

Text of silence

Title:

**Text of silence is for generating text from lip movement for
aphonia patient, Hearing-impaired people and in noisy
environments**

Under Supervisor

Dr. Mohammed Abdel Hameed

Students

1. Mamdoh Marof Abdulrazak
2. Mohammed Awadallah Bakir
3. Abdelrahman Mohammed Abdullah
4. Mahmoud Mohammed Fouad
5. Amr Ftouh Abbas

ABSTRACT

Lip reading is a visual way of “listening” to someone. This is done by looking at the speaker's face to follow their speech patterns in order to recognize what is being said. it is used to understand or interpret speech without hearing it, a technique especially mastered by people with hearing difficulties. The ability to lip-read enables a person with a hearing impairment to communicate with others and to engage in social activities, which otherwise would be difficult. Recent advances in the fields of computer vision, pattern recognition, and signal processing have led to a growing interest in automating this challenging task of lip reading. Indeed, automating the human ability to lip-read, a process referred to as visual speech recognition, could open the door for other novel applications.

it is also known as audio-visual recognition, has been considered as a solution for speech recognition tasks, especially when the audio is corrupted or when the conversation happened in noisy environments. It can also be an extremely helpful tool for people who are hearing-impaired to communicate through video calls. This task, however, is challenging, due to factors such as the variances in the inputs (facial features, skin colors, speaking speeds, etc.) and the one-to-many relationships between viseme and phoneme.

Lip-reading technology mainly includes face detection, lip localization, feature extraction, training the classifier through the corpus, and finally recognition of the word/sentence through lip movement.

Contents List

Chapter 1 Introduction	6
1.1 Abbreviation	7
1.2 Introduction	8-11
1.3 The Essential Question	12
1.4 Motivation and Justification	12
1.5 Related work	13-15
Chapter 2 Domain Analysis and Techniques	16
2.1 Domain Analysis	17-18
2.2 Techniques	19-22
Chapter 3 Proposed System and Methodology	23
3.1 System Use-Cases	24-25
3.2 Use Case Description (Use case scenario)	26-31
3.3 Analysis Class	32
3.3.1 State Diagram	32
3.4 Interaction Diagram (Sequence Diagram)	33
3.5 System Architecture	34-35
Chapter 4 Risk and Functional and Non-Functional Requirements	36
4.1 Problems/Constraints	37-38
4.2 Project Plan	39
4.3 Quality Assurance Plan	40
4.4 Requirements	41
4.4.1 Functional Requirements	41
4.4.2 Non-Functional Requirements	41
4.5 System Request	42-43

4.6 Project Key Objectives.....	44
4.7 Datasets.....	45-48
Chapter 5 System design	49
5.1 User interface	50-53
5.2 Samples.....	54-56
References.....	57-58



List of Figures

Figure 1. Convolutional Neural Networks.....	19
Figure 2. Recurrent Neural Network.....	20
Figure 3. Artificial Neural Network.....	21
Figure 4. Application Use case.....	24
Figure 5. Server Use case.....	25
Figure 6. User State diagram.....	32
Figure 7. User Sequence diagram.....	33
Figure 8. System Architecture.....	34
Figure 9. System Architecture Block Diagram.....	35
Figure 10. Project Plan.....	39
Figure 11. Datasets and Accuracy.....	45
Figure 12. Accuracy with Different Datasets	46
Figure 13. TOS in application drawer.....	50
Figure 14. Home Screen.....	51
Figure 15. Sentence Structure Screen.....	52
Figure 16. Capture Video.....	53
Figure 18. Lip region after segmentation.....	56

Chapter 1

Introduction



TOS
TEXT OF SILENCE

1.1 Abbreviations

Keyword	Meaning
AI	Artificial intelligence
ML	Machine learning
DL	Deep learning
ANN	Artificial neural network
RNN	Recurrent neural network
CNN	Convolutional neural network
ADAM	adaptive moment estimation

Table 1. Abbreviation

TOS
TEXT OF SILENCE

1.2 Introduction

To know how our project works we must first discuss what Artificial intelligence (AI) is and its types.

Artificial intelligence (AI):

AI is a field of computer science that studies how machines can imitate the intelligence of their human counterparts. Over the last decade, definitions of the term have become quite loose and refer to just about any computerized or automated function. However, the difference between an AI system and traditional software packages is the ability to make informed judgments and decisions by responding to patterns in data.

Applications of Artificial intelligence:

1. Smart homes, cities and infrastructure
2. Artificial intelligence against Covid-19
3. Machine translations
4. Transport
5. Health

Machine learning (ML):

Machine learning (ML) is a subset of artificial intelligence, which build a mathematical model based on sample data, known as “training data,” in order to make predictions or decisions without being explicitly programmed to perform the task. In machine learning, neural networks, support vector machines, and evolutionary computation, we are usually given a training set and a test set. By the training set, it will mean the union of the labeled set and the unlabeled set of examples available to machine learners. In comparison, test set consists of examples never seen before.

Deep learning (DL):

Deep Learning is an emerging field of Machine learning; that is, it is a subset of Machine Learning where learning happens from past examples or experiences with the help of 'Artificial Neural Networks'. Deep Learning uses deep neural networks, where the word 'deep' signifies the presence of more than 1 or 2 hidden layers apart from the input and output layer.

Deep learning algorithms:

1. Artificial neural network (ANN):

A series of algorithms that are trying to mimic the human brain and find the relationship between the sets of data.

2. Recurrent neural network (RNN):

Recurrent neural networks are designed to interpret temporal or sequential information. These networks use other data points in a sequence to make better predictions. They do this by taking in input and reusing the activations of previous nodes or later nodes in the sequence to influence the output. RNN works on the principle of saving the output of a particular layer and feeding this back to the input in order to predict the output of the layer.

3. Convolutional neural network (CNN):

A class of deep neural networks that extracts features from images, given as input, to perform specific tasks such as image classification, face recognition and semantic image system CNN has one or more convolution layers for simple feature extraction, which execute convolution operation (i.e. multiplication of a set of weights with input) while retaining the critical features (spatial and temporal information) without human supervision.

4. Adam optimizer:

The Adam optimization algorithm is an extension to stochastic gradient descent that has recently seen broader adoption for deep learning applications in computer vision and natural language processing.

And is an optimization algorithm that can be used instead of the classical stochastic gradient descent procedure to update network weights iterative based in training data.

What is lip reading?

The use of lip-reading has been documented since the 16th century and hearing-impaired people often use lip-reading as an adjunct to understanding fluent speech. When it comes to automating the process, there are many challenges compared to conventional audio recognition.¹

Automated lip-reading system refers to the systems which utilizes the visual information of the movement of the speech articulators such as the lips, teeth and somehow tongue of the speaker. The advantages are that such a system is not sensitive to ambient noise and change in acoustic conditions, does not require the user to make a sound, and provides the user with a natural feel of speech and dexterity of the mouth.

lip-reading plays a crucial role in human communication and speech understanding, as highlighted by the McGurk effect (McGurk & MacDonald, 1976)², where one phoneme's audio dubbed on top of a video of someone speaking a different phoneme results in a third phoneme being perceived. Lip-reading is a notoriously difficult task for humans, especially in the absence of context¹. Most lip-reading actuations, besides the lips and sometimes tongue and teeth, are latent and difficult to disambiguate without context (Fisher, 1968; Woodward & Barber, 1960)³. For example, Fisher (1968) gives 5 categories of visual phonemes (called visemes), out

of a list of 23 initial consonant phonemes, that are commonly confused by people when viewing a speaker's mouth. Many of these were asymmetrically confused, and observations were similar for final consonant phonemes. Consequently, human lip-reading performance is poor. Hearing-impaired people achieve an accuracy of only $17 \pm 12\%$ even for a limited subset of 30 monosyllabic words and $21 \pm 11\%$ for 30 compound words (Easton & Basala, 1982). An important goal, therefore, is to automate lip-reading. Machine lip readers have enormous practical potential, with applications in improved hearing aids, silent dictation in public spaces, security, and speech recognition in noisy environments, biometric identification, and silent-movie processing. Machine lip-reading is difficult because it requires extracting spatiotemporal features from the video (since both position and motion are important). Recent deep learning approaches attempt to extract those features end-to-end.⁴

Lip reading the ability to recognize what is being said from visual information alone, is an impressive skill, and very challenging for a novice. It is inherently ambiguous at the word level due to homophones – different characters that produce exactly the same lip sequence (e.g. 'p' and 'b'). However, such ambiguities can be resolved to an extent using the context of neighboring words in a sentence, and/or a language model. A machine that can lip read opens up a host of applications: 'dictating' instructions or messages to a phone in a noisy environment; transcribing and re-dubbing archival silent films; resolving multi-talker simultaneous speech; and, improving the performance of automated speech recognition in general. That such automation is now possible is due to two developments that are well known across computer vision tasks: the use of deep neural network models; and, the availability of a large-scale dataset for training. In this case the model is based on the recent sequence-to sequence (encoder-decoder with attention) translator

architectures that have been developed for speech recognition and machine translation.⁵

1.3 The Essential Question:

The essential question with relevance to the Vision and Mission of the Faculty of Computers and information at Luxor University would be: How can we, using our education, help the hearing-impaired and aphonia people nationally, regionally, and internationally to better serve individuals, society, and the environment.

1.4 Motivation and Justification:

We are encouraged to work on this project as it has the potential to help over ten million English-speaking deaf and hearing-impaired people all over the world 6 Based on our research, currently there is no product that could help solve the problem of lip reading at a reasonable cost, there is currently a solution being developed and is planned for release within the next year by a company called Liopa in the UK.

Seeing the lack of research in the field motivated us to explore how we can utilize our knowledge to implement a solution that could contribute to aiding the hearing-impaired.

The project is also relevant to our interests as it mainly involves two fields of computer engineering like Computer Vision, Neural Networks and Machine

Learning in general and Neural Networks in specific

1.5 Related work

Research on lip reading (a.k.a. visual speech recognition) has a long history. Many of the existing works in this field have followed similar pipelines which first extract spatiotemporal features around the lips (either motion-based, geometric-feature based or both), and then align these features with respect to a canonical template.

