

LIPSENSE

DEEP LEARNING BASED SILENT SPEECH RECOGNITION USING 3D CONVOLUTIONAL NEURAL NETWORKS AND GATED RECURRENT UNITS

1. Introduction

1.1 Overview

The field of speech recognition has witnessed significant advancements over the years, with a focus on auditory input for deciphering spoken language. However, in scenarios where auditory information is either unavailable or compromised, the comprehension of speech becomes challenging. This is particularly pertinent for individuals with hearing impairments or in noise-affected environments. In response to this challenge, the present study delves into the realm of silent speech recognition, capitalising on the visual modality of speech perception. The project seeks to develop an innovative solution by harnessing deep learning techniques, specifically, a fusion of 3D Convolutional Neural Networks (CNNs) and Gated Recurrent Units (GRUs), to facilitate automatic lip reading. The resultant model endeavours to transform lip movements into coherent textual representations, thereby bridging the gap between speech and visual cues.

1.2 Purpose

Speech is not always an effective form of communication under different circumstances. First, the environment could not be conducive to speech-based communication; for example, someone might be in a noisy restaurant or close to a crowded market where it is difficult to understand speech. Second, a person can be in a public place where others don't want to be disturbed, like a library or museum, and not want to talk aloud out of worries for their privacy and confidentiality. Finally, and most crucially, a large number of people struggle to talk or are completely speechless as a result of a variety of neurological and speech impairments.

The overarching purpose of this project is to engineer a sophisticated and accurate silent speech recognition system by means of a novel fusion of 3D CNNs and GRUs. This system has the potential to not only augment the quality of life for individuals with hearing impairments but also to transcend to diverse domains where audio data may be compromised or unavailable. In envisaging an innovative solution that leverages the inherent cues present in lip movements,

this project strives to empower individuals to communicate effectively through visual speech cues. Moreover, the fusion of 3D CNNs and GRUs promises to yield a comprehensive understanding of spatio-temporal aspects of lip motion, leading to enhanced accuracy and robustness in silent speech recognition. As such, this project seeks to break barriers in communication, revolutionising the way speech is perceived and interpreted in various spheres of human interaction and technology deployment.

2. Literature Survey

2.1 Existing Problem

This section presents an overview of prevailing methodologies in the domain of automated lipreading. The existing landscape encompasses a range of techniques, where the emphasis often lies on overcoming challenges associated with speech recognition in noise-affected, privacy-conscious, or speech-impaired scenarios. Prior endeavours delve into automated lipreading, utilising diverse methods that do not extensively harness deep learning. Many such approaches necessitate intricate preprocessing of frames for feature extraction, temporal preprocessing for video feature extraction, or employ manually crafted vision pipelines.

Goldschen et al. pioneered visual-only sentence-level lipreading through hidden Markov models (HMMs), albeit in a limited dataset. Neti et al. ventured into sentence-level audiovisual speech recognition, melding HMMs with engineered features. Akin to these endeavours, Potamianos et al. reported substantial speaker-independent and speaker-adapted improvements using similar techniques. Furthermore, Gergen et al. attained the then-state-of-the-art on the GRID corpus by employing speaker-dependent training in an HMM/GMM system, emphasising the challenges of generalisation across speakers and motion feature extraction.

Classification with deep learning has surfaced in recent years, often targeting word or phoneme classification. However, the transformative aspect of full sentence sequence prediction, a cornerstone of LipNet's approach, remains less explored. Notable attempts include Ngiam et al., Sui et al., Ninomiya et al., Petridis & Pantic, and others, who engage in multimodal audio-visual representations and integrate visual features into traditional speech-style processing pipelines. Spatial and spatiotemporal convolutional neural networks, proposed by Chung & Zisserman, demonstrate promise in word classification, yet lack the robustness for sentence-level sequence prediction. Similarly, other studies like Chung & Zisserman (audio-visual max-margin matching

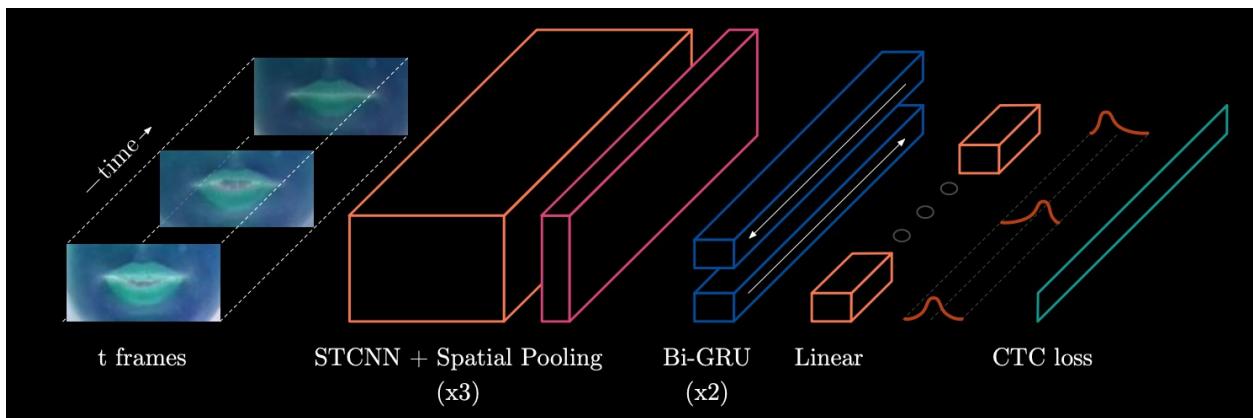
model), Wand et al. (LSTM recurrent neural networks), and Garg et al. (VGG pre-trained model) have yielded advancements within specific classifications, though not entirely encompassing sentence-level sequence prediction.

2.2 Proposed Solution

In light of the existing landscape, this project advocates for the adoption of LipNet—a groundbreaking end-to-end model that pioneers sentence-level sequence prediction for visual speech recognition. Unlike prior methodologies, LipNet efficiently processes sequences of images and yields output distributions over sequences of tokens. This innovative framework leverages the power of Connectionist Temporal Classification (CTC) loss, facilitating end-to-end training and obviating the need for intricate alignments. In a field that has often paralleled early automatic speech recognition (ASR) progress, LipNet emerges as a notable leap by transcending the limits of classification-based approaches, enhancing sentence-level sequence prediction.

3.Theoretical Analysis

3.1 Block Diagram



3.2 Hardware and Software Design

During the training phase, a powerful GPU is essential for accelerated deep learning model training. For this project, an Nvidia RTX 2070 Super was utilized. Specifications include:

- NVIDIA GeForce RTX 2070 Super
- CUDA Cores: 2560
- Base Clock: 1605 MHz
- Boost Clock: 1770 MHz
- VRAM: 8 GB GDDR6

- Memory Speed: 14 Gbps
- Memory Bus: 256-bit

For running the wep app a standard configuration of sufficient RAM and CPU is required.

