

End-to-end Lip Reading deep learning project.

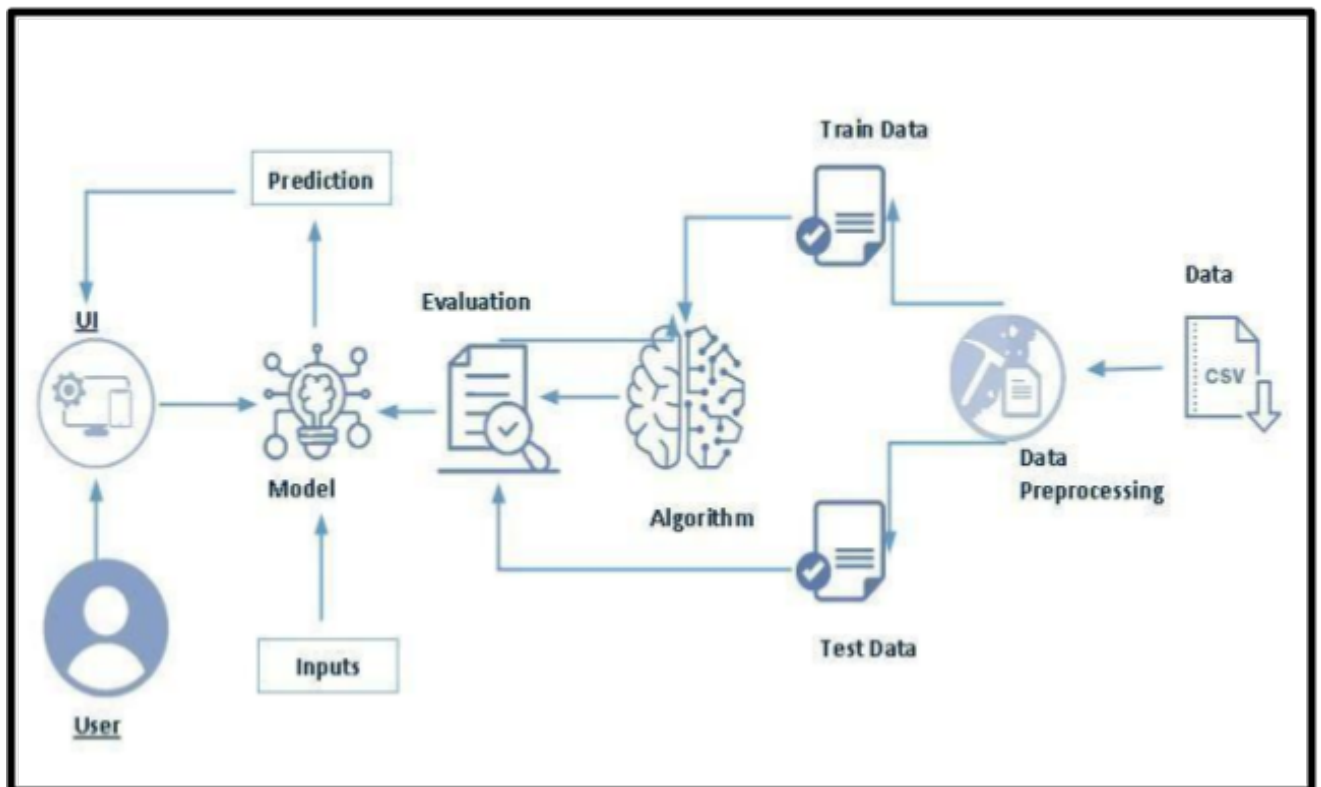
The objective of this project is to develop an end-to-end machine learning solution to detect words from a video of a person speaking. The proposed solution involves the use of Deep learning algorithms like LSTM, Neural Networks to predict the accurate output.