Title	Author	Year	abstract
Improved speaker independent lip-reading using speaker adaptive training and deep neural networks	<ul style="list-style-type: none">Ibrahim Almajai	2016	Furthermore, we show that error rates can be even further reduced by the additional use of Deep Neural Networks (DNN). We also find that there is no need to map phonemes to visemes for context-dependent visual speech transcription. ⁷
Out of Time: Automated Lip Sync in the Wild	<ul style="list-style-type: none">Joon Son ChungAndrew Zisserman	2017	They apply the network to two further tasks: active speaker detection and lip-reading. On both tasks, we set a new state-of-the-art on standard benchmark datasets. ⁸
Lip Reading Sentences in the Wild	<ul style="list-style-type: none">Joon Son ChungAndrew SeniorOriol VinyalsAndrew Zisserman	2017	The goal of this work is to recognize phrases and sentences being spoken by a talking face, with or without the audio. Unlike previous works that have focused on recognizing a limited number of words or phrases, we tackle lip reading as an open-world problem - unconstrained natural language sentences, and in the wild videos. Our key contributions are:

			<p>1.a Watch, Listen, Attend and Spell' (WLAS) network that learns to transcribe videos of mouth motion to characters;</p> <p>2.a curriculum learning strategy to accelerate training and to reduce overfitting;</p> <p>3.a 'Lip Reading Sentences' (LRS) dataset for visual speech recognition⁹</p>
Morse code application	<ul style="list-style-type: none"> • Morse Samsung 	2019	where blind people could use their mobile phones with Morse code taping instead of ordinary texting. ¹⁰
Combining Residual Networks with LSTMs for Lipreading	<ul style="list-style-type: none"> • Themos Stafylakis, • Georgios Tzimiropoulos 	2017	We propose an end-to-end deep learning architecture for word-level visual speech recognition. The system is a combination of spatiotemporal convolutional, residual and bidirectional Long Short-Term Memory networks. We train and evaluate it on the Lipreading In-The-Wild benchmark. ¹¹
Learning to lip read words by watching videos	<ul style="list-style-type: none"> • Joon Son Chung • Andrew Zisserman. 	2018	Our aim is to recognize the words being spoken by a talking face, given only the video but not the audio. Existing works in this area have focused on trying to recognize a small number of utterances in controlled environments (e.g. digits and alphabets), partially due to the shortage of suitable datasets. ¹²
Deep Learning of Mouth Shapes for Sign Language	<ul style="list-style-type: none"> • Oscar Koller • Hermann Ney • Richard Bowden 	2015	This paper deals with robust modelling of mouth shapes in the context of sign language recognition using deep convolutional neural networks. Sign language mouth shapes are difficult to annotate and thus hardly any

			publicly available annotations exist. As such, this work exploits related information sources as weak supervision. Humans mainly look at the face during sign language communication, where mouth shapes play an important role and constitute natural patterns with large variability.
--	--	--	---



Chapter 2

Domain Analysis and Techniques



2.1 Domain Analysis

Seeing the lack of research in the field motivated us to explore how we can utilize our knowledge to implement a solution that could contribute to aiding the hearing-impaired.

The project is also relevant to our interests as it mainly involves two fields of computer engineering like Computer Vision, Neural Networks and Machine Learning in general and Neural Networks in specific.

Text of Silence is a lip-reading mobile application that aims to help the hearing-impaired communicate and interact better with their surroundings.

It consists of a lip-extracting module-using image processing, a learning module using machine learning.

Description of Products and Services: The chief goal is to build a system that captures human face, detects the mouth position and traces its movements regarding lips positions and movements in order to predict the words said or willing to be said by the user.

Mobile camera captures the user's video, passing it to a back-end where all the lips tracing will be processed by an already trained module to detect the desired word or group of words.

The actual text then propagates to the user's application interface.

Many new initiatives and users can rely on this project.

Hearing-impaired users can rely on this application to communicate with those who cannot understand sign language.

Moreover, some organizations can use this functionality in order to detect words out of videos with no sound like football games.

Technology Consideration: Considering the technological issues that might happen during the live-phase of the prototype testing and subsequently commercial operation, problems may appear with the mobile camera quality regarding resolution, night sight mode, and similar environmental conditions.

Therefore, it might be prudent to operate on devices with medium to high camera capabilities.

Another issue that might arise is the internet speed as well as back-end server speed; since the internet speed in some countries is limited below, certain boundaries thus the application version operated in low-speed internet countries will have the most restricted video time.

this project will facilitate the communication between people as well as other field's requirements like detecting any unethical words used in a sports event.



2.2 Techniques

1. Convolutional Neural Networks (CNN)

A CNN is one of the variants of neural networks used heavily in the field of Computer Vision. It derives its name from the type of hidden layers it consists of. The hidden layers of a CNN typically consist of convolutional layers, pooling layers, fully connected layers, and normalization layers. Here it simply means that instead of using the normal activation functions, convolution and pooling functions are used as activation functions ¹³. Figure 1 illustrates a Convolutional Neural Networks.

We use it to extract features from images that we infer from video

“used to transform images of the lips region to its vector representation.”

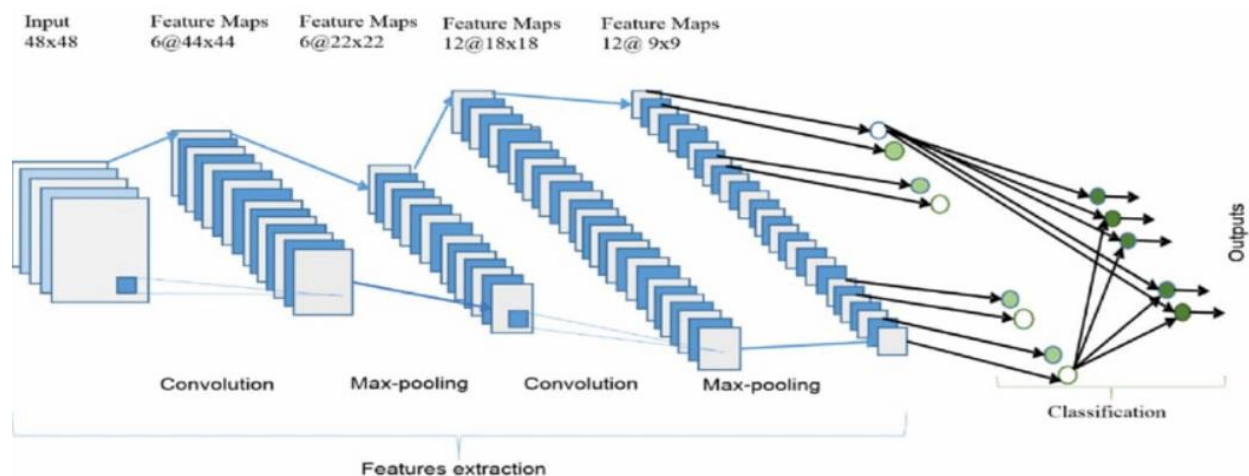


Figure 1. Convolutional Neural Networks

2. Recurrent Neural Network (RNN)

Humans do not start their thinking from scratch every second. As you read this report, you understand each word based on your understanding of previous words.

You do not throw everything away and start thinking from scratch again. Your thoughts have persistence. Traditional artificial neural networks cannot do this, and it seems like a major shortcoming. For example, imagine you want to classify what kind of event is happening at every point in a movie. It is unclear how a traditional neural network could use its reasoning about previous events in the film to inform later ones.

Recurrent neural networks address this issue by allowing for a feedback between layers. They are networks that allow information to persist, where the output from previous step are fed as input to the current step⁸. Figure 2 shows a Recurrent Neural Network and how it's 'Recurrent' nature.

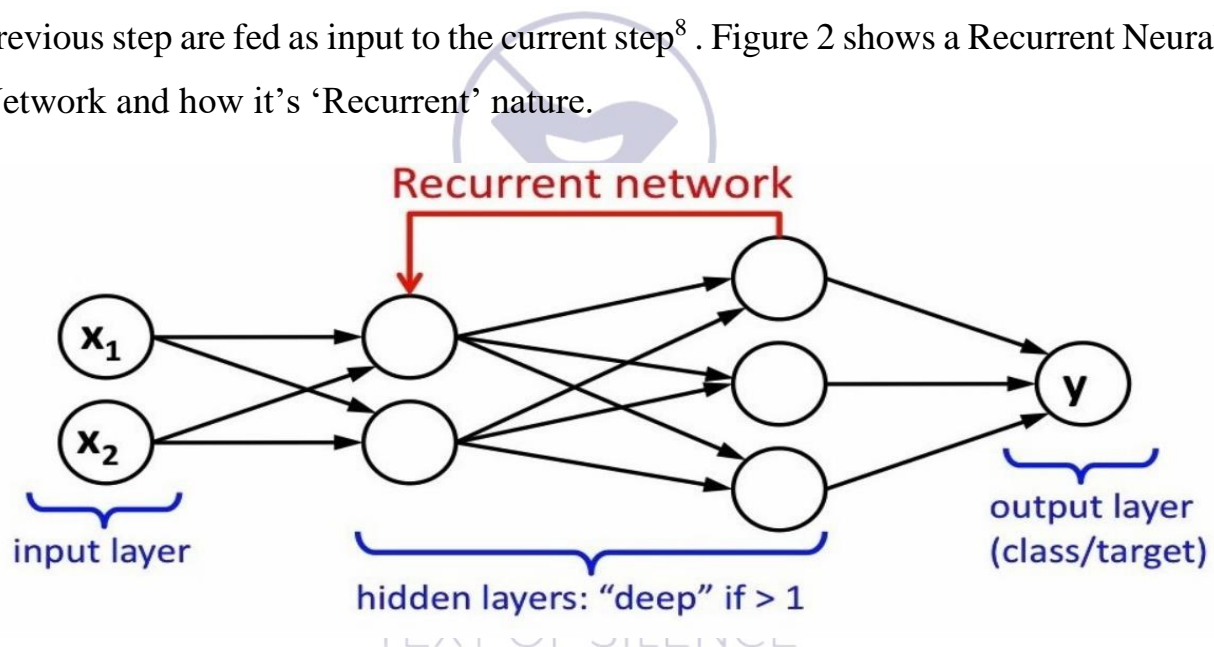


Figure 2. Recurrent Neural Network

3. Artificial Neural Networks (ANN)

They are one of the main tools used in machine learning. As the “neural” part of their name suggests, they are brain-inspired systems, which are intended to replicate the way that we humans learn. Neural networks consist of input and output layers, as well as hidden layers consisting of units that transform the input into something that the output layer can use. They are excellent tools for finding patterns which are

far too complex or numerous for a human programmer to extract and teach the machine to recognize ¹⁴.