Software frameworks used:

- IBM Cloud
- IBM Watson Machine Learning API Client
- Python
- Numpy
- Tensorflow
- OpenCV
- Matplotlib
- ImageIO
- Streamlit

4.Experimental Investigations

4.1 Data collection

Evaluating which dataset would be the most suitable for training the deep learning model was analysed. The Grid corpus dataset was chosen. The dataset was then split to training, testing and validation sections.

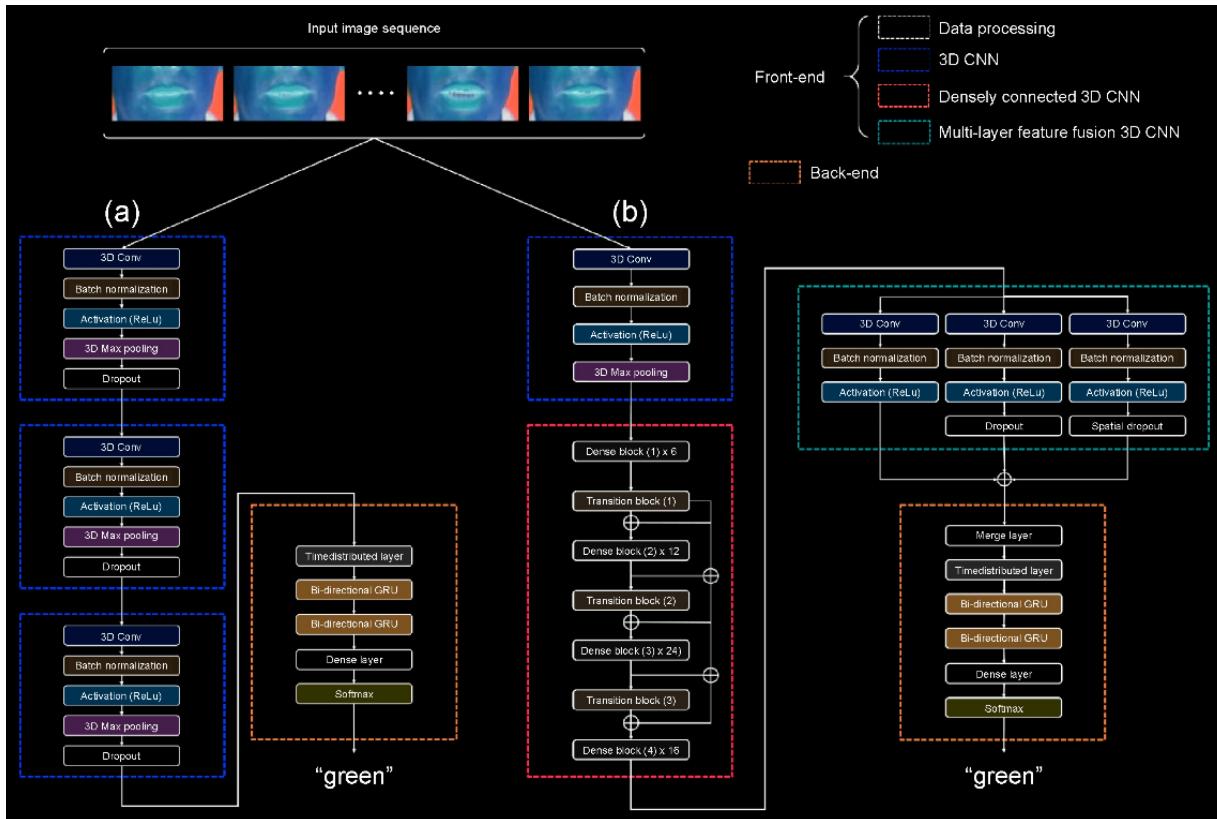
4.2 Neural Architecture

Then the design and architecture of the silent speech recognizer was chosen. The Lipnet architecture was strategically chosen to holistically capture both spatial and temporal features embedded in the lip motion sequences. The architecture is based upon 3D CNN and Bi-GRUs.

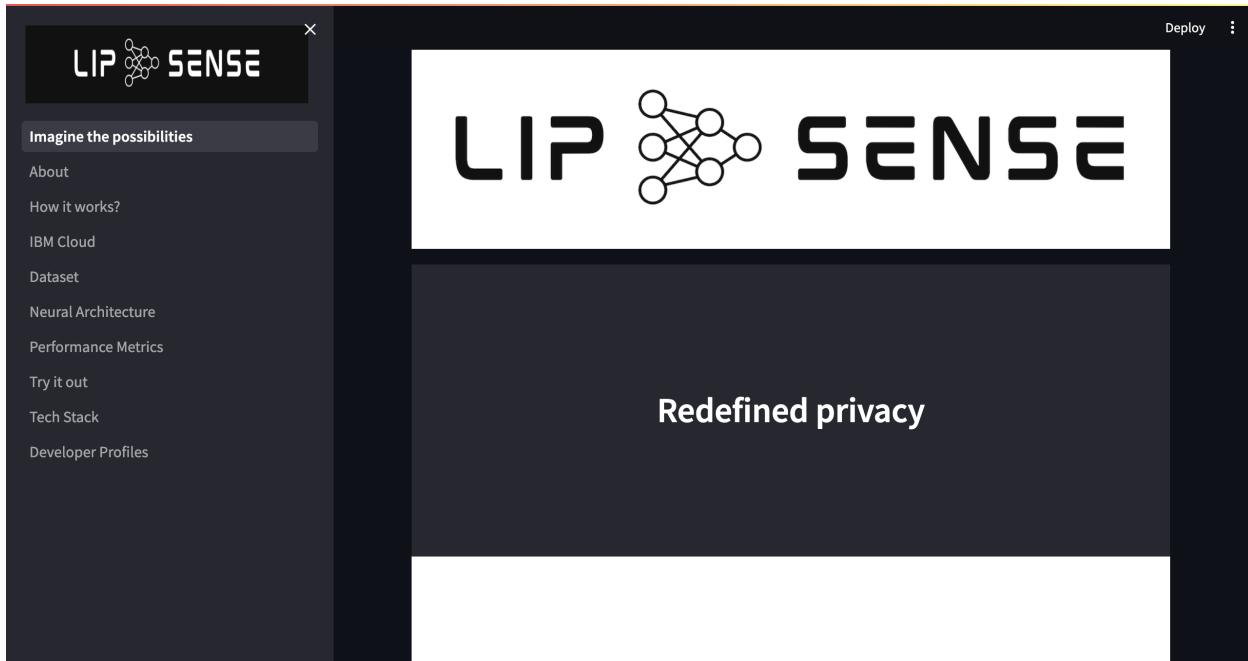
4.3 Training and Evaluation Metrics

The training process was refined. We periodically performed validation to monitor accuracy and mitigate overfitting concerns. Thorough hyperparameter optimization was carried out. This involved fine-tuning parameters like learning rate, batch size, and dropout rates to maximise model performance. Word error rate and character error rate were then analysed to evaluate the performance of the model.