Lip reading using machine learning can offer several benefits:

1. **Improved Speech Recognition:** Lip reading can complement audio-based speech recognition systems, especially in noisy environments or scenarios where the audio signal is unclear. Integrating lip reading with traditional speech recognition can enhance accuracy and robustness.
2. **No Need for Audio Data:** Traditional speech recognition models require large amounts of transcribed audio data for training. In contrast, an end-to-end lip reading system can be trained solely on video data, eliminating the need for transcribed audio, which can be expensive and time-consuming to obtain.
3. **Multi-Modal Applications:** End-to-end lip reading can be combined with audio-based systems to create multi-modal applications. For example, in video conferencing, it can help improve real-time communication by providing more accurate transcriptions.
4. **Accessibility for Hearing-Impaired Individuals:** Lip reading can be an essential communication tool for individuals with hearing impairments. An accurate lip reading system can enhance their ability to understand spoken language and participate in conversations.

Let us look at the Technical Architecture of the project.

Technical Architecture:



Project Flow:

- The user interacts with the UI to select the file.
- Selected input is analyzed by the model which is integrated/developed by you.
- Once the model analyzes the input, the prediction is showcased on the UI. To accomplish this, we have to complete all the activities listed below,
- Define Problem / Problem Understanding
 - 0 Specify the business problem
 - o Business requirements
 - o Literature Survey.
 - o Social or Business Impact.
- Data Collection & Preparation
 - o Collect the dataset
 - o Data Preparation
- Exploratory Data Analysis
 - 0 Descriptive statistical
 - o Visual Analysis
- Model Building
 - o Building a model.
 - 0 Training the model.
 - o Testing the model
- Model Deployment
 - 0 Save the best model
 - o Integrate with Web Framework
- Project Demonstration & Documentation
 - 0 Record explanation Video for project end to end solution

- Project Documentation-Step by step project development procedure

Prior Knowledge:

To complete this project, you must require following software's , concepts and packages

- VS Code :
 - Refer to the link below to download VS Code.
 - Link : <https://code.visualstudio.com/download>
- Create Project:
 - Open VS Code.
 - Create a Folder and name it "Lip_Reading" .
- Machine Learning Concepts
 - Deep Learning:
<https://towardsdatascience.com/introducing-deep-learning-and-neural-networks-deep-learning-for-r-ookies-1-bd68f9cf5883>
 - NLP Models :
<https://medium.com/voice-tech-podcast/an-overview-of-rnn-lstm-gru-79ed642751c6>
- Web concepts:
 - Get the gist on streamlit :
<https://www.geeksforgeeks.org/a-beginners-guide-to-streamlit/>

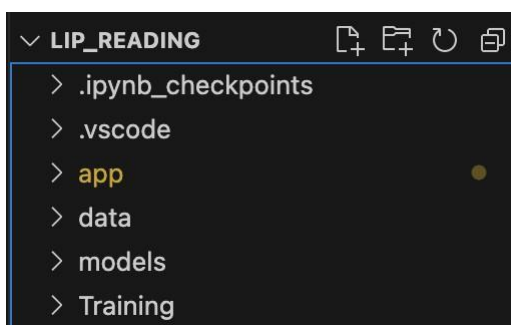
Project Objectives:

By the end of this project you will:

- Know fundamental concepts and techniques of Deep Learning.
- Gain a broad understanding of Lip Reading.
- Know how to train models in an efficient way.
- Know how to build a web application using the Streamlit framework.

Project Structure:

Create the Project folder which contains files as shown below



- We are building a Streamlit application which needs an app folder for a website.
- models folder contains your saved models.
- The Training folder contains the code for building/training/testing the model.
- The Dataset folder contains the videos and alignments.

Milestone 1: Define Problem / Problem Understanding

Activity 1: Specify the business problem

Refer Project Description

Activity 2: Business requirements

Here are some potential business requirements for lip reading using deep learning:

1. Accurate prediction: The predictor must be able to accurately predict the words. The accuracy of the prediction is crucial for the client. The client could be a corporate or a business person or anyone.
2. User-friendly interface: The predictor must have a user-friendly interface that is easy to navigate and understand. The interface should present the results of the predictor in a clear and concise manner.
3. Scalability: The predictor must be able to scale up based on the prediction from our product. The model should be able to handle any size of data without compromising on its accuracy or efficiency.

Activity 3: Literature Survey (Student Will Write)

A literature survey would involve researching and reviewing existing studies, articles, and other publications on the topic of the project. The survey would aim to gather information on current systems, their strengths and weaknesses, and any gaps in knowledge that the project could address. The literature survey would also look at the methods and techniques used in previous projects, and any relevant data or findings that could inform the design and implementation of the current project.