While neural networks (also called “perceptron”) have been around since the 1940s, it is only in the last several decades where they have become a major part of artificial intelligence. This is due to the arrival of a technique called “backpropagation,” which allows networks to adjust their hidden layers of neurons in situations where the outcome does not match what the creator is hoping for — like a network designed to recognize dogs, which misidentifies a cat, for example. Figure 3 shows an example of an Artificial Neural Network showing its input, hidden, and output layers where the input is selected based on the number of input features while the output is based on the number of required classes to classify

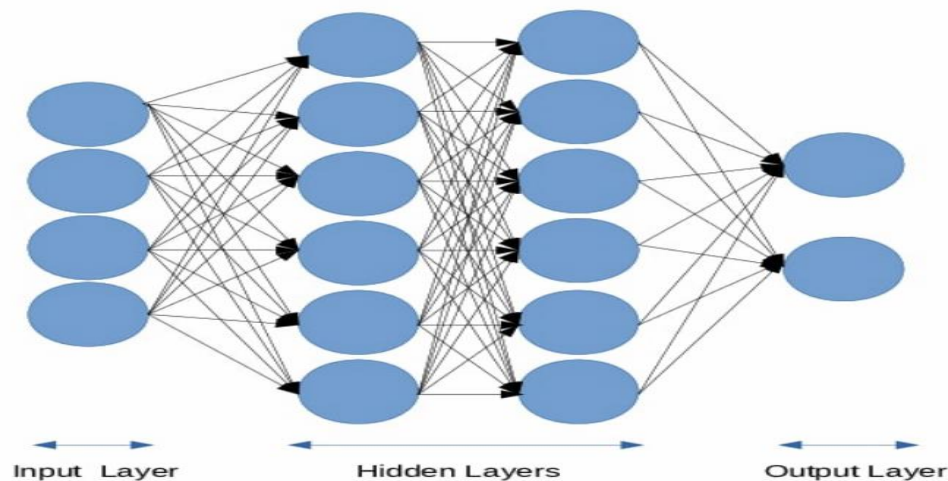


Figure 3. Artificial Neural Network

4. Adam Optimizer

Adaptive Moment Estimation (ADAM) is a method that calculates the individual adaptive learning rate for each parameter from estimates of first and second moments of the gradients. It can be viewed as a combination of Adagrad, which works well on sparse gradients and RMSprop, which works well in online and nonstationary settings. Adam implements the exponential moving average of the gradients to scale the learning rate instead of a simple average as in Adagrad. It keeps an exponentially decaying average of past gradients while staying computationally efficient and having very little memory requirement. Adam optimizer is one of the most popular gradient descent optimization algorithms.



