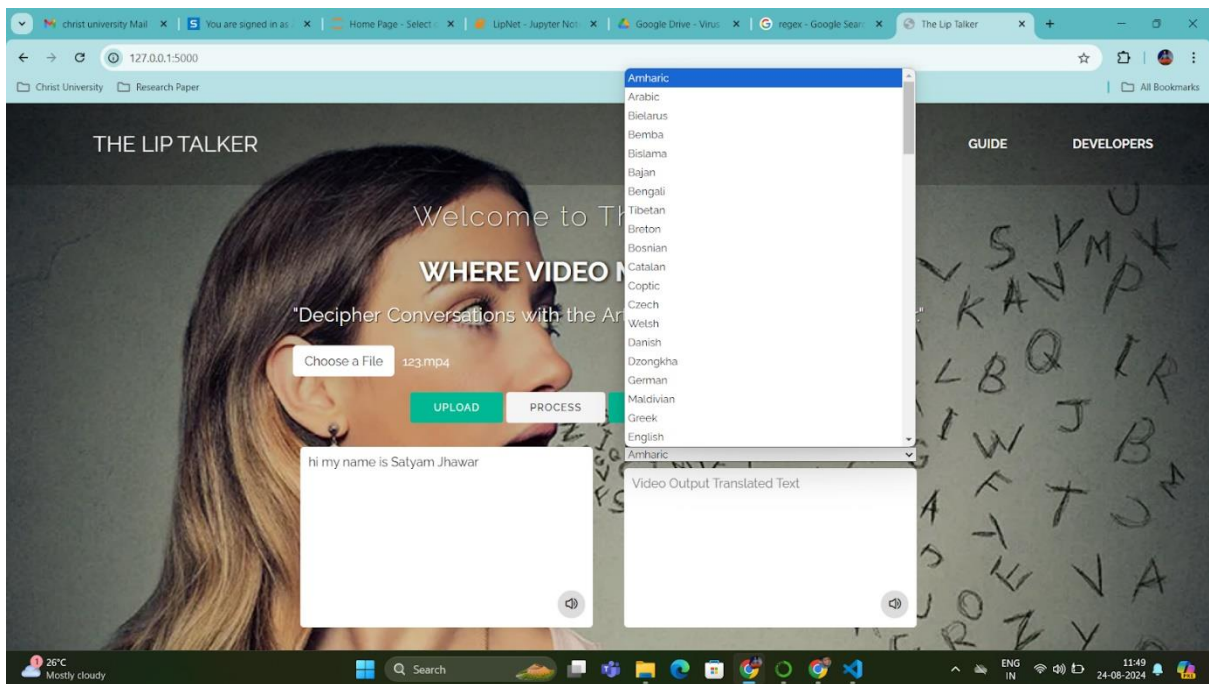*Figure 13: Processing the uploaded video*