Activity 4: Social or Business Impact.

1. **Privacy and Security**: Audio-based speech recognition systems may raise privacy concerns, as they capture and process audio data, potentially infringing on individuals' privacy. Lip reading systems, on the other hand, rely on visual information and might be considered less intrusive in this regard.
2. **Use in Noisy Environments**: In environments with high background noise, audio-based speech recognition can be challenging. Lip reading can help provide context and improve accuracy in these noisy scenarios.
3. **Cross-Lingual Applications**: Lip reading is language-agnostic, which means the same model can potentially be applied to lip reading in different languages without requiring language-specific training data.

Milestone 2: Data Collection & Preparation

DL depends heavily on data. It is the most crucial aspect that makes algorithm training possible. So, this section allows you to acquire the required dataset.

NOTE : If you don't have a GPU, then train your model on kaggle.

Activity 1: Loading the data.

Link for the required dataset : <https://www.kaggle.com/datasets/rishisrddy/lipreading>

Download the provided data and load the data into the data folder

Activity 2: Data Preparation

As we have understood how the data is, let's pre-process the collected data.

Let's import the required packages

```
import os
import cv2
import tensorflow as tf
import numpy as np
from typing import List
from matplotlib import pyplot as plt
import imageio
```

We require 5 functions to process the data

- load_video()
- char_to_num
- num_to_char
- load_alignments()
- load_data()

- `load_video()`:

```
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
```

- `char_to_num`:

```
char_to_num = tf.keras.layers.StringLookup(vocabulary=vocab, oov_token="")
```

- `num_to_char`:

```
num_to_char = tf.keras.layers.StringLookup(
    vocabulary=char_to_num.get_vocabulary(), oov_token="", invert=True
)
```

- `load_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:]
```

- **load_data():**

```
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('/kaggle/input/lipreading/data', 's1', f'{file_name}.mpg')
    alignment_path = os.path.join('/kaggle/input/lipreading/data', 'alignments', 's1', f'{file_name}.align')
    frames = load_video(video_path)
    alignments = load_alignments(alignment_path)

    return frames, alignments
```

Let's convert the data into tensors using the above functions

```
test_path = '/kaggle/input/lipreading/data/s1/bbal6n.mpg'
tf.convert_to_tensor(test_path).numpy().decode('utf-8').split('/')[ -1].split('.')[0]
```

```
frames, alignments = load_data(tf.convert_to_tensor(test_path))
frames
```

The output will be :

```
<tf.Tensor: shape=(75, 46, 140, 1), dtype=float32, numpy=
array([[[[ 1.460374  ],
          [ 1.460374  ],
          [ 1.4209044 ],
          ...,
          [ 0.15787826],
          [ 0.19734783],
          [ 0.07893913]],

        [[ 1.460374  ],
          [ 1.460374  ],
          [ 1.4209044 ],
          ...,
          [ 0.15787826],
          [ 0.07893913],
          [ 0.03946957]],

        [[ 1.4209044 ],
          [ 1.4209044 ],
          [ 1.3814349 ],
          ...,
          [ 0.07893913],
          [ 0.07893913],
          [ 0.07893913]],