Chapter 3

Proposed System and Methodology



TOS
TEXT OF SILENCE

1.1 System Use-Cases

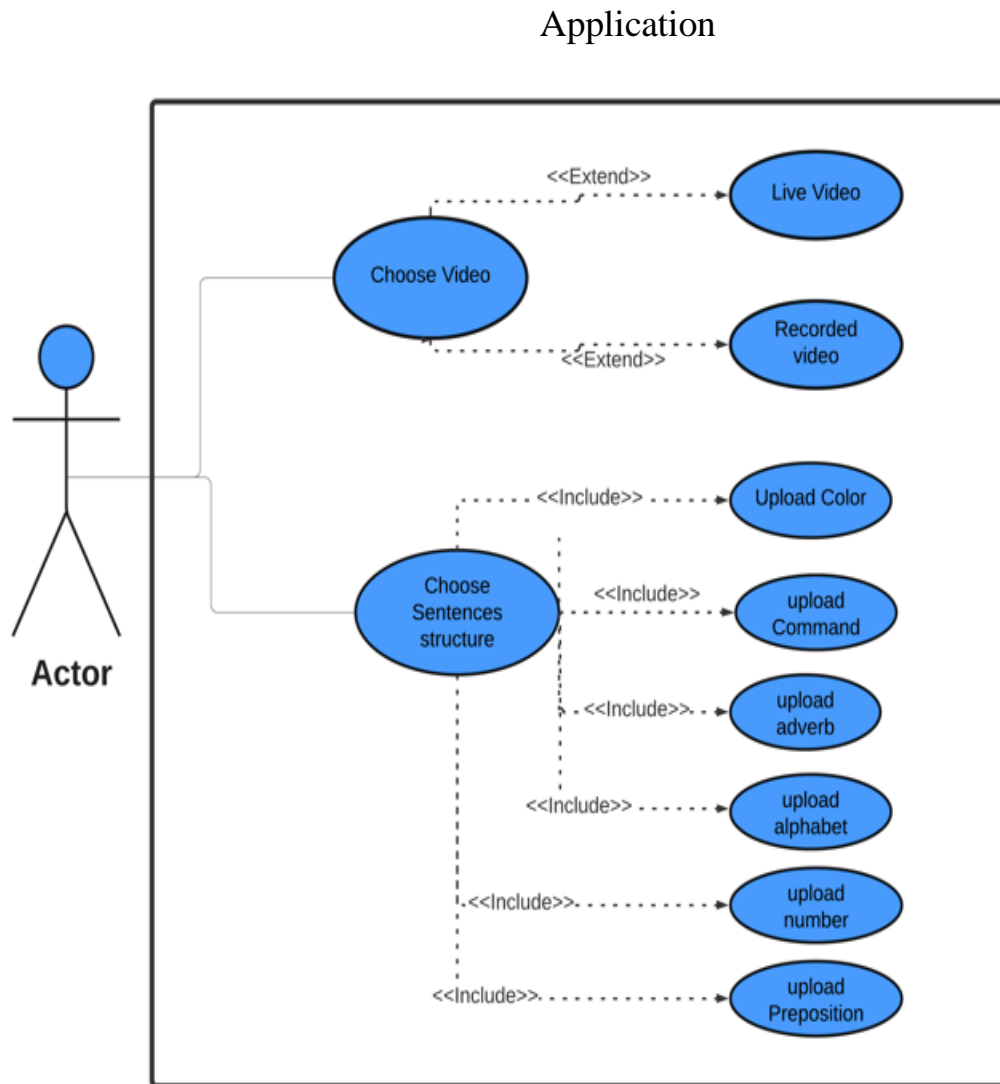


Figure 4. Application Use case

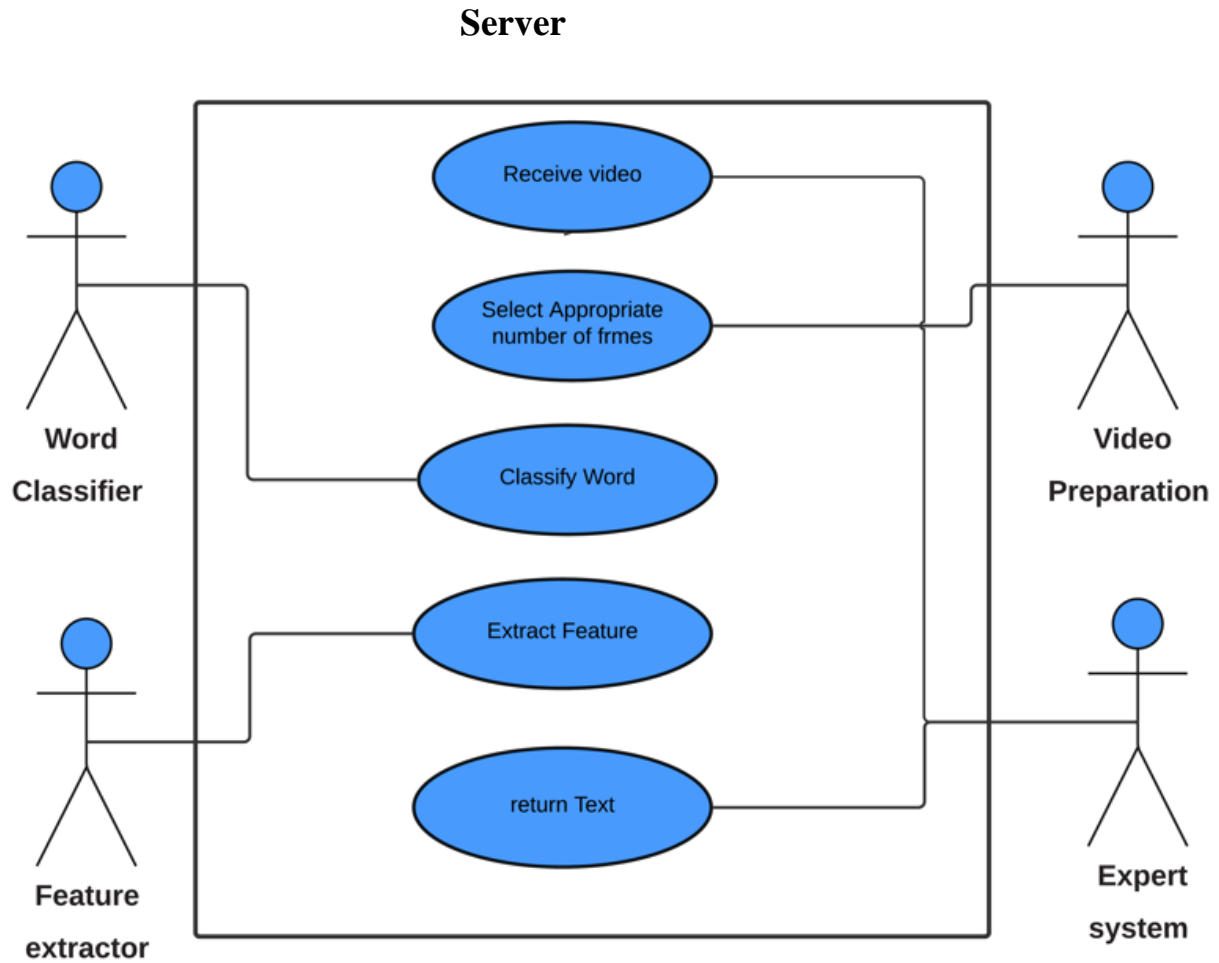


Figure 5. Server Use case

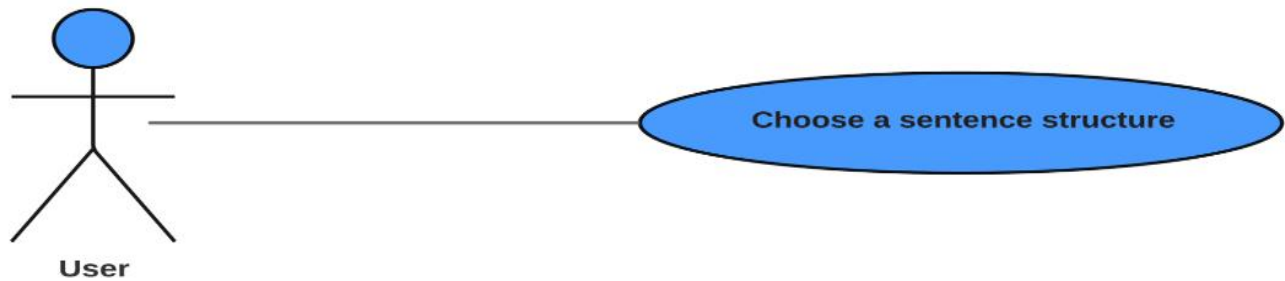
TEXT OF SILENCE

1.2 Use Case Description (Use case scenario)

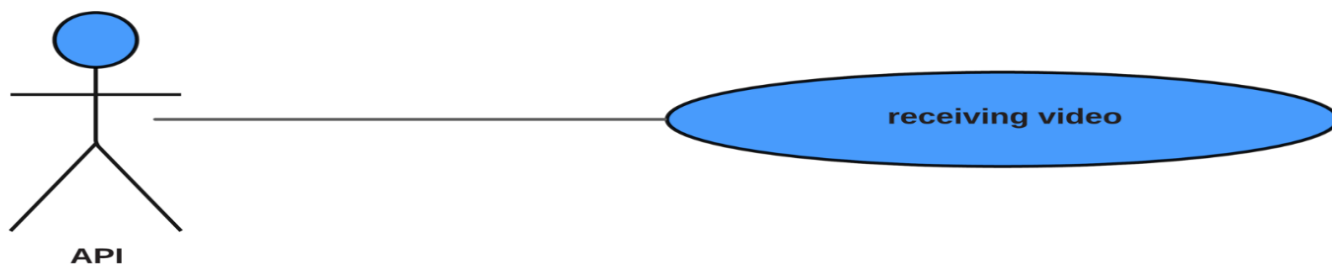


WordClassifier

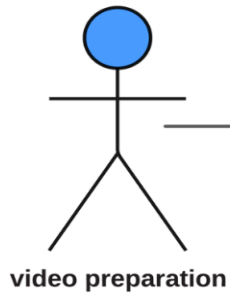
Use case name	Choose a video	
Unique ID	user 001	
Area	Application	
Actor(s)	user	
Description	User adds a video to translate lip movement into words	
Triggering Event	The user clicks the video selection button	
Preconditions	The user needs to download the app The user needs access to the Internet Choose the sentence structure to represent in a video	
Postconditions	The user uploaded the video to the server successfully Waiting for the expected word	
Assumptions	The user predicted the words he made	
Steps Performed		Information for Steps
Open the application and the Internet Choose the sentence structure you want to translate Click Add Video Choose whether you want to record a live video or a video stored in your phone		Step 3: Choose what you want to translate (alphabet - color - command - number -.....)
Extensions (Alternative Flows)	If you don't choose anything and he can't add videos, you will get a warning message	



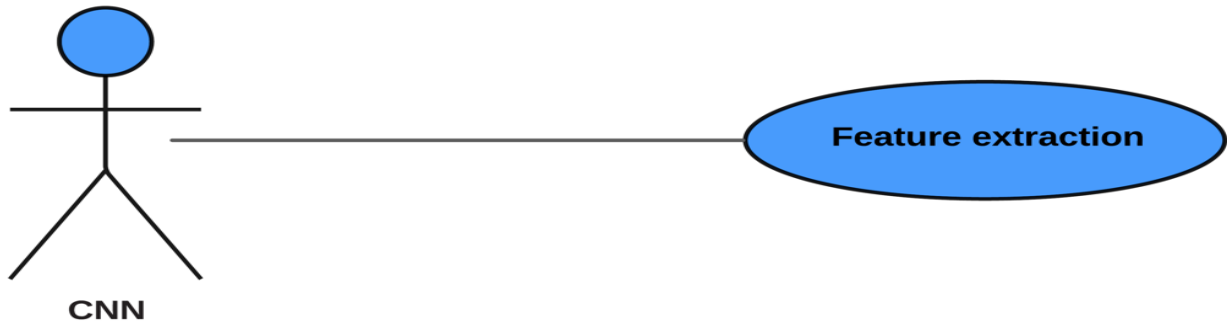
Use case name	Choose a sentence structure	
Unique ID	User 002	
Area	Application	
Actor(s)	user	
Description	Choose the appropriate types of words you want to translate	
Triggering Event	Choose sentence structure by clicking on the checkboxes	
Preconditions	The user needs to download the application The user needs access to the Internet	
Postconditions	User can add videos	
Assumptions	User can record videos for uploading	
Steps Performed	Information for Steps	
Open the app Choosing the right types of words Add videos	Step 3: (alphabet – color -command – number -)	
Extensions (Alternative Flows)	If he doesn't choose any things and can't adding videos, he should get a warning message	



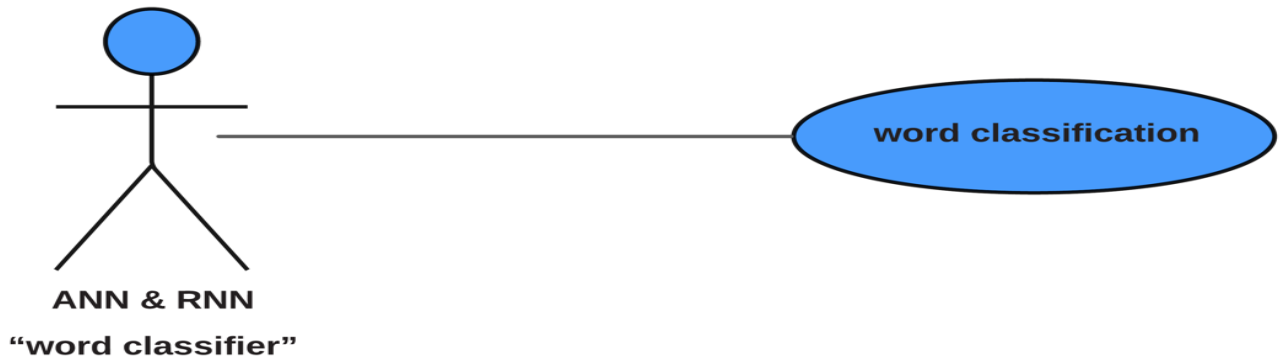
Use case name	receiving video
Unique ID	server 001
Area	server
Actor(s)	API
Description	Receiving video from the user after recording it
Triggering Event	Receive video from API
Preconditions	No preconditions
Postconditions	successful receiving
Assumptions	User can upload videos to server
Steps Performed	Information for Steps
Recording videos Get videos from API Send videos to prepare the video	Step 2: From the user
Extensions (Alternative Flows)	If he doesn't choose any things and can't adding videos, he should get a warning message



Use case name	Select the appropriate number of frames	
Unique ID	server 002	
Area	server	
Actor(s)	video preparation	
Description	Prepare the appropriate number of frames from the video	
Triggering Event	Converts the video preparation to the appropriate number of frames and sends it to the feature extraction process	
Preconditions	The server should succeed in receiving the videos	
Postconditions	Returns the appropriate number of frames from the received video	
Assumptions	The video can enter the feature extraction process	
Steps Performed		Information for Steps
Receive video from API Select the appropriate number of frames from the videos Send the frames to the second step		Step 3: Feature extraction
Extensions (Alternative Flows)	If he doesn't choose any things and can't adding videos, he should get a warning message	



Use case name	Feature extraction	
Unique ID	Server 003	
Area	Server	
Actor(s)	CNN	
Description	Getting the vector of features from the frames	
Triggering Event	Receiving frames from video preparation process and returning feature vector	
Preconditions	The video must be converted to the appropriate number of frames	
Postconditions	Return Vector of Features	
Assumptions	The vector can enter the following process "word classification"	
Steps Performed		Information for Steps
Receive the appropriate number of frames Create a vector of feature Return Vector of Features		Step 1: Video Preparation
Extensions (Alternative Flows)	If he doesn't choose any things and can't adding videos, he should get a warning message	



Use case name	word classification	
Unique ID	server 004	
Area	server	
Actor(s)	ANN & RNN “word classifier”	
Description	Get the correct word that fits the feature vector	
Triggering Event	ANN & RNN classify the word	
Preconditions	The feature vector must be valid	
Postconditions	classified the word	
Assumptions	The predicted word can return to the API	
Steps Performed		Information for Steps
Receiving Vector of Features classified the word return classified the word		Step 1: CNN Step 2: ANN & RNN
Extensions (Alternative Flows)	If he doesn't choose any things and can't adding videos, he should get a warning message	