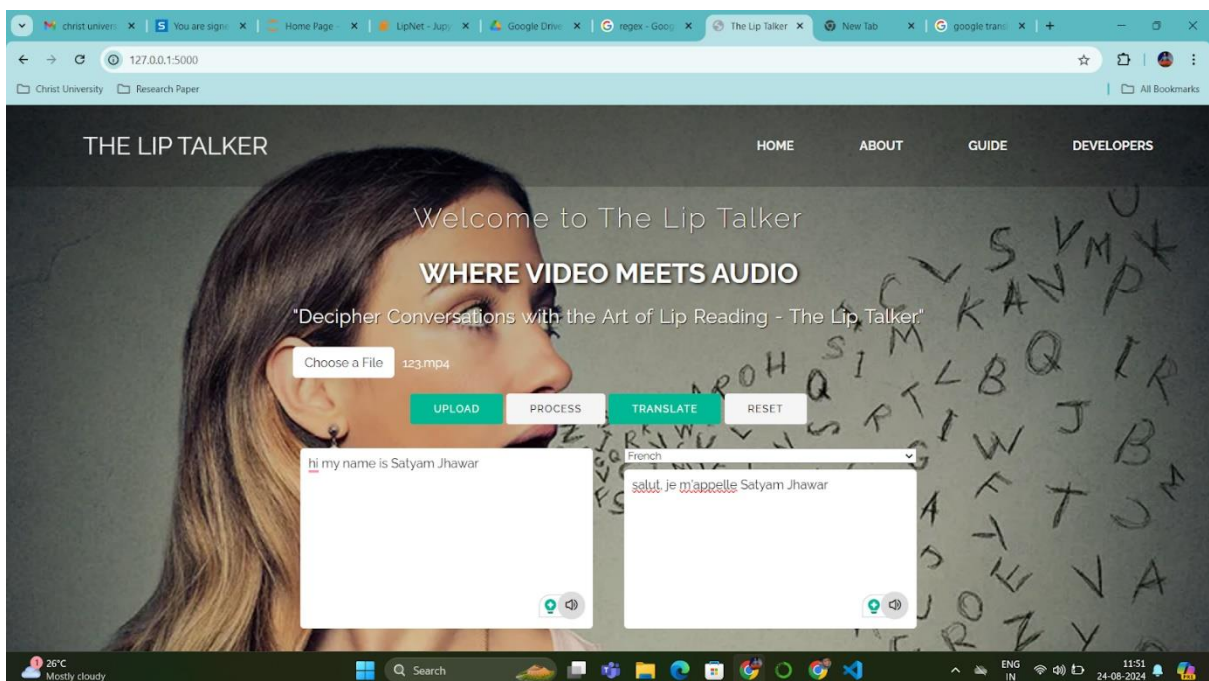*Figure 14: Selecting the language for translation*



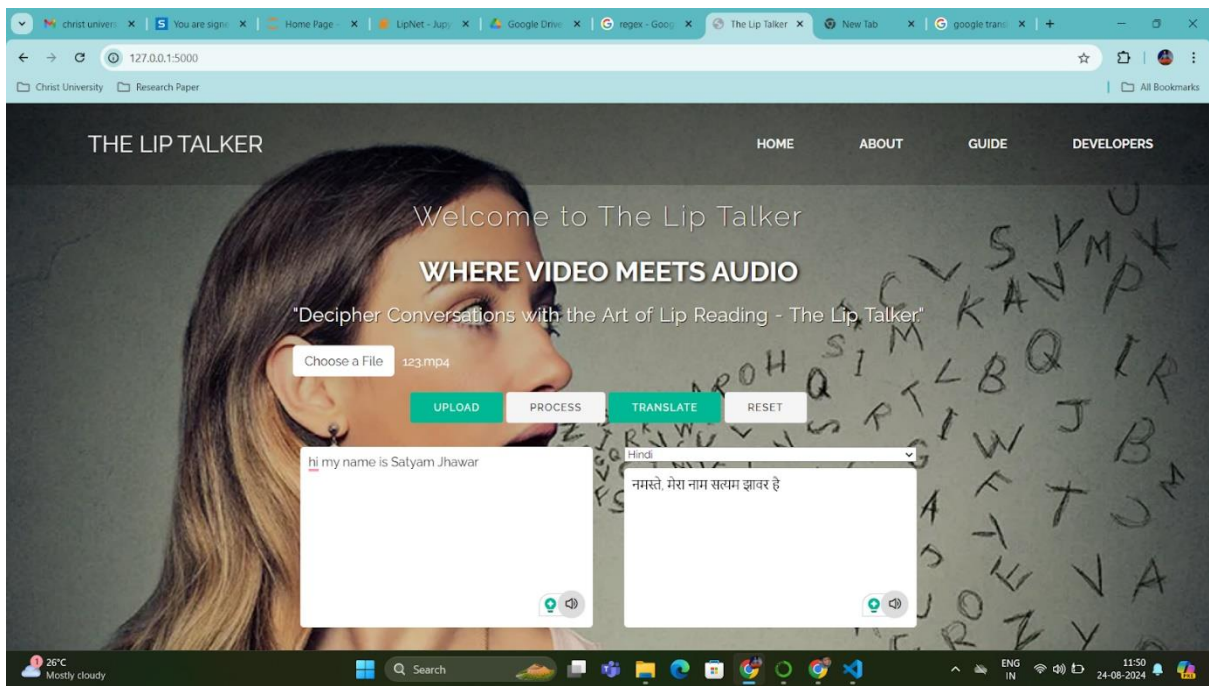*Figure 15: Translation of the transcription of the video in French*