        ...,
        ...,
        ...]])
```

and the output goes on....

Milestone 3: Exploratory Data Analysis

Activity 1: Descriptive statistical

Descriptive analysis is to study the basic features of data with the statistical process. With this describe function we can understand the unique, top and frequent values of categorical features. And we can find mean, std, min, max and percentile values of continuous features. Which are not suitable for our dataset.

As our data consists of videos and alignments, There is no unnecessary data which can be eliminated.

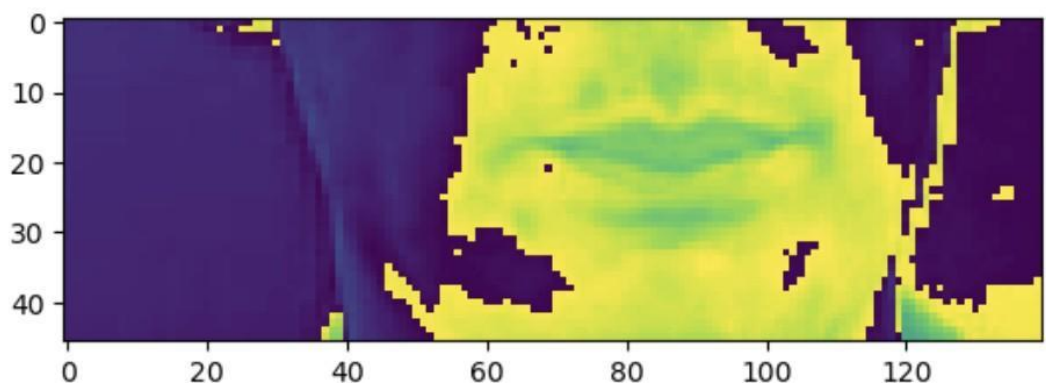
Activity 2: Visual analysis

Visual analysis is the process of using visual representations, such as charts, plots, and graphs, to explore and understand data. It is a way to quickly identify patterns, trends, and outliers in the data, which can help to gain insights and make informed decisions.

Let's plot the converted frames

```
plt.imshow(frames[15])
```

<matplotlib.image.AxesImage at 0x79962055e200>



Activity 3: Splitting data into train and test and validation sets

First we will create a mappable function

```
def mappable_function(path:str) ->List[str]:  
    result = tf.py_function(load_data, [path], (tf.float32, tf.int64))  
    return result
```

The following code is to divide the preprocessed data into train and test data equally.

```
data = tf.data.Dataset.list_files('preparing/data/s1/*.mpg')  
data = data.shuffle(500, reshuffle_each_iteration=False)  
data = data.map(mappable_function)  
data = data.padded_batch(2, padded_shapes=([75, None, None, None], [40]))  
data = data.prefetch(tf.data.AUTOTUNE)  
## Added for split  
train = data.take(450)  
test = data.skip(450)
```

The below codes are to print the preprocessed data

```
frames, alignments = data.as_numpy_iterator().next()
```

```
sample = data.as_numpy_iterator()
```

```
val = sample.next(); val[0]
```

Milestone 4: Model Building

Activity 1: Importing necessary libraries

We'll import some necessary libraries using the following code

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Conv3D, LSTM, Dense, Dropout, Bidirectional, MaxPool3D, Activation, Reshape, SpatialDropout3D, BatchNormalization, TimeDistributed, Flatten
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.callbacks import ModelCheckpoint, LearningRateScheduler
```

Activity 2: Defining callbacks, Loss function and building the model

Let's define the following:

- Callbacks :
 - scheduler()

```
def scheduler(epoch, lr):
    if epoch < 30:
        return lr
    else:
        return lr * tf.math.exp(-0.1)
```

- ProduceExample()

```
class ProduceExample(tf.keras.callbacks.Callback):
    def __init__(self, dataset) -> None:
        self.dataset = dataset.as_numpy_iterator()

    def on_epoch_end(self, epoch, logs=None) -> None:
        data = self.dataset.next()
        yhat = self.model.predict(data[0])
        decoded = tf.keras.backend.ctc_decode(yhat, [75,75], greedy=False)[0][0].numpy()
        for x in range(len(yhat)):
            print('Original:', tf.strings.reduce_join(num_to_char(data[1][x])).numpy().decode('utf-8'))
            print('Prediction:', tf.strings.reduce_join(num_to_char(decoded[x])).numpy().decode('utf-8'))
            print('~'*100)
```

- Loss Function:
 - CTCLoss() :

```
def CTCLoss(y_true, y_pred):
    batch_len = tf.cast(tf.shape(y_true)[0], dtype="int64")
    input_length = tf.cast(tf.shape(y_pred)[1], dtype="int64")
    label_length = tf.cast(tf.shape(y_true)[1], dtype="int64")

    input_length = input_length * tf.ones(shape=(batch_len, 1), dtype="int64")
    label_length = label_length * tf.ones(shape=(batch_len, 1), dtype="int64")

    loss = tf.keras.backend.ctc_batch_cost(y_true, y_pred, input_length, label_length)
    return loss
```