1.3 Analysis Class

4.3.1 State Diagram

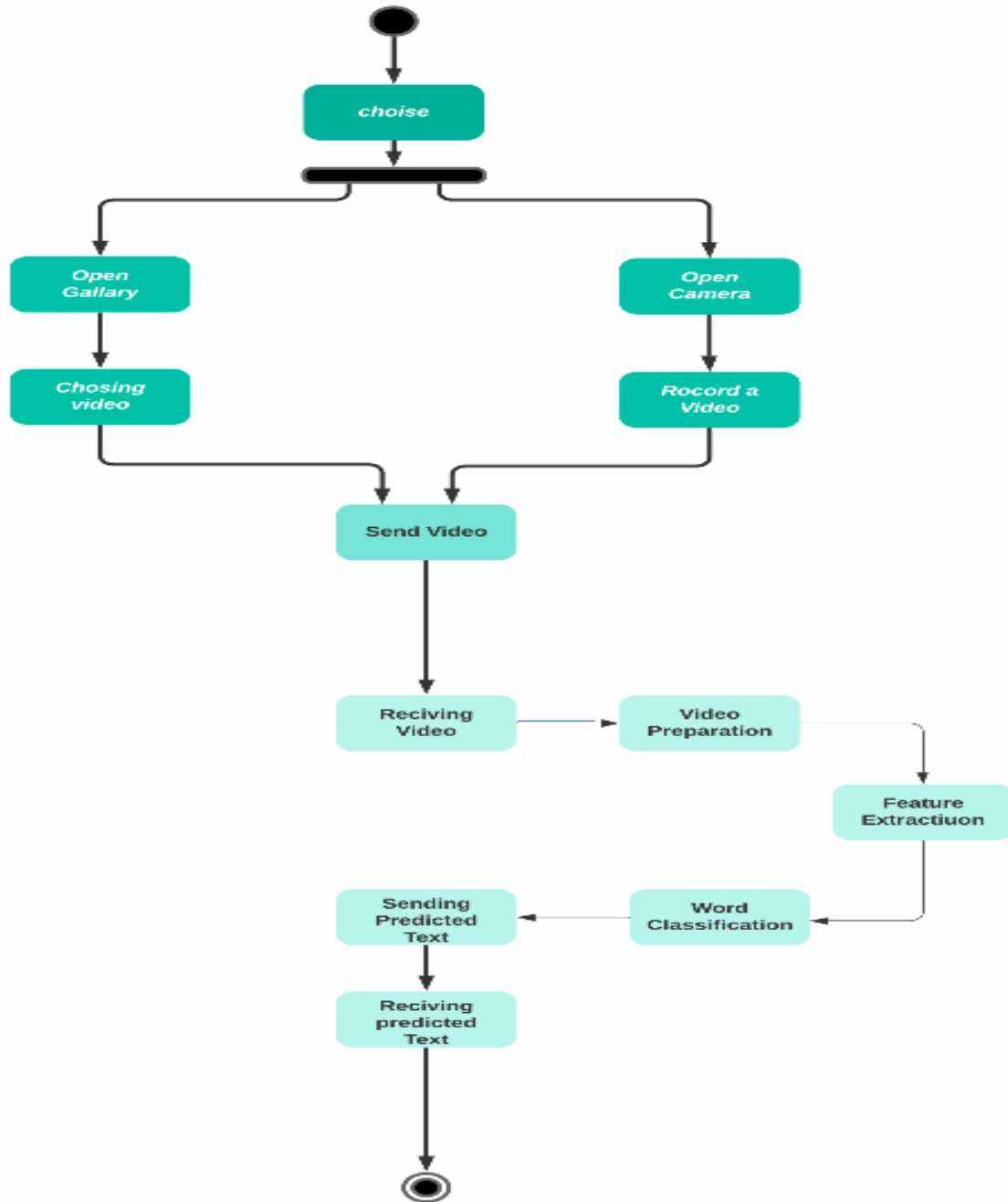


Figure 6. User State diagram

1.4 Interaction Diagram (Sequence Diagram)

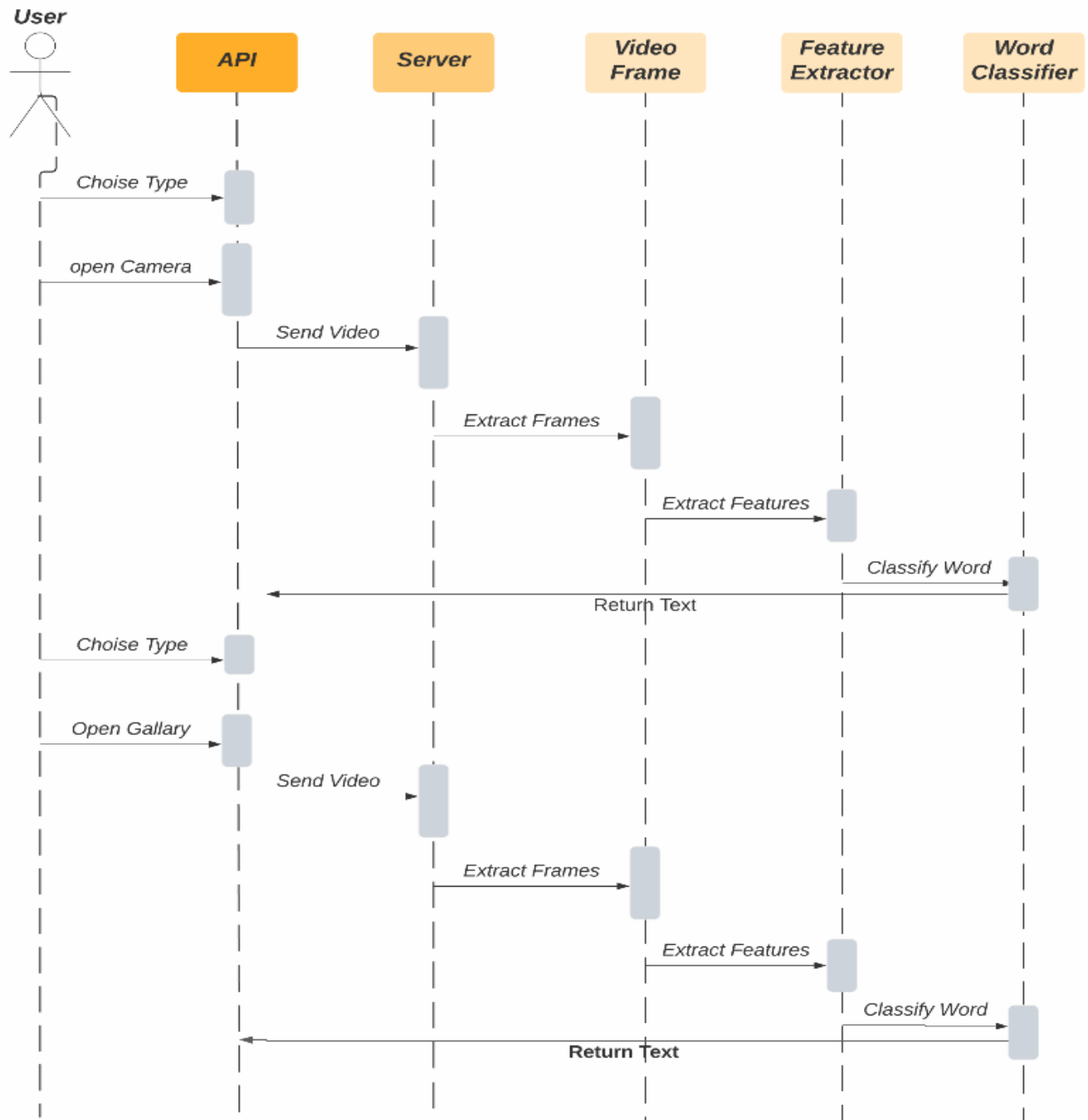


Figure 7. User Sequence diagram

1.4.1 System Architecture

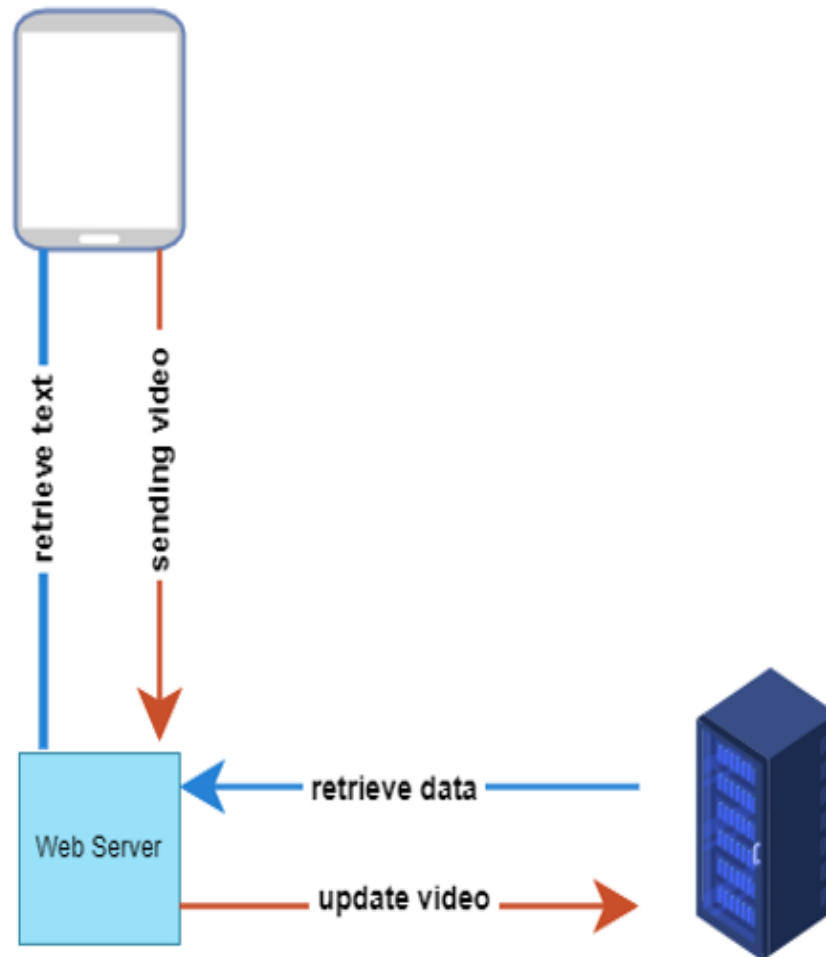


Figure 8. System Architecture

System Architecture Block Diagram.

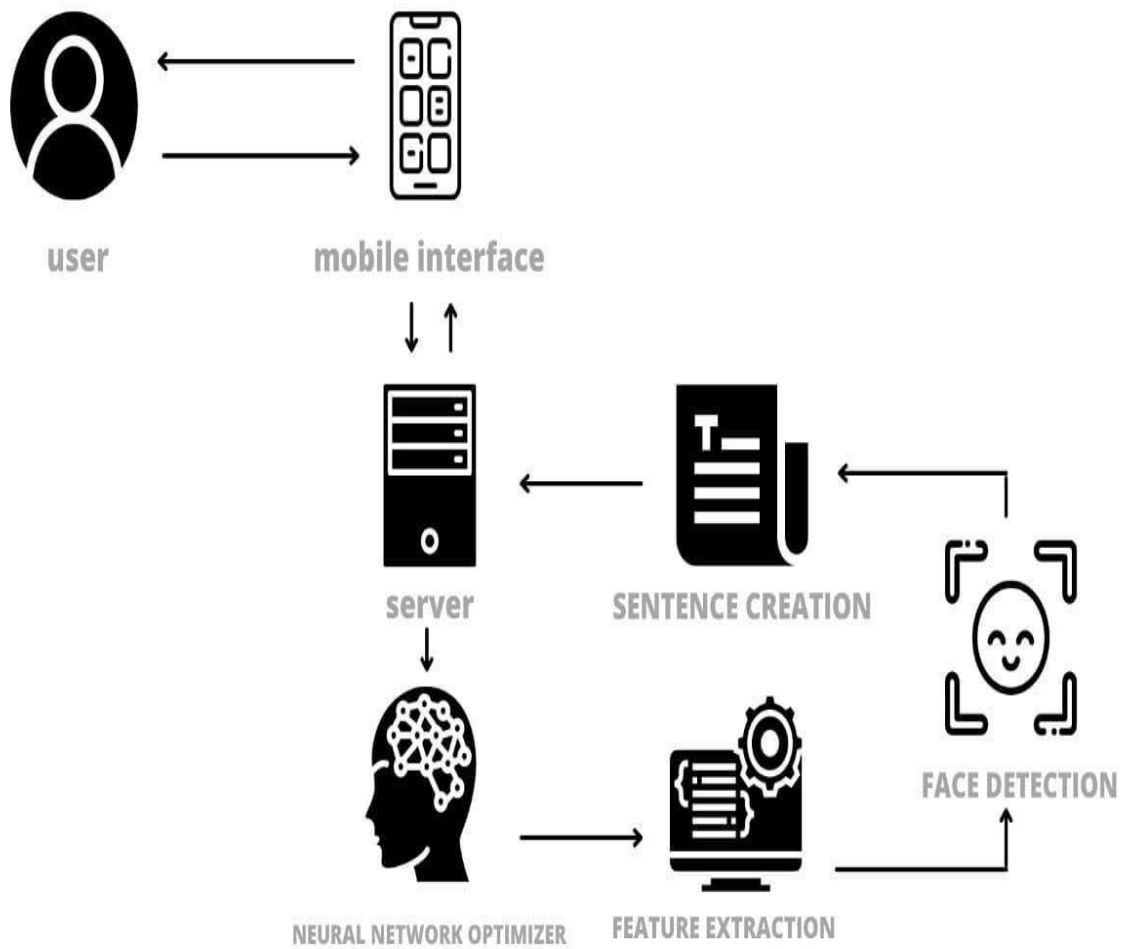


Figure 9. System Architecture Block Diagram.