5. Flowchart



6.Result



LIP SENSE

Imagine the possibilities

About

How it works?

IBM Cloud

Dataset

Neural Architecture

Performance Metrics

Try it out

Tech Stack

Developer Profiles

Deploy :

Enhanced Accessibility

Shaping the future of Human Computer Interaction

LIP SENSE

Imagine the possibilities

About

How it works?

IBM Cloud

Dataset

Neural Architecture

Performance Metrics

Try it out

Tech Stack

Developer Profiles

Deploy :

Description

Lip reading, the ability to understand spoken language by observing lip movements, is an essential skill for individuals with hearing impairments and in scenarios where audio information is compromised. Automatic lip reading systems, powered by deep learning techniques, have shown promising results in converting lip movements into text. This project proposes an advanced lip reading model that combines 3D Convolutional Neural Networks (CNNs) and Gated Recurrent Units (GRUs) to improve the accuracy and robustness of lip reading. The project's approach involves leveraging the spatio-temporal characteristics of lip motion using 3D CNNs. This allows the model to analyze sequential frames of lip movements and capture both spatial and temporal features simultaneously. The GRU component is employed to learn long-term dependencies and capture temporal dynamics in the lip movement sequences.

Made with Streamlit

LIP SENSE

Imagine the possibilities

About

How it works?

IBM Cloud

Dataset

Neural Architecture

Performance Metrics

Try it out

Tech Stack

Developer Profiles

How it works?

- The LipNet architecture operates through a meticulously designed sequence of operations.
- It begins with a series of three stages, each involving spatiotemporal convolutions, channel-wise dropout, and spatial max-pooling.
- This initial stage extracts crucial information from the input, encapsulating essential features.
- The architecture introduces two Bidirectional Gated Recurrent Units (Bi-GRUs), a pivotal element in the process.
- These Bi-GRUs facilitate the effective aggregation of outputs from the earlier spatiotemporal convolutions.
- Following this, each time-step undergoes a linear transformation.
- Subsequently, a softmax operation covers a vocabulary expanded to encompass the CTC blank symbol.
- The architecture's learning process is driven by the application of the Connectionist Temporal Classification (CTC) loss.
- Throughout the architecture, the Rectified Linear Unit (ReLU) serves as the activation function in all layers.
- This consistent choice of activation function contributes to the model's efficiency and performance.
- The orchestrated sequence of operations culminates in the LipNet architecture's proficiency in lip reading tasks.

LIP SENSE

Imagine the possibilities

About

How it works?

IBM Cloud

Dataset

Neural Architecture

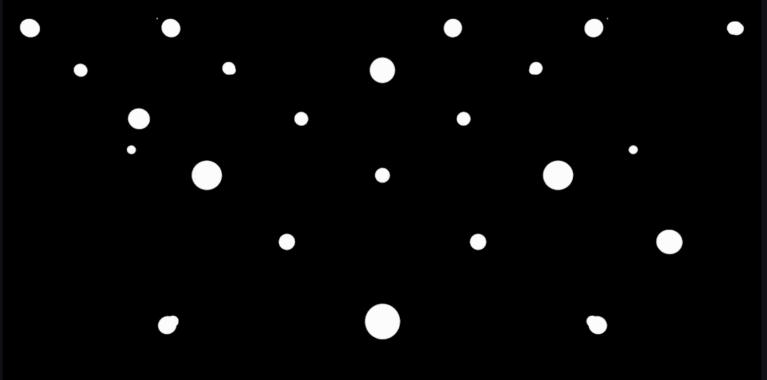
Performance Metrics

Try it out

Tech Stack

Developer Profiles

Deployment made easy with IBM Cloud



IBM Watson Machine Learning is a full-service IBM Cloud offering that makes it easy for developers and data scientists to work together to integrate predictive capabilities with their applications. The Watson Machine Learning service is a set of REST APIs that you can call from any programming language to develop applications that make smarter decisions, solve tough problems, and improve user outcomes.

LIP SENSE

Imagine the possibilities

- About
- How it works?
- IBM Cloud**
- Dataset
- Neural Architecture
- Performance Metrics
- Try it out
- Tech Stack
- Developer Profiles

IBM Watson Machine Learning is a full-service IBM Cloud offering that makes it easy for developers and data scientists to work together to integrate predictive capabilities with their applications. The Watson Machine Learning service is a set of REST APIs that you can call from any programming language to develop applications that make smarter decisions, solve tough problems, and improve user outcomes.

```

model_details=client.repository.store_model(model="lipbud.tgz",meta_props={
    client.repository.ModelMetaNames.NAME:"lipbud.h5",
    client.repository.ModelMetaNames.TYPE:"tensorflow_2.9",
    client.repository.ModelMetaNames.SOFTWARE_SPEC_UID:software_space_uid
})
[74] ✓ 19.0s
Loading...
model_id=client.repository.get_model_id(model_details)
model_id
[75] ✓ 0.0s
... '928ca4b3-02aa-4c90-9eaf-72dc26a15ff8'

client.repository.download(model_id,'lipbud3.tgz')
[77] ✓ 22.0s
... Successfully saved model content to file: 'Lipbud3.tgz'

```

Using the IBM Watson Machine Learning API client

LIP SENSE

Imagine the possibilities

- About
- How it works?
- IBM Cloud**
- Dataset**
- Neural Architecture
- Performance Metrics
- Try it out
- Tech Stack
- Developer Profiles

Dataset

We employ the GRID corpus for training due to its sentence-level nature and extensive dataset. The sentences are structured using a straightforward grammar pattern: command(4) + color(4) + preposition(4) + letter(25) + digit(10) + adverb(4). Each number indicates the available word choices within the corresponding word category. These categories encompass {bin, lay, place, set}, {blue, green, red, white}, {at, by, in, with}, {A,...Z}{W}, {zero, ..., nine}, and {again, now, please, soon}, resulting in a pool of 64,000 potential sentences. For instance, exemplars from this dataset include 'set blue by A four please' and 'place red at C zero again'. Model training was conducted on the first 450 videos of speaker one.