*Figure 16: Translation of the transcription of the video in Hindi*

# 6. BENEFITS OF THE PROPOSED SYSTEM

**Lip Talker Upgrade:** Efficiency Improvement in Translation and Usage

Lip Talker project aims to address these limitations by integrating advanced deep learning techniques and Google Translate, bringing several key upgrades:

**Improved Accuracy:** Taking advantage of the State-of-art machine learning Lip Talker helps in improving the lip-reading, even in conditions when seeing a speaker's lips is somewhat difficult. For this reason, our models operate on various datasets to identify many lip movements with high accuracy.

**Multilingual Support**: In contrast to most of the existing systems, Lip Talker incorporates the possibilities of translating lip-read English text into other languages. This particular feature is managed by Google translate meaning that communicating becomes very easy even in completely different linguistic environments.

**User-Friendly Interface:** Lip Talker is very easy to use as it comes with a simple and graphical user interface that can be activated by using either the mouse or the keyboard by any user regardless of his or her IT literacy level of the computer. Usability is at the core of the design and therefore the target user is able to quickly and efficiently navigate through the tool.

Taking into account the shortcomings of the current systems mentioned above and the incorporation of such improvements into Lip Talker, it is suggested that the application will serve as an effective and efficient means of communication and access.

# 7. ADVANTAGES OF THE PROPOSED SYSTEM

1. **Advancements in the way the Deaf and Hard-of-Hearing Can and Do Communicate.:** Lip Talker helps the hearing-impaired folk communicate much better, thus eradicating the possibility of socially awkward situations and improving their lives. In this way, the lip movements are converted to text and translated, where the translations make for meaningful interaction in different contexts.

2. **Multilingual Support:** The system also uses Google API to translate current texts to other languages in real-time. This feature helps include all the participants in terms of communication and is especially useful in multicultural applications.

3. **User-Friendly Interface:** The idea behind the web application layout is to guarantee that users ranging from low-tech to high-tech users can conveniently upload videos, request text-to-video translations, and explore the desirable functions of the system. This makes the application easily accessible and opens the door to more users, resulting in a better user experience.

4. **Superior Security and Defence Features:** The actual real-time lip-reading feature in the system can be very effective in increasing situational awareness in defence and surveillance. The current work offers a reliable means of observing and analysing the communication processes in security operations and enhancing safety and intelligence gains.

5. **High accuracy and saving are functions of robust performance.:** As a result, it delivers a high level of accuracy when referring to lip gesture interpretation by applying advanced machine learning algorithms in Lip Talker. This reliability avoids misunderstanding and allows communication with higher efficiency. As a result, users get accurate translations.

6. **Social and Economic Impact:** Thus, Lip Talker positively affects the area of social inclusion and equal opportunities for members of the deaf and hard-of-hearing community. Such can ultimately result in better education, employment, and general social assimilation among other aspects of life.

7. **Technological Advancement:** That is why the proposal shows that integrating deep learning with translation technologies is feasible and can open many new possibilities in assistive technologies and tools for multilingual communications.

8. Lip Talker aims to create a transformative impact on how people with hearing impairments communicate, breaking down barriers and fostering a more inclusive and connected world.

# 8. CHALLENGES OF THE PROPOSED SYSTEM

1. **Data Privacy:** Maintaining the basic tenets of user uploaded videos, especially the privacy of the videos.

2. **Accuracy of Lip-Reading:** Besides, it is important to achieve high accuracy in various conditions such as the type of lighting, direction of the camera, the environment noise, etc.

3. **Real-Time Processing:** Being able to translate videos in real time and have the system be capable of handling a large number of videos at a time so that the procedures do not take a long time.