Chapter 4

Risk and Functional and Non-Functional Requirement



TOS
TEXT OF SILENCE

1.1 Problem Statement /Constraint

3.1.1 Problem Statement

1. Different countries have their own sign language standards, and even the differences between British Sign Language and American Sign Language are significant.
2. Speech recognition in noisy environments (e.g. cars):

In noisy environments, it's difficult to hear the sound of video hence you can generate speak from lip movements by our project.

3. Silent dictation in public spaces:

for example, in reading places like libraries or any public place you can't talk with sound so you can use our project.

- Therefore, our aim to build a smart and intelligent mobile application that is able to see the user's lips and transform their words into text.

3.1.2 Constraint

Constraint classification	Constraint	Effects
Event classification	1-light may be low and brightness may be very high	1-2-3: May leads to inaccurate predicted text.

	2-Video quality may be very low	4- If internet is low so it will take more time and May leads to incomplete process.
	4-position of the speaker	
	3-Internet connection speed may be slow	
Ambiguity constrains	1-similar letters	1-It may leads to inaccurate predicted text.
	2-more than one person in video	2-It may leads to unknown where speaker is?
Implementation constrains	1-Lack of data	1- Lack of enough data leads to cannot find all the words we want it to expect.

3.2 Project Plan

3.2.1 Project Plan

phase	Oct	Nov	Dec	Jan	Feb	Mar	Apr	May	Jun
Gathering Information									
Define Requirements									
analysis									
design									
implementation									
Develop AI System									
Testing and Final Discussion									

Figure 10. Project Plan

3.3 Quality Assurance Plan:

- **Black Box:**

In this stage, we will use test inputs in our dataset in our system to ensure the accuracy of output.

- **White box:**

- **Unity Testing:**

In this stage of testing, we will take every component of our AI TEXT OF SILENCE system such as web service, CNN machine learning model, mobile application to test them separately.

- **Validation Testing:**

Validation testing is the process of ensuring if the tested and developed application satisfies its functionality requirements. The business requirement logic or scenarios have to be tested in detail. All the critical functionalities of an application must be tested here.

- **Alpha:**

In this part, a group of testers in our team test the product in a laboratory environment to ensure the efficiency of the product and fix errors.

- **Beta:**

In this stage of testing, the application has been sent to some users that can't speak or have a hearing impairment to test the system efficiency and its outputs are correct or not and retrieve feedback to our team

3.4 Requirements

3.4.1 Functional Requirements

1. A functional mobile app that controls workflow with servers, APIs, users, in organized, accurate, and quick methods.
2. A Server contains data and Machine Learning Models.
3. API to guarantee communication between server and application.
4. A functional and tested Machine Learning Model that analyzes lip movement to predict output text.
5. A system that specifies output text and generates an audio clip that helps user to communicate with other users.
6. The Application must forward the recorded video to the model for analysis and prediction text through an API
7. Web service to be used as a mobile app.

3.4.2 Non-Functional Requirements

1. Fast APIs and retrieve data from servers within seconds.
2. The speed of video analysis and the speed of word prediction from the movement of the lips within seconds
3. A clear, attractive, responsive, less cluttered and efficient user interface that ensures a good user experience and handles a variety of tasks.
4. The application works on most versions of Android and IOS
5. A well-organized and managed system architecture that facilitates maintenance, restoration and scalability.
6. The system correctly predicts most words from videos

3.5 System Request

This kind of problem is a very good challenge for AI/machine learning because it expands capabilities a machine to describe a human's innate and highly innate ability. In addition, there is a fair amount of data that can be used to train a machine to do this.

This project focuses on lip reading, lip reading is not a new science. Lip reading has been around in the local community as long as the language has been around. However, learning takes many years and can be difficult especially for people who have lost their hearing later in life or who have difficulty making sounds. Lip reading has deep roots in the intelligence community. There are entire audio and video recordings where the voice is too muted or too noisy to understand speech. Including multiple camera face directions¹⁵.

1. Project Sponsor

- All people who suffer from hearing problems that happened to them in later life, as well as people who have speech problems, that is, they cannot get the sound out of their mouth.
- And also, people interested in improving the quality of video calls.

2. Business need

- It will facilitate communication between people who have lost their hearing in later life.
- Also, people who have trouble getting the sound out of the mouth.
- And assistance in communicating in video calls in noisy environments and in cases of sound interruption.

3. Business Requirements

- The application must be located on a mobile phone or device that can access the web pages
- It must be connected to the internet
- And that the phone has a camera to capture videos in which the movement of the lips will be translated
- And shoot the video in a place with good lighting so that the filming is done properly.

4. Business value

- Facilitating the process of communication between people who have a hearing problem or people who have problems removing sound from the mouth due to surgery or psychological problems.
- Improving future video calls, which currently depend on audio output only, but with the lip-reading technique, audio output will depend on audio and video outputs together, which leads to a significant improvement in the quality of video calls.

5. Special issue

- The main problem that we will face is the quality of video capture and the surrounding lighting in the shooting environment.
- And also the confusion that occurs when similar letters are pronounced in the movement of the mouth
- There is also a major problem in implementing the project, which is the lack of enough data to be able to train the model so that it can find all the words we want it to expect.

3.6 Project Key Objectives

1. To take advantage of the current technological advance to make the lives of the hearing-impaired and who have problems in speak better.
2. It also aims at facilitating the everyday activities these people perform with their loved ones.
3. One of the major problems facing the hearing-impaired is communication with people who do not know sign language and may misunderstand them and to enhance voice when the voice is absent or corrupted by external noise. Or who suffers from a loss of ability to speak as a result of exposure to any surgical or other reason. Sign language is far from universal and leads to lots of misunderstandings; it is not a universal dialect.



3.7 Datasets:

Lip-reading datasets (AVICar, AVLetters, AVLetters2, BBC TV, CUAVE, OuluVS1, OuluVS2), but most only contain single words or are too small. One exception is the GRID corpus, which has audio and video recordings of 34 speakers who produced 1000 sentences each, for a total of 28 hours across 34000 sentences.

data set	output	accuracy
AVICAR	Digits	37.9%
AVLetter	Alphabet	64.6%
CUAVE	Digits	83.0%
OuluVS	Phrases	91.4%
OuluVS2	Phrases	94.1%
BBC TV	Words	65.4%
GRID	Words	86.4%
GRID	Sentences	95.2%

Figure 11. Datasets and Accuracy

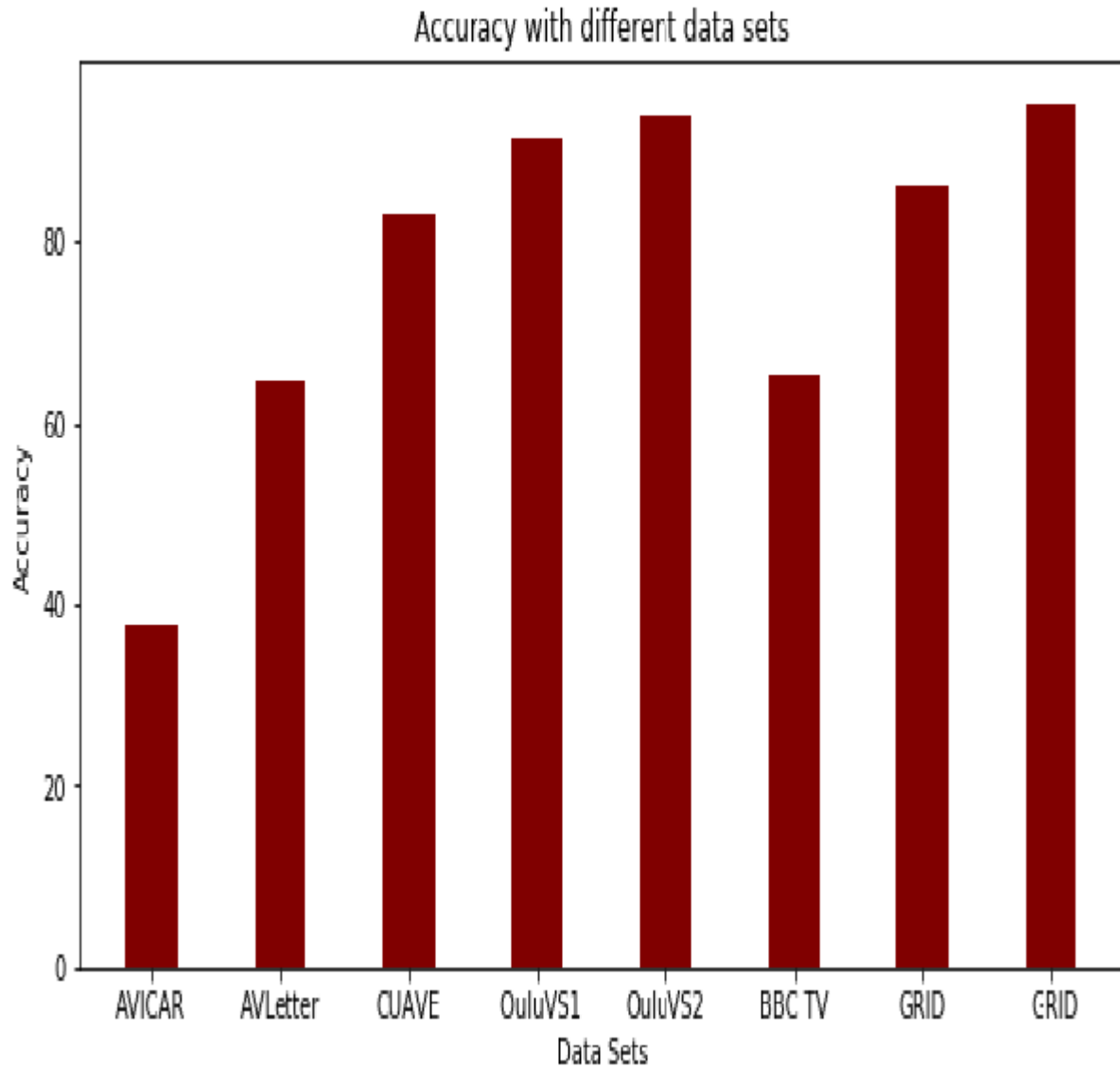


Figure 12. Accuracy with Different Datasets

We use the GRID corpus to evaluate our project because it is sentence-level and has the most data. The sentences are drawn from the following simple grammar: command + color + preposition + letter+ digit + adverb, where the number denotes how many word choices there are for each of the 6 word categories. The categories consist of, respectively, {bin, lay, place, set}, {blue, green, red, white}, {at, by, in, with}, {A, . . . , Z}\{W}, {zero, . . . , nine}, and {again, now, please, soon}, yielding

64000 possible sentences. For example, two sentences in the data are “set blue by a four please” and “place red at C zero again”.