	#1	#11	#21	#31	#41	#51	#61	#71
Ground truth								
[Zhou et al., 2019]								
[Chen et al., 2019]								
irs								

LIP SENSE

Imagine the possibilities

- About
- How it works?
- IBM Cloud
- Dataset
- Neural Architecture**
- Performance Metrics
- Try it out
- Tech Stack
- Developer Profiles

Neural Architecture

Layer (type)	Output Shape	Param #
0 conv3d_3 (Conv3D)	(None, 75, 46, 140, 128)	3584
1 activation_3 (Activation)	(None, 75, 46, 140, 128)	0
2 max_pooling3d_3 (MaxPooling3D)	(None, 75, 23, 70, 128)	0
3 conv3d_4 (Conv3D)	(None, 75, 23, 70, 256)	884992
4 activation_4 (Activation)	(None, 75, 23, 70, 256)	0
5 max_pooling3d_4 (MaxPooling3D)	(None, 75, 11, 35, 256)	0
6 conv3d_5 (Conv3D)	(None, 75, 11, 35, 75)	518475
7 activation_5 (Activation)	(None, 75, 11, 35, 75)	0

LIP SENSE

Imagine the possibilities

- About
- How it works?
- IBM Cloud
- Dataset
- Neural Architecture
- Performance Metrics**
- Try it out
- Tech Stack
- Developer Profiles

Performance Metrics

Outperforms human by	4.1x	Word Error Rate	4.8%	Sentence accuracy	95.2%
Training Videos	450	Unseen CER	6.40%	Vocabulary tokens	41
Trainable parameters	8471924	STCNN Layers	x3	Bi-GRU	x2

LIP SENSE

Imagine the possibilities

About

How it works?

IBM Cloud

Dataset

Neural Architecture

Performance Metrics

Try it out

Tech Stack

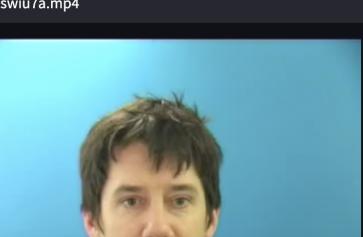
Developer Profiles

LIP SENSE

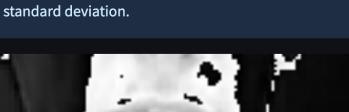
Try it out

Choose from a set of unseen videos that the model hasn't seen before

swiu7a.mp4



This is representative of what the machine learning model sees when making a prediction. The video has been preprocessed to have zero mean and unit standard deviation.



The screenshots illustrate the project's user interface and technology stack:

- Tech Stack Page:** Shows the project's dependencies: NumPy, OpenCV, Streamlit, IBM Watson, Matplotlib, and IBM Cloud imageio.
- Developer Profiles Page:** Displays the "Team LipBuddies" section, listing the following team members:
 - Chaitanya Tandon
 - Mohammad Sarim Shamim
 - Aditi Anand
 - Aashray Gupta

Made with Streamlit

7. Advantages and Disadvantages

7.1 Advantages

1. Enhanced Accessibility: The project significantly improves accessibility for individuals with speech impairments or in scenarios where vocalization is challenging.

2. Noise proof system: In noisy environments where traditional speech recognition struggles, this system remains effective. It is capable of interpreting lip movements even in acoustically challenging settings.
3. Privacy and Discretion: The system offers a discreet communication mode for situations where speaking aloud may disrupt the environment or raise privacy concerns, such as in libraries or public transportation.
4. Human-Machine Interaction: Integrating this technology into voice-controlled virtual assistants and technology interfaces enhances human-machine interaction, making it more intuitive and accessible.

7.2 Disadvantages

1. Hardware Requirements: Real-time processing and accuracy may demand substantial computational resources, potentially limiting the system's accessibility on low-end devices.
2. Lighting Dependence: Variations in lighting conditions can impact system performance. A consistent light source may be required for optimal accuracy.
4. Multilingual Training: Training the system for multilingual support is resource-intensive and time-consuming, requiring a diverse dataset for each language.
6. Accuracy Challenges: While the system can be highly accurate, challenges may arise in accurately interpreting subtle lip movements or dialect-specific variations.

8. Applications

The silent speech recognition system offers versatile applications, transcending speech limitations in various scenarios:

1. Noisy Environments: Thrives in bustling places where noise hinders speech recognition, ensuring effective communication.
2. Privacy and Security: Enables discreet communication in settings like libraries, preserving privacy.
3. Speech and Neurological Disorders: A lifeline for those with speech difficulties, fostering inclusivity.
4. AAC Enhancement: Enhances Augmentative and Alternative Communication (AAC) by providing natural communication via lip movements.
5. Human-Machine Interaction: Revolutionizes voice-controlled technology interaction, boosting accessibility.
6. Cross-Cultural Communication: Breaks language barriers, facilitating global communication.

7. Entertainment and Gaming: Enhances gaming experiences by responding to users' lip movements.

In summary, this system redefines communication across diverse contexts, offering inclusive, efficient, and engaging interactions

9. Conclusion

Our project employs a model to apply deep learning to end-to-end learning of a model that maps sequences of image frames of a speaker's mouth to entire sentences. The end-to-end model eliminates the need to segment videos into words before predicting a sentence. This model requires neither hand-engineered spatiotemporal visual features nor a separately-trained sequence model.

However, to extend the range of potential applications of LipNet, we aim to apply this approach to a jointly trained audio-visual speech recognition model, where visual input assists with robustness in noisy environments.

10. Future Scope

1. Individual Calibration: Tailoring the model to individual users can significantly enhance its accuracy and personalization. An adaptive calibration mechanism, accounting for variations in lip movements and speech patterns across individuals, promises to refine the system's performance.
2. Integration with Other Expressions: Expanding beyond speech recognition, integrating the system with facial expressions could enable the interpretation of a broader spectrum of communication cues.
3. Spelling Correction: Extending the model to offer spelling correction based on character-level predictions is a logical extension. Such a feature would not only enhance the accuracy of the recognized words but also contribute to the user's communication fluency.
4. Light Enhancement Module: Incorporating a module for light enhancement could mitigate challenges posed by varying lighting conditions. This enhancement would ensure consistent performance across diverse environments, further boosting the reliability of the system.
5. Multilingual Support: The project's potential can be broadened by incorporating multilingual support. By training the model on diverse languages, the silent speech recognition system can accommodate a global user base, transcending linguistic barriers.
6. Gesture Interpretation: Expanding the system's capabilities to include interpretation of hand gestures and body language could result in a comprehensive communication platform that