- Model Building :

```
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(char_to_num.vocabulary_size()+1, kernel_initializer='he_normal', activation='softmax'))
```

Activity 3: Train the models

First initialize the callback functions

```
checkpoint_callback = ModelCheckpoint(os.path.join('models','checkpoint'), monitor='loss', save_weights_only=True)
schedule_callback = LearningRateScheduler(scheduler)
example_callback = ProduceExample(test)
```

Then we compile the model

```
model.compile(optimizer=Adam(learning_rate=0.0001), loss=CTCLoss)
```

Now it's time to train the model

```
model.fit(train, validation_data=test, epochs=100, callbacks=[checkpoint_callback, schedule_callback, example_callback])
```

```

Epoch 1/100
340/450 [=====>.....] - ETA: 1:25 - loss: 88.6656
[mpeg1video @ 0x79970014ae00] ac-tex damaged at 22 17
[mpeg1video @ 0x79970014ae00] Warning MVs not available
450/450 [=====] - ETA: 0s - loss: 84.8949
[mpeg1video @ 0x7996f0036c00] ac-tex damaged at 22 17
[mpeg1video @ 0x7996f0036c00] Warning MVs not available
450/450 [=====] - 621s 1s/step - loss: 84.8949 - val_loss: 69.6304 - lr: 1.0000e-04
Epoch 2/100
450/450 [=====] - ETA: 0s - loss: 71.1230
[mpeg1video @ 0x7996d408a800] ac-tex damaged at 22 17
[mpeg1video @ 0x7996d408a800] Warning MVs not available
450/450 [=====] - 575s 1s/step - loss: 71.1230 - val_loss: 66.3392 - lr: 1.0000e-04
Epoch 3/100
106/450 [=====>.....] - ETA: 4:17 - loss: 69.5230
[mpeg1video @ 0x7996f29137c0] ac-tex damaged at 22 17
[mpeg1video @ 0x7996f29137c0] Warning MVs not available
450/450 [=====] - ETA: 0s - loss: 67.6537
[mpeg1video @ 0x7996f80ab900] ac-tex damaged at 22 17
[mpeg1video @ 0x7996f80ab900] Warning MVs not available
450/450 [=====] - 574s 1s/step - loss: 67.6537 - val_loss: 62.6463 - lr: 1.0000e-04
Epoch 4/100
20/450 [>.....] - ETA: 5:17 - loss: 65.4729

```

Let's wait for the model to train until 100 epochs.

Milestone 5: Model Deployment

Activity 1: Save the model

Your model is already saved in checkpoints

Zip the latest checkpoint and upload it to your google drive for future use.

Activity 2: Make a Prediction

```

url = 'https://drive.google.com/uc?id=1vWscXs4Vt0a_1IH1-ct2TCgXAZT-N3_Y'
output = 'checkpoints.zip'
gdown.download(url, output, quiet=False)
gdown.extractall('checkpoints.zip', 'models')

```

If you are having difficulties in training the model, you can download the weights from the above URL.

```

model.load_weights('/kaggle/working/LipNet.h5')

```

Now, Let's make the prediction on text data.

```

test_data = test.as_numpy_iterator()

```

```

sample = test_data.next()

```



```
yhat = model.predict(sample[0])
```

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

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

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

Now, Let's make prediction on video

```
sample = load_data(tf.convert_to_tensor('/kaggle/input/lipreading/data/s1/bbaf5a.mpg'))
```

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

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

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

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

If your prediction text is matching real text then move forward for the next step.

If not, train your model again.

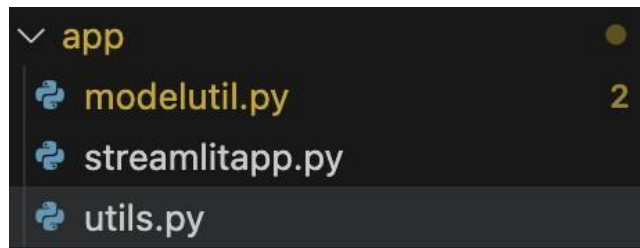
Activity 3: Integrate with Web Framework

In this section, we will be building a web application that is integrated to the model we built.