Examples of Dataset:





TOS
TEXT OF SILENCE

Chapter 5

System Design



5.1 User Interface

TOS's user guide on how to install, run, and use the application. The user manual is designed to be as easy and simple and possible to enable all types of users to easily access the application. After downloading and installing the application, the application will show in your app drawer as shown in Figure 13.

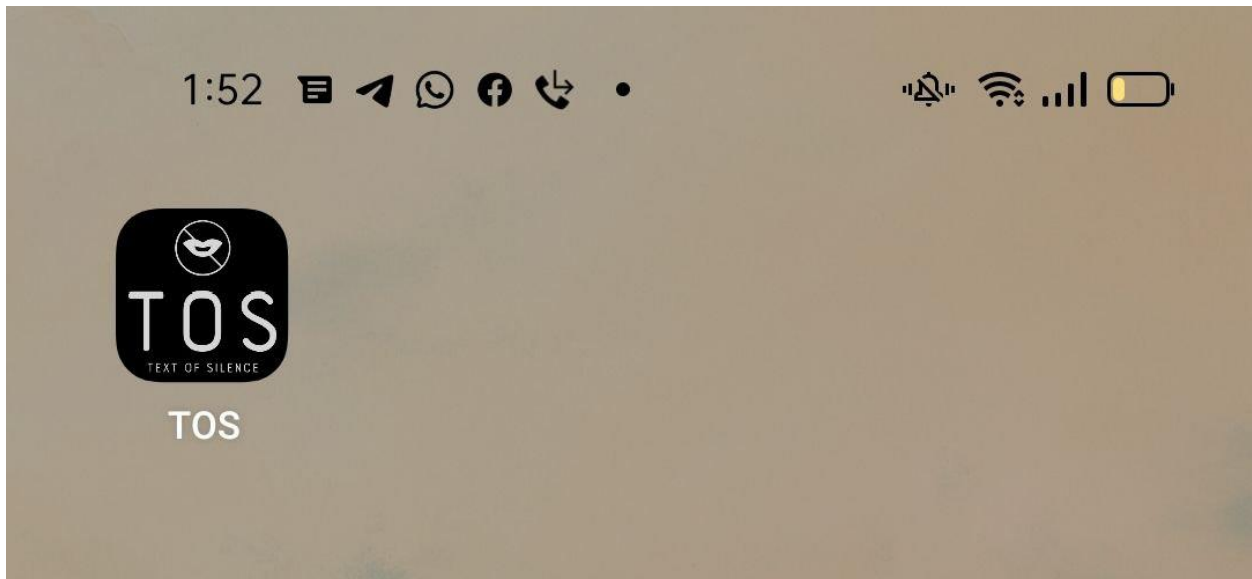


Figure 13. TOS in application drawer

Opening the app navigates to the Home Screen shown in Figure 14. This is just a welcome screen that you just need to click on the “Sentence Structure” button.

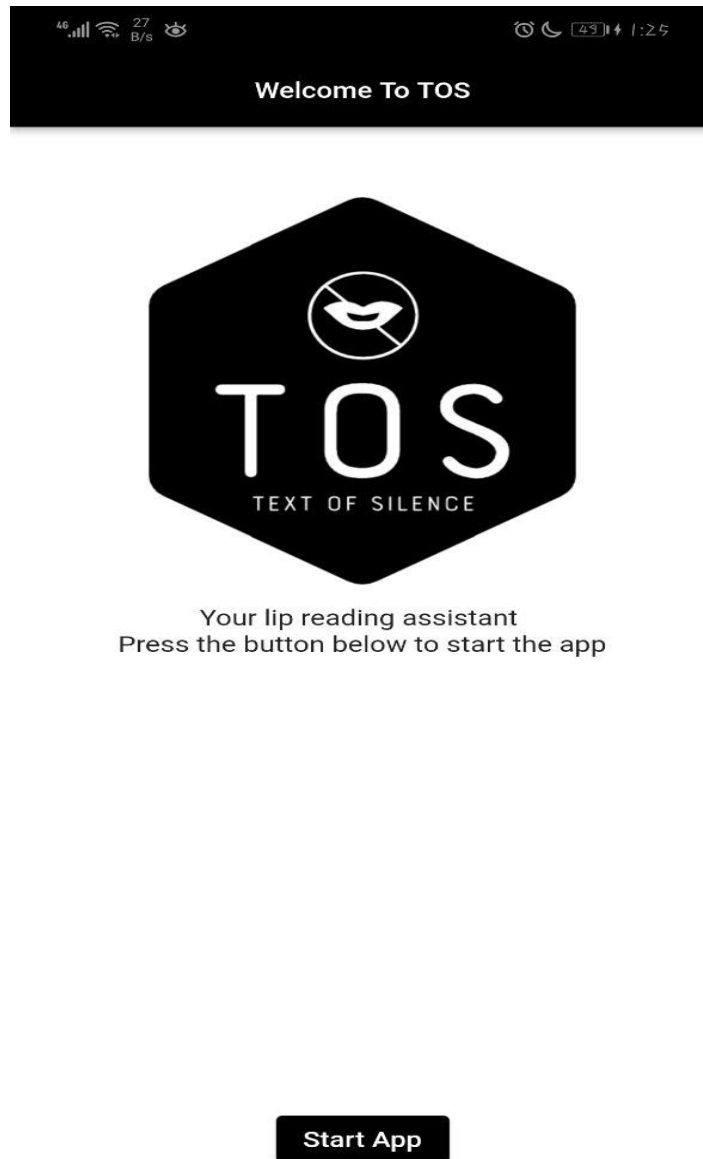


Figure 14. Home Screen

After selecting desired categories as shown in Figure 15. click on the Open Camera button to go to the Camera Screen.

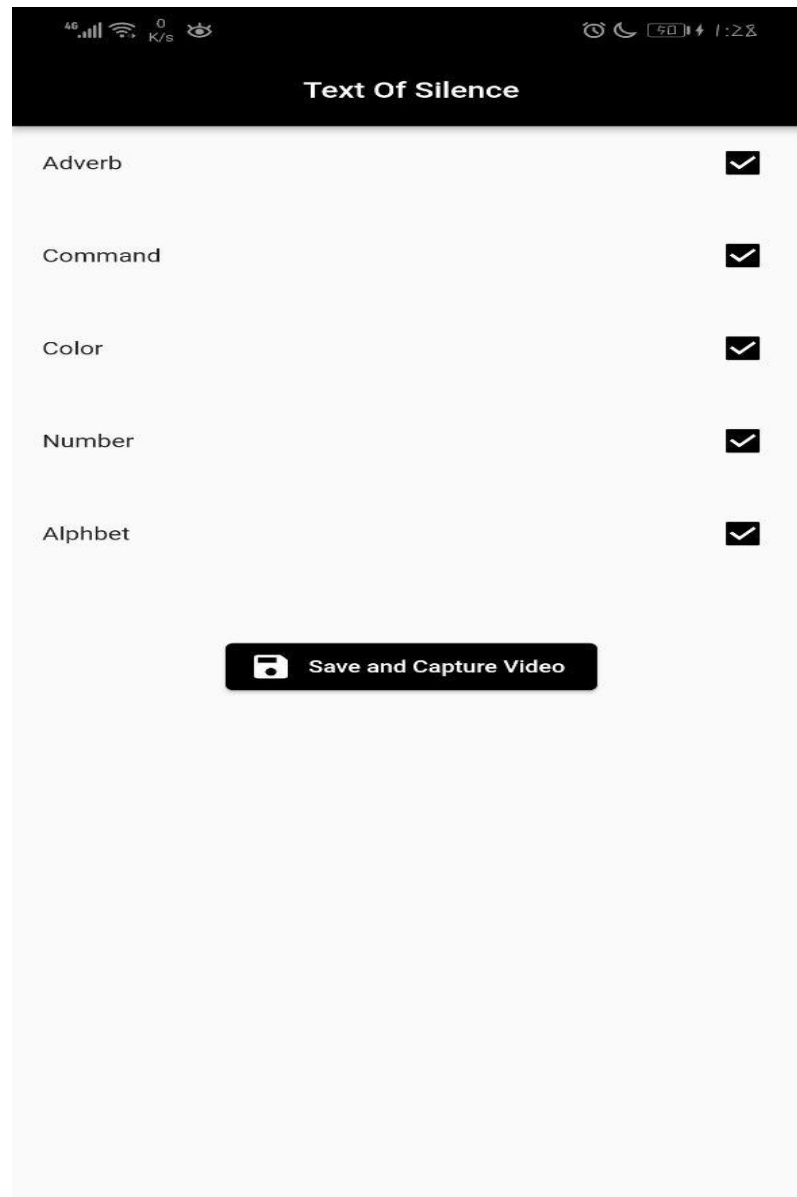


Figure 15. Sentence Structure Screen

A person can capture live video or from pre-recorded videos as shown in Figure 16.