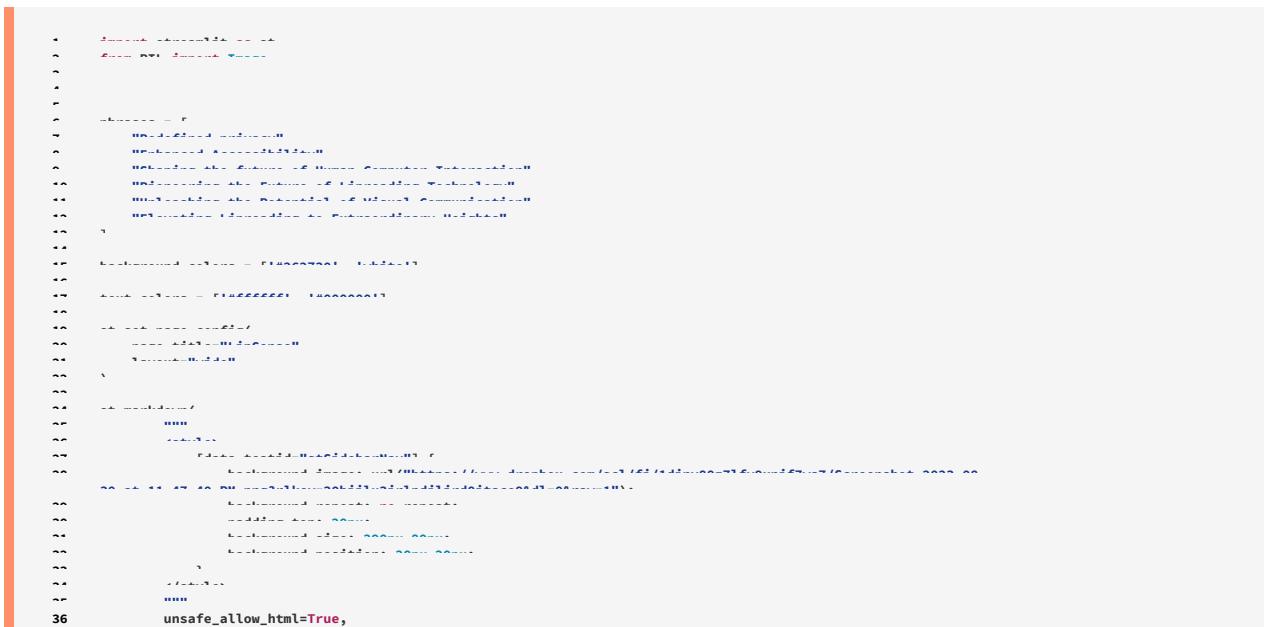
encompasses various non-verbal cues.

7. Continuous Learning: Implementing mechanisms for continuous learning can equip the system with the ability to adapt and evolve over time, ensuring that it remains relevant and effective in dynamic communication scenarios.
8. Mobile Application Integration: Developing a user-friendly mobile application can democratize access to the system, allowing users to engage with it seamlessly on their smartphones and tablets.

11. Bibliography

1. [LipNet: End-to-End Sentence-level Lipreading](<https://arxiv.org/abs/1611.01599>)
2. [CS230 - Deep Learning Course](<https://stanford.edu/~shervine/teaching/cs-230/>)
3. [LipType: A Silent Speech Recognizer Augmented with an Independent Repair Model](URL)
4. [YouTube Video](<https://youtu.be/qFJeN9V1Zsl>)
5. (<https://github.com/nicknochnack/LipNet/blob/main/LipNet.ipynb>)
6. [YouTubeVideo](https://www.youtube.com/watch?v=gZmobeGL0Yg&list=PLZbbT5o_s2xq7Lwl2y8_QtvuXZedL6tQU&ab_channel=deeplizard)
7. [YouTube Video](<https://www.youtube.com/watch?v=Zm8-qssPE58&t=698s>)
8. [YouTube Video](<https://youtu.be/TysuP3KgSz?si=-MMwMK3AmoVBvqtX>)