We will be using the streamlit package for our website development.

Streamlit is a free and open-source framework to rapidly build and share beautiful machine learning and data science web apps.

Activity 4.1: Create an app folder and create following .py files in it:



Let's write the code for “utils.py”

- Importing necessary packages:

```
import tensorflow as tf
from typing import List
import cv2
import os
```

- Defining char_to_num, num_to_char, load_video(), load_alignments(), load_data() functions:

```
vocab = [x for x in "abcdefghijklmnopqrstuvwxyz'?!123456789 "]
char_to_num = tf.keras.layers.StringLookup(vocabulary=vocab, oov_token="")
# Mapping integers back to original characters
num_to_char = tf.keras.layers.StringLookup(
    vocabulary=char_to_num.get_vocabulary(), oov_token="", invert=True
)
```

```
def load_video(path:str) -> List[float]:
    #print(path)
    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_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:]
```

```
def load_data(path: str):
    path = bytes.decode(path.numpy())
    file_name = path.split('/')[-1].split('.')[0]
    tf.convert_to_tensor(path).numpy().decode('utf-8').split('/')[-1].split('.')[0]
    # File name splitting for windows
    # file_name = path.split('\\')[-1].split('.')[0]
    video_path = os.path.join('./data/', 's1', f'{file_name}.mpg')
    #data/alignments/s1/bbaf4p.align
    alignment_path = os.path.join('./data/', 'alignments', 's1', f'{file_name}.align')
    frames = load_video(video_path)
    alignments = load_alignments(alignment_path)

    return frames, alignments
```

Let's write code for “modelutils.py”

- Import all the necessary packages:

```
import os
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Conv3D, LSTM, Dense, Dropout, Bidirectional, MaxPool3D, Activation, Reshape, SpatialDropout3D, BatchNormalization, TimeDistributed, Flatten
```

- Defining load_model() function :

```
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 '..', 'models', 'checkpoint'))

    return model
```

Let's write code for “**streamlitapp.py**”

- **Importing necessary packages and functions:**

```
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')
```

NOTE : If you get any error like “streamlit not found”. You have to make sure streamlit is installed on your system.

- Running the command “pip install streamlit” would do that job for you.

- **Creating sidebar:**

```
with st.sidebar:
    st.image('https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcSnGP00Pa-c3bWfl-E32Pkg_otJrA')
    st.title('LipBuddy')
    st.info('This application is originally developed from the LipNet deep learning.')
```

You can choose any image of your wish.

- **Drop down menu for selecting the input :**

```
options = os.listdir(os.path.join('../', 'data', 's1'))
options.sort()
options.remove('Thumbs.db')
selected_video = st.selectbox('Choose Video', options)
```

- **Creating two columns :**

```
col1, col2 = st.columns(2)
```

- **col1 and col2 :**

```
if options:
    with col1:
        st.info('The video below displays the converted video in mp4 format')
        file_path = os.path.join('/Users/rishi/BTECH/Programming/My_Projects/Lip_Reading/data/', 's1', selected_video)
        os.system(f'ffmpeg -i {file_path} -vcodec libx264 test_video.mp4 -y')

        # Rendering inside of the app
        video = open('test_video.mp4', 'rb')
        video_bytes = video.read()
        st.video(video_bytes)

    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))
        imageio.mimsave('animation.gif', video, fps=10)
        st.image('animation.gif', width=400)

        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)

        # Convert prediction to text
        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)
```

- col1 will display the selected video and play it.
- col2 will display :
 - animation of lips.
 - raw output of your model
 - result prediction

To run your website go to your terminal with your respective directory and run the command :

- “streamlit run streamlitapp.py”

You can now view your Streamlit app in your browser.