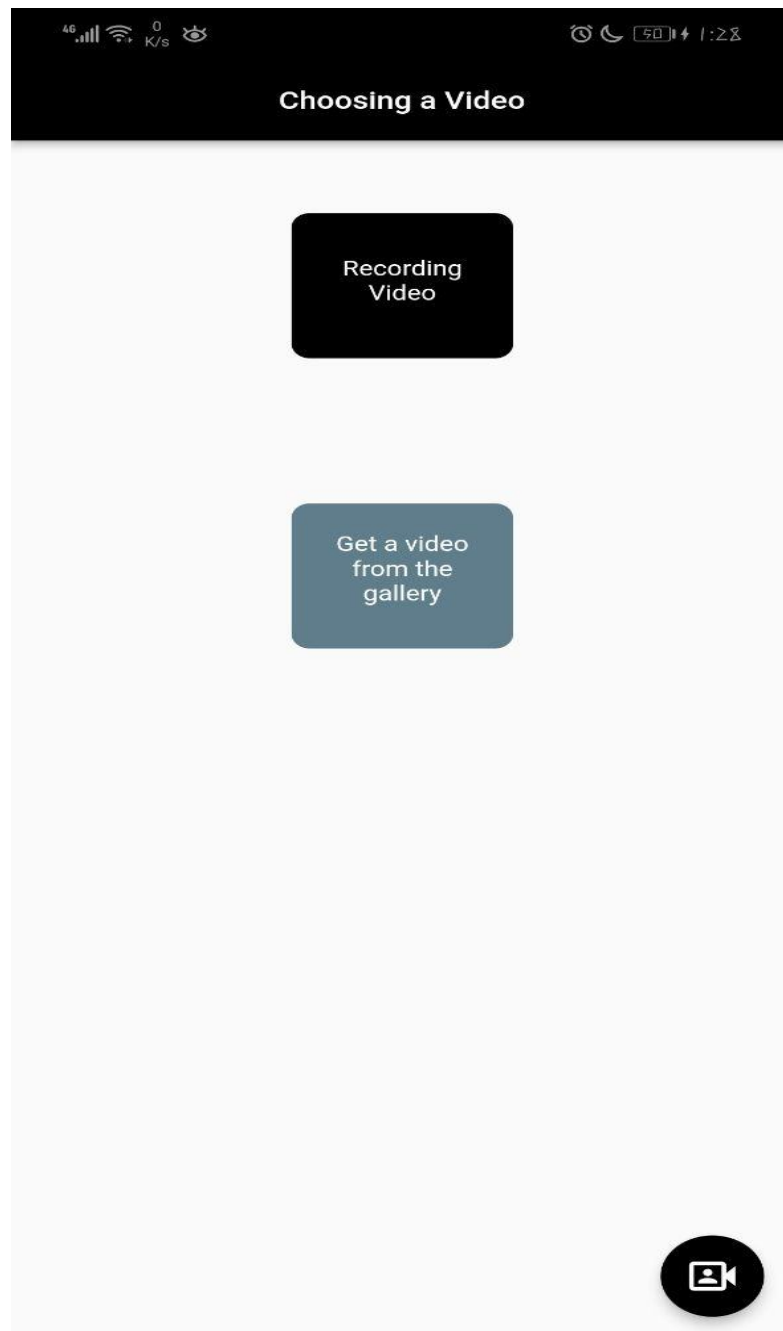


Figure 16. Capture Video

5.2 Samples:

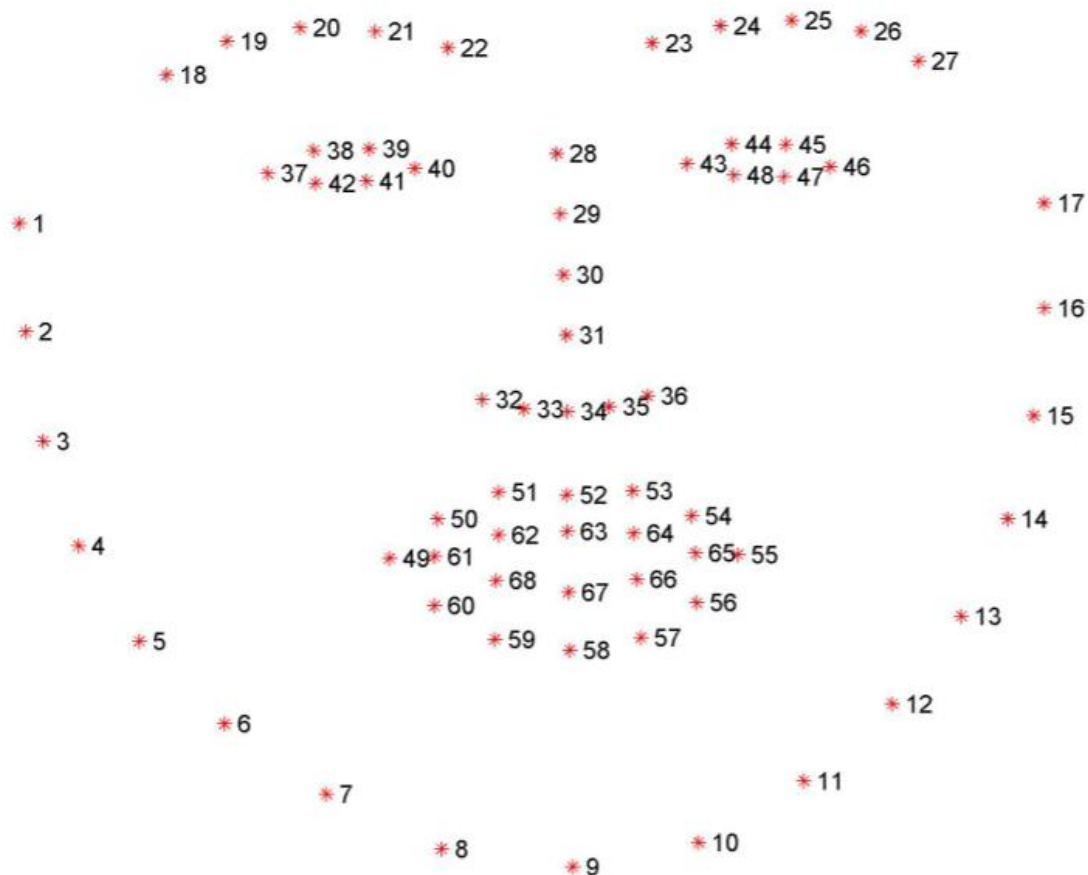
Data Preparation

We have to make some transformations and approximations in order to make it usable.

Capturing Lip

several approaches to use in order to capture the lips from our data set that will be used to train and test the network model

- Skin Segmentation Algorithms to detect skin region of speaker then applying contours on segmented parts to detect face followed by using human face aspect ratio to detect lips region.
- Applying Haar cascade on the face to extract face interest points which includes the mouth interest points that is used in feature extraction phase.
- Apply Facial landmarks on dlib, the mouth can be accessed through points [48, 68]



After pre-processing our dataset and normalizing the videos to 30 frames each, we started extracting the mouth regions from the videos to prepare our data for CNN training. We started by extracting video frames and then extracting the mouth region, in this part we compared the performance of various extraction techniques

	DLIB	Haar Cascade
Time to extract per video	72 Seconds	2.5 Seconds
Accuracy in extraction	~100%	~85%

We found that using DLIB produced an accuracy of almost 100% but made the segmentation process extremely slow and not viable for a large dataset.

That is why we decided to go with using Haar Cascades in the segmentation process as it produces acceptable results in the shortest time

Sample video and output:



Figure. Sample video

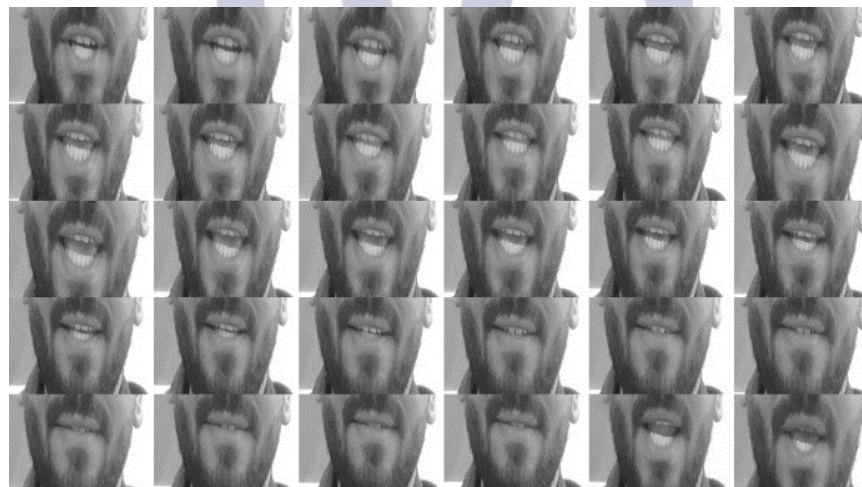


Figure 18. Lip region after segmentation.

References

- ¹Yuxuan Lan¹, Richard Harvey¹, Barry-John Theobald¹, Eng-Jon Ong² and Richard Bowden² <http://www2.cmp.uea.ac.uk/~bjt/avsp2009/proc/papers/paper-35.pdf>
- ²McGurk and J. MacDonald. Hearing lips and seeing voices. *Nature*, 264:746–748, 1976.
<http://wexler.free.fr/library/files/mcgurk%20%281976%29%20hearing%20lips%20and%20seeing%20voices.pdf>
- ³ G. Fisher. Confusions among visually perceived consonants. *Journal of Speech, Language, and Hearing Research*, 11(4):796–804, 1968.
- ⁴ Yannis M. Assael^{1,†}, Brendan Shillingford^{1,†}, Shimon Whiteson¹ & Nando de Freitas^{1,2,3} Department of Computer Science, University of Oxford, Oxford, UK ¹Google DeepMind, London, UK ²CIFAR, Canada ³<https://arxiv.org/pdf/1611.01599.pdf>
- ⁵IEEE : Chung, J. S., Senior, A., Vinyals, O., & Zisserman, A.(2017). Lip Reading Sentences in the Wild. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). doi:10.1109/cvpr.2017.367
<https://ieeexplore.ieee.org/abstract/document/8099850>
- ⁶ British Deaf Association (2015, Sept 7) Fast facts about the Deaf community. Last Accessed 20 March 2020:
<https://bda.org.uk/fast-facts-about-the-deaf-community/>
- ⁷ IMPROVED SPEAKER INDEPENDENT LIP READING USING SPEAKER ADAPTIVE TRAINING AND DEEP NEURAL NETWORKS for Ibrahim Almajai
<https://sci-hub.se/10.1109/ICASSP.2016.7472172>

⁸ Out of Time: Automated Lip Sync in the Wild for Joon Son Chung and Andrew Zisserman

https://link.springer.com/chapter/10.1007%2F978-3-319-54427-4_19

⁹ 1.Lip Reading Sentences in the Wild

<https://sci-hub.se/10.1109/CVPR.2017.367>

¹⁰ 2. Combining Residual Networks with LSTMs for Lipreading
Themos Stafylakis, Georgios Tzimiropoulos

https://www.isca-speech.org/archive/pdfs/interspeech_2017/stafylakis17_interspeech.pdf

¹¹ 3. Learning to lip read words by watching videos
Joon Son Chung*, Andrew Zisserman

<https://sci-hub.se/https://doi.org/10.1016/j.cviu.2018.02.001>

¹² Deep Learning of Mouth Shapes for Sign Language

<https://sci-hub.se/10.1109/ICCVW.2015.69>

¹⁵ Computer Vision Lip Reading Grace Tilton,

http://cs229.stanford.edu/proj2019aut/data/assignment_308832_raw/26646023.pdf