Appendix



The screenshot shows a Jupyter Notebook cell with Python code. The code defines a class `LipNet` with several methods: `__init__`, `createModel`, `train`, `predict`, and `predictBatch`. It uses TensorFlow's Keras API, including imports for `Sequential`, `TimeDistributed`, `Conv2D`, `MaxPooling2D`, `Flatten`, `Dense`, and `Lambda`. The `predict` method includes a parameter `unsafe_allow_html=True` at the bottom.

```
class LipNet(object):  
    def __init__(self, n_timesteps, n_input, n_output, n_filters=32, kernel_size=(3, 3), pool_size=(2, 2), hidden_units=128):  
        self.n_timesteps = n_timesteps  
        self.n_input = n_input  
        self.n_output = n_output  
        self.n_filters = n_filters  
        self.kernel_size = kernel_size  
        self.pool_size = pool_size  
        self.hidden_units = hidden_units  
  
        self.model = Sequential()  
        self.model.add(TimeDistributed(Conv2D(n_filters, kernel_size, activation='relu'),  
                                      input_shape=(n_timesteps, n_input, 1)))  
        self.model.add(TimeDistributed(MaxPooling2D(pool_size=pool_size)))  
        self.model.add(Flatten())  
        self.model.add(Dense(hidden_units, activation='relu'))  
        self.model.add(Dropout(0.2))  
        self.model.add(Dense(n_output, activation='softmax'))  
  
    def createModel(self):  
        self.model.compile(loss='categorical_crossentropy',  
                           optimizer='adam',  
                           metrics=['accuracy'])  
  
    def train(self, X_train, y_train, batch_size=16, epochs=10):  
        self.model.fit(X_train, y_train, batch_size=batch_size, epochs=epochs)  
  
    def predict(self, X):  
        return self.model.predict(X)  
  
    def predictBatch(self, X):  
        return self.model.predict(X)  
  
        unsafe_allow_html=True,
```

```
1 import streamlit as st
2 st.markdown(
3     """
4         <style>
5             [data-testid="stSidebarNav"] {
6                 background-image: url("https://www.dropbox.com/scl/fi/1diny80z7lfy9xnif7wa7/Screenshot_2023-08-
30_at_11.47.40_PM.png?rlkey=30biiiu3jrlpdj1ird0itaco9&dl=0&raw=1");
7                     background-repeat: no-repeat;
8                     padding-top: 20px;
9                     background-size: 290px 80px;
10                    background-position: 20px 20px;
11                }
12            </style>
13        """,
14        unsafe_allow_html=True,
15    )
16 #st.image('/Users/chaitanyatandon/Desktop/LIPBUDDIES/Screenshot_2023-08-30_at_11.47.40_PM.png')
17 st.header('Description', divider='blue')
18
19 st.write("Lip reading, the ability to understand spoken language by observing lip movements, is an essential skill for individuals with hearing impairments and in scenarios where audio information is compromised. Automatic lip reading systems, powered by deep learning techniques, have shown promising results in converting lip movements into text. This project proposes an advanced lip reading model that combines 3D Convolutional Neural Networks (CNNs) and Gated Recurrent Units (GRUs) to improve the accuracy and robustness of lip reading. The project's approach involves leveraging the spatio-temporal characteristics of lip motion using 3D CNNs. This allows the model to analyze sequential frames of lip movements and capture both spatial and temporal features simultaneously. The GRU component is employed to learn long-term dependencies and capture temporal dynamics in the lip movement sequences.")
1 import streamlit as st
2
3
4 st.markdown(
5     """
6         <style>
7             [data-testid="stSidebarNav"] {
8                 background-image: url("https://www.dropbox.com/scl/fi/1diny80z7lfy9xnif7wa7/Screenshot_2023-08-
30_at_11.47.40_PM.png?rlkey=30biiiu3jrlpdj1ird0itaco9&dl=0&raw=1");
9                     background-repeat: no-repeat;
10                    padding-top: 20px;
11                    background-size: 290px 80px;
12                    background-position: 20px 20px;
13                }
14            </style>
15        """,
16        unsafe_allow_html=True,
17    )
18 st.header('How it works?', divider='blue')
19
20 points = [
21     "The LipNet architecture operates through a meticulously designed sequence of operations.",
22     "It begins with a series of three stages, each involving spatiotemporal convolutions, channel-wise dropout, and spatial max-pooling.",
23     "This initial stage extracts crucial information from the input, encapsulating essential features.",
24     "The architecture introduces two Bidirectional Gated Recurrent Units (Bi-GRUs), a pivotal element in the process.",
25     "These Bi-GRUs facilitate the effective aggregation of outputs from the earlier spatiotemporal convolutions.",
26     "Following this, each time-step undergoes a linear transformation.",
27     "Subsequently, a softmax operation covers a vocabulary expanded to encompass the CTC blank symbol.",
28     "The architecture's learning process is driven by the application of the Connectionist Temporal Classification (CTC) loss.",
29     "Throughout the architecture, the Rectified Linear Unit (ReLU) serves as the activation function in all layers.",
30     "This consistent choice of activation function contributes to the model's efficiency and performance.",
31     "The orchestrated sequence of operations culminates in the LipNet architecture's proficiency in lip reading tasks."
32 ]
33
34 for point in points:
35     st.write(f"- {point}")
36
37 st.image('/Users/chaitanyatandon/Desktop/LIPBUDDIES/image (1).jpg')
38
1 import streamlit as st
2
3 st.header('Deployment made easy with IBM Cloud',divider="blue")
4 st.markdown(
5     """
6         <style>
7             [data-testid="stSidebarNav"] {
8                 background-image: url("https://www.dropbox.com/scl/fi/1diny80z7lfy9xnif7wa7/Screenshot_2023-08-
30_at_11.47.40_PM.png?rlkey=30biiiu3jrlpdj1ird0itaco9&dl=0&raw=1");
9                     background-repeat: no-repeat;
10                    padding-top: 20px;
11                    background-size: 290px 80px;
```

```

12         background-position: 20px 20px;
13     }
14 
```

```

35         "(None, 75, 256)", "(None, 75, 41)"],
36     "Param #": [3584, "0", "0", "884992", "0", "0", "518475", "0", "0", "0", "6660096", "0", "394240",
37     "0", "10537"]
38 }
39
40 df = pd.DataFrame(data)
41
42 st.table(df)
43
44 total_params = "8471924"
45 trainable_params = "8471924"
46 non_trainable_params = "0"
47
48 st.header("Model Parameters",divider='blue')
49 st.write("Total params:", total_params, "(" + str(int(total_params) / (1024 * 1024)) + ".2f" + "MB)")
50 st.write("Trainable params:", trainable_params, "(" + str(int(trainable_params) / (1024 * 1024)) + ".2f" + "MB)")
51 st.write("Non-trainable params:", non_trainable_params, "(" + str(int(non_trainable_params) / (1024 * 1024)) + ".2f" + "Byte")")
52
53 import streamlit as st
54 st.set_page_config(layout="wide")
55 st.markdown(
56     """
57     <style>
58     [data-testid="stSidebarNav"] {
59         background-image: url("https://www.dropbox.com/scl/fi/1diny80z7lfy9xnif7wa7/Screenshot_2023-08-
60     30_at_11.47.40_PM.png?rlkey=30bihilu3jrlpdjlird0itaco9&dl=0&raw=1");
61         background-repeat: no-repeat;
62         padding-top: 20px;
63         background-size: 290px 80px;
64         background-position: 20px 20px;
65     }
66     </style>
67     """,
68     unsafe_allow_html=True,
69 )
70
71 #st.image('/Users/chaitanyatandon/Desktop/LIPBUDDIES/Screenshot_2023-08-30_at_11.47.40_PM.png')
72 st.header('Performance Metrics', divider='blue')
73
74 col1, col2, col3 = st.columns(3)
75 col1.metric("Outperforms human by", "4.1x")
76 col2.metric("Word Error Rate", "4.8%")
77 col3.metric("Sentence accuracy", "95.2%")
78
79 col1, col2, col3 = st.columns(3)
80 col1.metric("Training Videos", "450")
81 col2.metric("Unseen CER", "6.4%")
82 col3.metric("Vocabulary tokens", "41")
83
84 col1, col2, col3 = st.columns(3)
85 col1.metric("Trainable parameters", "8471924")
86 col2.metric("STCNN Layers", "x3")
87 col3.metric("Bi-GRU", "x2")
88
89 st.image('/Users/chaitanyatandon/Desktop/LIPBUDDIES/ezgif.com-video-to-gif.gif',use_column_width="always")
90
91 import streamlit as st
92 from PIL import Image
93 st.image('/Users/chaitanyatandon/Desktop/LIPBUDDIES/Screenshot_2023-08-30_at_11.47.40_PM.png')
94 st.header('Try it out', divider='blue')
95
96
97
98 import streamlit as st
99 import os
100 import imageio
101 import time
102
103 import tensorflow as tf
104 from typing import List
105 import cv2
106
107
108 from tensorflow.keras.models import Sequential
109 from tensorflow.keras.layers import Conv3D, LSTM, Dense, Dropout, Bidirectional, MaxPool3D, Activation, Reshape, SpatialDropout3D,
110     BatchNormalization, TimeDistributed, Flatten
111
112 st.markdown(
113     """
114     <style>
115     [data-testid="stSidebarNav"] {
116         background-image: url("https://www.dropbox.com/scl/fi/1diny80z7lfy9xnif7wa7/Screenshot_2023-08-
117     30_at_11.47.40_PM.png?rlkey=30bihilu3jrlpdjlird0itaco9&dl=0&raw=1");
118         background-repeat: no-repeat;
119         padding-top: 20px;
120         background-size: 290px 80px;