Local URL: <http://localhost:8501>

Network URL: <http://192.168.1.11:8501>

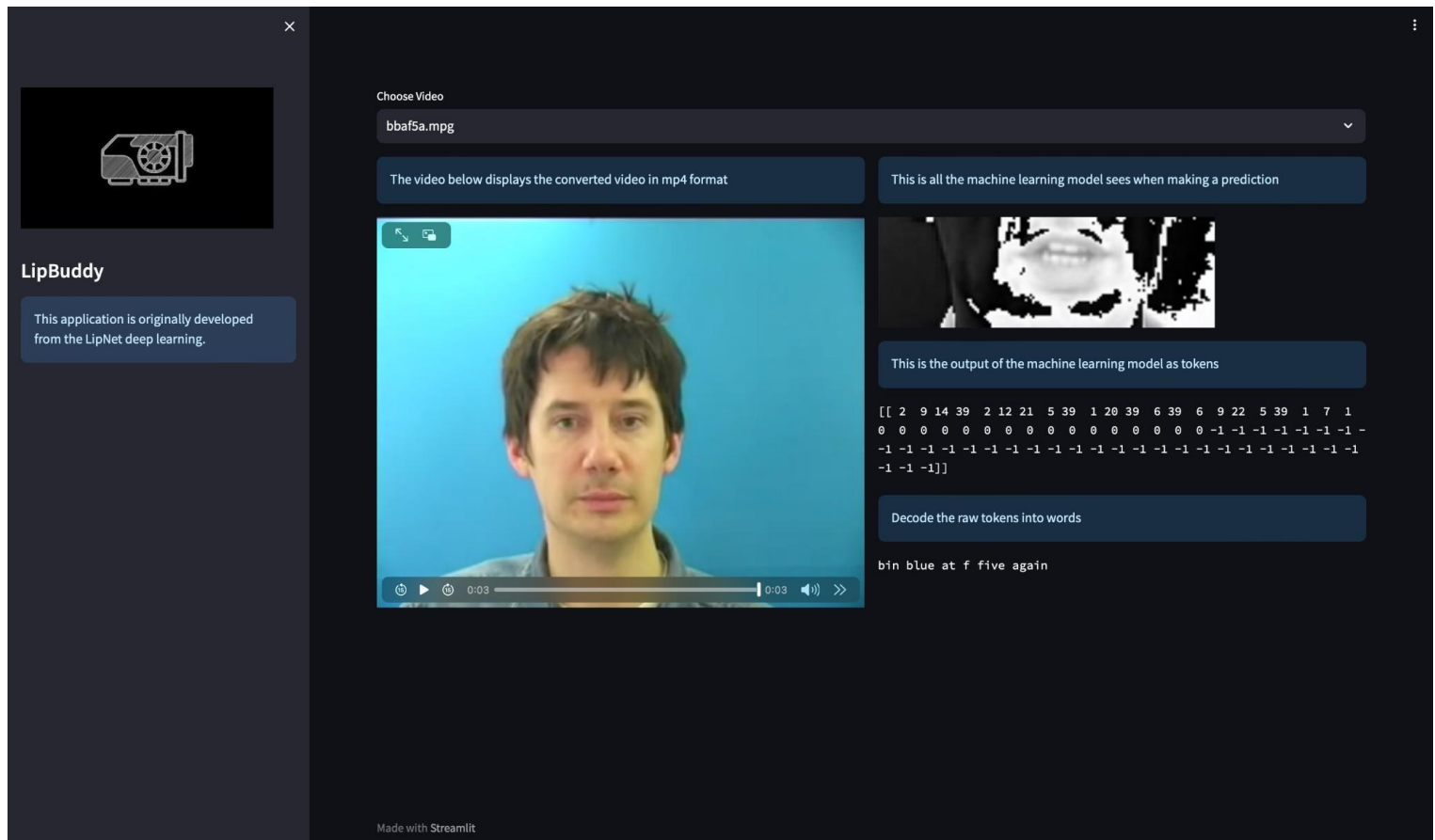
For better performance, install the Watchdog module:

```
$ xcode-select --install
```

```
$ pip install watchdog
```

On successful execution you'll get to see this in your terminal.

HOW DOES YOUR WEBSITE LOOK LIKE ?



Milestone 6: Project Demonstration & Documentation

Below mentioned deliverables to be submitted along with other deliverables.

Activity 1: - Record explanation Video for project end to end solution. Activity 2: - Project Documentation-Step by step project development procedure.

Create a document as per the template provided.