# 9.FUTURE SCOPE

1. **Model Enhancement:**

- Architecture Variations: Experiment with different architectures, such as attention mechanisms or transformer models, to potentially improve performance.

- Pretrained Models: Consider integrating pretrained models or transfer learning techniques to leverage existing knowledge and improve accuracy.

2. **Data Augmentation:**

- Diverse Data Sources: Incorporate more diverse video data to enhance model robustness and generalization.

- Augmentation Techniques: Apply data augmentation techniques such as random cropping, rotations, or color jittering to improve model performance and resilience.

3. **Performance Metrics:**

- Evaluation Metrics: Explore additional evaluation metrics specific to sequence-to-sequence tasks, such as Word Error Rate (WER) or Character Error Rate (CER), to better assess model accuracy.

4. **Real-World Applications:**

- Integration: Develop applications that utilize the trained model for real-time transcription of videos or other multimedia content.

- User Interface: Create user-friendly interfaces to allow non-experts to use the model for transcription tasks.

5. **Scalability and Deployment:**

- Scalability: Ensure the model and data pipeline are scalable to handle larger datasets or more complex video content.

- Deployment: Explore deployment options such as cloud services or mobile platforms to make the model accessible in various environments.

6. **Ethical Considerations:**

- Bias and Fairness: Assess and mitigate any biases in the model to ensure fair and ethical use in diverse applications.

- Privacy: Ensure that the data used for training and testing complies with privacy regulations and ethical guidelines.

# 10. CONCLUSION

Our project has achieved significant milestones in implementing a robust data pipeline for loading, preprocessing, and batching video data and alignments, crucial for training a deep learning model for sequence-to-sequence tasks. The model architecture, a combination of CNN-LSTM, effectively extracts features from video frames and predicts sequences by capturing spatial and temporal dependencies in the data. Additionally, the use of a custom loss function, CTC, enables the handling of variable-length sequences, such as transcriptions of spoken words from video frames, improving the model's ability to handle sequences without explicitly provided alignments. The successful training and evaluation processes, including the incorporation of pre-trained weights and making predictions, demonstrate the model's generalization to unseen video data. Looking forward, there are opportunities for refining and enhancing the current implementation to achieve even more promising results.

# REFERENCES

[1] N.Durga Sri, R. Akhil, S. Venkat Durga Prasad, V.Jayanth, and K.Krishna Jyothi, "Lip Reading Using Neural Networks and Deep Learning," International Journal of Scientific Research in Engineering and Management, vol. 7, no. 4, pp. 1-4, Apr. 2023.

[2] J. S. Chung, A. Senior, O. Vinyals, and A. Zisserman, "Lip Reading Sentences in the Wild," arXiv:1611.05358v2 [cs.CV], Nov. 2016.

[3] Y. M. Assael, B. Shillingford, S. Whiteson, and N. De Freitas, "LIPNET: End-to-end sentence-level lipreading," in *Proceedings of the 30th International Conference on Neural Information Processing Systems*, 2016.

[4] J. K. Chorowski, A. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, "Attention-based models for speech recognition," in Advances in Neural Information Processing Systems 28, Montreal, Canada, Dec. 7-12, 2015, pp. 577-585.

[5] J. K. Chorowski, A. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, "Hybrid CTC/Attention Model Based on Conformers," arXiv e-prints, arXiv:2102.06657, Feb. 2021.

[6] A. Koumparoulis and G. Potamianos, "Accurate and resource efficient lipreading with EfficientNetV2 and transformers," in ICASSP, 2022, pp. 8467-8471.

[7] P. Ma, S. Petridis, and M. Pantic, "End-to-end audio-visual speech recognition with conformers," in Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2021.

[8] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. In *Proceedings of the Nature*, 521(7553), 436-444. doi: 10.1038/nature14539.

[9] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. In *Neural Computation*, 9(8), 1735-1780. doi: 10.1162/neco.1997.9.8.1735.

[10] Schuster, M., & Paliwal, K. K. (1997). Bidirectional recurrent neural networks. In *IEEE Transactions on Signal Processing*, 45(11), 2673-2681. doi: 10.1109/78.650093.