```

```

29         background-position: 20px 20px;
30     }
31 
```

`</style>`

```

32     """,
33     unsafe_allow_html=True,
34 )
35
36
37
38 vocab = [x for x in "abcdefghijklmnopqrstuvwxyz?!123456789 "]
39 char_to_num = tf.keras.layers.StringLookup(vocabulary=vocab, oov_token="")
40 # Mapping integers back to original characters
41 num_to_char = tf.keras.layers.StringLookup(
42     vocabulary=char_to_num.get_vocabulary(), oov_token="", invert=True
43 )
44
45 def load_model() -> Sequential:
46     model = Sequential()
47
48     model.add(Conv3D(128, 3, input_shape=(75,46,140,1), padding='same'))
49     model.add(Activation('relu'))
50     model.add(MaxPool3D((1,2,2)))
51
52     model.add(Conv3D(256, 3, padding='same'))
53     model.add(Activation('relu'))
54     model.add(MaxPool3D((1,2,2)))
55
56     model.add(Conv3D(75, 3, padding='same'))
57     model.add(Activation('relu'))
58     model.add(MaxPool3D((1,2,2)))
59
60     model.add(TimeDistributed(Flatten()))
61
62     model.add(Bidirectional(LSTM(128, kernel_initializer='Orthogonal', return_sequences=True)))
63     model.add(Dropout(.5))
64
65     model.add(Bidirectional(LSTM(128, kernel_initializer='Orthogonal', return_sequences=True)))
66     model.add(Dropout(.5))
67
68     model.add(Dense(41, kernel_initializer='he_normal', activation='softmax'))
69
70     model.load_weights(os.path.join('..','models','checkpoint'))
71
72     return model
73
74 def load_video(path:str) -> List[float]:
75     #print(path)
76     cap = cv2.VideoCapture(path)
77     frames = []
78     for _ in range(int(cap.get(cv2.CAP_PROP_FRAME_COUNT))):
79         ret, frame = cap.read()
80         frame = tf.image.rgb_to_grayscale(frame)
81         frames.append(frame[190:236,80:220,:])
82     cap.release()
83     mean = tf.math.reduce_mean(frames)
84     std = tf.math.reduce_std(tf.cast(frames, tf.float32))
85     return tf.cast((frames - mean), tf.float32) / std
86
87 def load_alignments(path:str) -> List[str]:
88     #print(path)
89     with open(path, 'r') as f:
90         lines = f.readlines()
91         tokens = []
92         for line in lines:
93             line = line.split()
94             if line[2] != 'sil':
95                 tokens = [*tokens, ' ',line[2]]
96     return char_to_num(tf.reshape(tf.strings.unicode_split(tokens, input_encoding='UTF-8'), (-1)))[1:]
97
98 def load_data(path: str):
99     path = bytes.decode(path.numpy())
100    file_name = path.split('/')[-1].split('.')[0]
101    video_path = os.path.join('..','data','s1',f'{file_name}.mp4')
102    alignment_path = os.path.join('..','data','alignments','s1',f'{file_name}.align')
103    frames = load_video(video_path)
104    alignments = load_alignments(alignment_path)
105    return frames, alignments
106
107
108
109 # Generating a list of options or videos
110 options = os.listdir(os.path.join('..', 'data', 's1'))
111 selected_video = st.selectbox('Choose from a set of unseen videos that the model hasn\'t seen before', options)
112

```

```

113 # Generate two columns
114 col1, col2 = st.columns(2)
115
116 if options:
117
118     # Rendering the video
119     with col1:
120         file_path = os.path.join('..','data','s1', selected_video)
121         video = open(file_path, 'rb')
122         video_bytes = video.read()
123         st.video(video_bytes)
124         #st.info('The video that you have selected')
125
126     with col2:
127         st.info('This is representative of what the machine learning model sees when making a prediction. The video has been preprocessed to have zero mean and unit standard deviation.')
128         video, annotations = load_data(tf.convert_to_tensor(file_path))
129         st.image('/Users/chaitanyatandon/Desktop/LIPBUDDIES/animation.gif',use_column_width="always")
130
131     model = load_model()
132     yhat = model.predict(tf.expand_dims(video, axis=0))
133     with st.spinner('Predicting'):
134         time.sleep(5)
135         st.success('Success')
136     decoder = tf.keras.backend.ctc_decode(yhat, [75], greedy=True)[0].numpy()
137
138     # Convert prediction to text
139     st.header('Decoded raw tokens as words',divider='blue')
140     converted_prediction = tf.strings.reduce_join(num_to_char(decoder)).numpy().decode('utf-8')
141     st.info(converted_prediction)
142
143     st.header('This is the output of the machine learning model as tokens',divider='blue')
144     st.text(decoder)
145
146 import streamlit as st
147 st.markdown(
148     """
149     <style>
150         [data-testid="stSidebarNav"] {
151             background-image: url("https://www.dropbox.com/scl/fi/1diny80z7lfy9xnif7wa7/Screenshot_2023-08-30_at_11.47.40_PM.png?rlkey=30biilu3jrlpdjlird0itaco9&dl=0&raw=1");
152             background-repeat: no-repeat;
153             padding-top: 20px;
154             background-size: 290px 80px;
155             background-position: 20px 20px;
156         }
157     </style>
158     """,
159     unsafe_allow_html=True,
160 )
161 st.image('/Users/chaitanyatandon/Desktop/LIPBUDDIES/Screenshot 2023-08-31 at 2.22.51 AM.png')
162 st.header('Tech Stack',divider='blue')
163
164 st.image('/Users/chaitanyatandon/Desktop/LIPBUDDIES/tech stack.png')
165 import streamlit as st
166 st.markdown(
167     """
168     <style>
169         [data-testid="stSidebarNav"] {
170             background-image: url("https://www.dropbox.com/scl/fi/1diny80z7lfy9xnif7wa7/Screenshot_2023-08-30_at_11.47.40_PM.png?rlkey=30biilu3jrlpdjlird0itaco9&dl=0&raw=1");
171             background-repeat: no-repeat;
172             padding-top: 20px;
173             background-size: 290px 80px;
174             background-position: 20px 20px;
175         }
176     </style>
177     """,
178     unsafe_allow_html=True,
179 )
180 #st.image('/Users/chaitanyatandon/Desktop/LIPBUDDIES/Screenshot 2023-08-31 at 2.22.51 AM.png')
181 st.header('Developer Profiles',divider='blue')
182 st.title('Team LipBuddies')
183 st.write('Chaitanya Tandon')
184 st.write('Mohammad Sarim Shamim')
185 st.write('Aditi Anand')
186 st.write('Aashray Gupta')
187
188 !pip3 install tensorflow
189 !pip3 list
190 import tensorflow as tf
191 from ibm_watson_machine_learning import APIClient
192
193 wml_credentials={

```

```

7      "url":"https://us-south.ml.cloud.ibm.com",
8      "apikey":nF2Bu5KRYc9z2ZnMhNuzVcghozZaznLHue1h5h0crcVk"
9  }
10 client=APIClient(wml_credentials)
11 client
12 client.spaces.list()
13 space_uid=f9a08d7a-1730-4d9f-b753-50c531422a76"
14 client.set_default_space(space_uid)
15 client.software_specifications.list()
16 software_space_uid=client.software_specifications.get_uid_by_name("runtime-22.2-py3.10")
17 software_space_uid
18 from tensorflow.keras.models import Sequential
19 from tensorflow.keras.layers import Conv3D, LSTM, Dense, Dropout, Bidirectional, MaxPool3D, Activation, Reshape, SpatialDropout3D,
20 BatchNormalization, TimeDistributed, Flatten
21 from tensorflow.keras.optimizers import Adam
22 from tensorflow.keras.callbacks import ModelCheckpoint, LearningRateScheduler
23 model = Sequential()
24 model.add(Conv3D(128, 3, input_shape=(75,46,140,1), padding='same'))
25 model.add(Activation('relu'))
26 model.add(MaxPool3D((1,2,2)))
27
28 model.add(Conv3D(256, 3, padding='same'))
29 model.add(Activation('relu'))
30 model.add(MaxPool3D((1,2,2)))
31
32 model.add(Conv3D(75, 3, padding='same'))
33 model.add(Activation('relu'))
34 model.add(MaxPool3D((1,2,2)))
35
36 model.add(TimeDistributed(Flatten()))
37
38 model.add(Bidirectional(LSTM(128, kernel_initializer='Orthogonal', return_sequences=True)))
39 model.add(Dropout(.5))
40
41 model.add(Bidirectional(LSTM(128, kernel_initializer='Orthogonal', return_sequences=True)))
42 model.add(Dropout(.5))
43
44 model.add(Dense(char_to_num.vocabulary_size()+1, kernel_initializer='he_normal', activation='softmax'))
45 vocab = [x for x in "abcdefghijklmnopqrstuvwxyz?!123456789 "]
46 char_to_num = tf.keras.layers.StringLookup(vocabulary=vocab, oov_token="")
47 num_to_char = tf.keras.layers.StringLookup(
48     vocabulary=char_to_num.get_vocabulary(), oov_token="", invert=True
49 )
50 model.summary()
51 model.save("lipbud.h5")
52 !tar -zcvf lipbud.tgz lipbud.h5
53 model_details=client.repository.store_model(model="lipbud.tgz",meta_props={
54     client.repository.ModelMetaNames.NAME:"lipbud.h5",
55     client.repository.ModelMetaNames.TYPE:"tensorflow_2.9",
56     client.repository.ModelMetaNames.SOFTWARE_SPEC_UID:software_space_uid
57 })
58 model_id=client.repository.get_model_id(model_details)
59 model_id
60 client.repository.download(model_id,'lipbud3.tgz')
61 model_id=client.repository.get_model_id(model_details)
62 model_id
63 model.load_weights('models/checkpoint')
64 model.save("lipbud.h5")
65 url = 'https://drive.google.com/uc?id=1vWscXs4Vt0a_1IH1-ct2TCgXAZT-N3_Y'
66 output = 'checkpoints.zip'
67 gdown.download(url, output, quiet=False)
68 gdown.extractall('checkpoints.zip', 'models')
69 !pip install opencv-python matplotlib imageio gdown tensorflow
70 import os
71
72 import tensorflow as tf
73 import numpy as np
74 from typing import List
75 from matplotlib import pyplot as plt
76 !pip3 install opencv-python matplotlib imageio gdown tensorflow
77 import os
78 import cv2
79 import tensorflow as tf
80 import numpy as np
81 from typing import List
82 from matplotlib import pyplot as plt
83 import imageio
84 import gdown
85 model.load_weights('models/checkpoint')
86 model.save("lipbud.h5")
87 model.summary()
88 !tar -zcvf lipbud.tgz lipbud.h5
89 model_details=client.repository.store_model(model="lipbud.tgz",meta_props={
```

```

90     client.repository.ModelMetaNames.NAME:"lipbud.h5",
91     client.repository.ModelMetaNames.TYPE:"tensorflow_2.9",
92     client.repository.ModelMetaNames.SOFTWARE_SPEC_UID:software_space_uid
93   ))
94 model_id=client.repository.get_model_id(model_details)
95 model_id
96 client.repository.download(model_id,'lipbud2.tgz')
97 import numpy as np
98 import os
99 from tensorflow.keras.models import load_model
100 from tensorflow.keras.preprocessing import image
101 from flask import Flask , request, render_template
102 !pip3 install flask
103 def load_video(path:str) -> List[float]:
104
105   cap = cv2.VideoCapture(path)
106   frames = []
107   for _ in range(int(cap.get(cv2.CAP_PROP_FRAME_COUNT))):
108     ret, frame = cap.read()
109     frame = tf.image.rgb_to_grayscale(frame)
110     frames.append(frame[190:236,80:220,:])
111   cap.release()
112   mean = tf.math.reduce_mean(frames)
113   std = tf.math.reduce_std(tf.cast(frames, tf.float32))
114   return tf.cast((frames - mean, tf.float32) / std
115 vocab = [ for x in "abcdefghijklmnopqrstuvwxyz?!123456789 "]
116 char_to_num = tf.keras.layers.StringLookup(vocabulary=vocab, oov_token="")
117 num_to_char = tf.keras.layers.StringLookup(
118   vocabulary=char_to_num.get_vocabulary(), oov_token="", invert=True
119 )
120
121 print(
122   f"The vocabulary is: {char_to_num.get_vocabulary()}"
123   f"(size ={char_to_num.vocabulary_size()})"
124 )
125 char_to_num.get_vocabulary()
126 def load_alignments(path:str) -> List[str]:
127   with open(path, 'r') as f:
128     lines = f.readlines()
129   tokens = []
130   for line in lines:
131     line = line.split()
132     if line[2] != 'sil':
133       tokens = [*tokens, ' ',line[2]]
134   return char_to_num(tf.reshape(tf.strings.unicode_split(tokens, input_encoding='UTF-8'), (-1)))[1:]
135 def load_data(path: str):
136   path = bytes.decode(path.numpy())
137   file_name = path.split('/')[-1].split('.')[0]
138   # File name splitting for windows
139   #file_name = path.split('\\')[-1].split('.')[0]
140   video_path = os.path.join('data','alignments','sil',f'{file_name}.mpg')
141   alignment_path = os.path.join('data','alignments','sil',f'{file_name}.align')
142   frames = load_video(video_path)
143   alignments = load_alignments(alignment_path)
144   return frames, alignments
145 test_path = '/Users/chaitanyatandon/Desktop/LIPBUDDIES/data/s1/bbal6n.mpg'
146 tf.convert_to_tensor(test_path).numpy().decode('utf-8').split('/')[1].split('.')[0]
147 frames, alignments = load_data(tf.convert_to_tensor(test_path))
148 plt.imshow(frames[40])
149 alignments
150 tf.strings.reduce_join([bytes.decode(x) for x in num_to_char(alignments.numpy())])
151 def mappable_function(path:str) ->List[str]:
152   result = tf.py_function(load_data, [path], (tf.float32, tf.int64))
153   return result
154 from matplotlib import pyplot as plt
155 sample = load_data(tf.convert_to_tensor('/Users/chaitanyatandon/Desktop/LIPBUDDIES/data/s1/bras9a.mpg'))
156 print('*'*100, 'REAL TEXT')
157 [tf.strings.reduce_join([num_to_char(word) for word in sentence]) for sentence in [sample[1]]]
158
159 yhat = model.predict(tf.expand_dims(sample[0], axis=0))
160 decoded = tf.keras.backend.ctc_decode(yhat, input_length=[75], greedy=True)[0][0].numpy()
161 print('*'*100, 'PREDICTIONS')
162 [tf.strings.reduce_join([num_to_char(word) for word in sentence]) for sentence in decoded]
163
164 model.save("lipbud.h5")
165
```