



**KLE** Technological  
University  
Creating Value  
Leveraging Knowledge

**School  
of  
Electronics and Communication Engineering**

**SDP Project Report  
on  
Speaker Verification for Multilingual  
Scenarios**

**By:**

- |                       |                   |
|-----------------------|-------------------|
| 1. Daivarth B N       | USN: 01FE20BEC009 |
| 2. Arvind Morbad      | USN: 01FE20BEC012 |
| 3. Nagaraj S Hiremath | USN: 01FE20BEC016 |
| 4. Vinay Hegde        | USN: 01FE20BEC043 |

**Semester: VII, 2022-2023**

Under the Guidance of

**Prof.Satish**

K.L.E SOCIETY'S  
KLE Technological University,  
HUBBALLI-580031  
2022-23



SCHOOL OF ELECTRONICS AND COMMUNICATION  
ENGINEERING

## CERTIFICATE

This is to certify that project entitled “ **Speaker Verification for Multilingual Scenarios** ” is a bonafide work carried out by the student team of “ **Daivarath B N (01FE20BEC009), Arvind Morbad (01FE20BEC012), Nagaraj Hiremath (01FE20BEC016), Vinay Hegde (01FE20BEC043)** ”. The project report has been approved as it satisfies the requirements with respect to the Minor project work prescribed by the university curriculum for BE (VI semester) in School of Electronics and Communication Engineering of KLE Technological University for the academic year 2022-23.

Satish C Chikkamath  
Guide

Dr.Suneeta V B  
Head of School

Prof. B. S. Anami  
Registrar

External Viva:

Name of Examiners

Signature with date

- 1.
- 2.

## ACKNOWLEDGMENT

Without mentioning the names of the people who kindly assisted us in the process of fulfilling the project's aims and with their persistent directives and support, brought about its fulfillment, these come with its effective completion that would have been incomplete. We express our sincere appreciation and gratitude to the School of Electronics and Communication and the Center for Artificial Intelligence Research[CAIR] lab. Dr. Suneeta V Budihal , the Head of the Department, and Prof. Nirmala and Assistant Prof. Satish C Chikkamath for their guidance and encouragement as we completed this project. We are grateful to our esteemed institution KLE Technological University, Hubballi which has given us this opportunity to realize the most cherished desire to achieve our goal. Finally, we would like to express our gratitude to everyone who assisted us, directly or indirectly, in fulfilling the project's goals and objectives.

-The project team

## ABSTRACT

This paper showcases the efficiency of several algorithms in case of speaker verification using several features. Mel-frequency cepstral coefficients (MFCCs) are used with great care in the feature extraction process of our speaker recognition system development. MFCCs are well known for their effectiveness in encapsulating the inherent properties of speech, and they form the basis of our model's capacity to identify distinct speaker qualities. We build and train four different models: Convolutional Neural Network (CNN), Artificial Neural Network (ANN), SincNet, and VGG16. The architecture of each model is deliberately designed to navigate the complex complexities found in speech patterns in the English language. Training these models on a variety of multilingual data sets guarantees their flexibility and capacity to identify characteristics unique to individual speakers. Our dedication to accuracy and resilience is demonstrated throughout this procedure by the use of state-of-the-art techniques.

# Contents

<b>1</b>	<b>Introduction</b>	<b>8</b>
1.0.1	Overview and scope : . . . . .	9
1.0.2	Applications of CNN: . . . . .	9
1.0.3	Application of ANN : . . . . .	10
1.0.4	Application of Vgg16: . . . . .	10
1.0.5	Application of SincNet : . . . . .	11
1.1	Motivation . . . . .	12
1.2	Objectives . . . . .	12
1.3	Literature survey . . . . .	13
1.3.1	Speaker Recognition with ResNet and VGG Networks . . . . .	13
1.3.2	CNN based speaker recognition in language and text-independent small scale system . . . . .	13
1.3.3	Analytical and iterative approaches to the equalisation of sub-band errors in speech and speaker recognition . . . . .	13
1.3.4	Improving text-independent speaker recognition with GMM . . . . .	14
1.3.5	Conception of Speaker Recognition Methods: A Review . . . . .	14
1.3.6	A method of noisy Tibetan speakers verification based on SKA-TDNN . . . . .	14
1.3.7	Deep CNNs With Self-Attention for Speaker Identification . . . . .	14
1.3.8	An improved Tibetan speaker recognition method based on ResNet . . . . .	15
1.3.9	An Adaptive ResNet Based Speaker Recognition in Radio Commu- nication . . . . .	15
1.3.10	Speaker Identification Using Custom SincNet Layer and Deep Learn- ing . . . . .	15
1.3.11	Speaker Recognition from Raw Waveform with SincNet . . . . .	15
1.4	Problem statement . . . . .	16
1.5	Organization of the report . . . . .	16
<b>2</b>	<b>System design</b>	<b>17</b>
2.1	Methodology for CNN,ANN,VGG16 and SincNET . . . . .	17
<b>3</b>	<b>Implementation details</b>	<b>19</b>
3.1	Specifications and system architecture . . . . .	19
3.2	Algorithms . . . . .	19
3.2.1	Artificial Neural Network [ANN]: . . . . .	19
3.2.2	Convolutional Neural Network [CNN] . . . . .	20
3.2.3	SincNet: . . . . .	21
3.2.4	VGG16: . . . . .	21
3.3	Flowchart . . . . .	22

<b>4</b>	<b>Results and discussions</b>	<b>24</b>
4.1	Result Of CNN,ANN,VGG-16,Sinc-Net . . . . .	24
4.1.1	Result Analysis Of Device-D01 . . . . .	24
4.1.2	Result Analysis Of Device-H01 . . . . .	27
4.1.3	Result Analysis Of Device-M01 . . . . .	29
4.1.4	Result Analysis Of Device-M02 . . . . .	31
4.1.5	Result Analysis Of Device-T01 . . . . .	33
<b>5</b>	<b>Conclusions and future scope</b>	<b>35</b>
5.1	Conclusion . . . . .	35
5.2	Future scope . . . . .	36
5.2.1	Application in the societal context . . . . .	36
	<b>References</b>	<b>36</b>

# Chapter 1

## Introduction

This project employs creative filter design using Convolutional Neural Networks (CNNs), SincNet, Artificial Neural Networks (ANNs), and VGG-16 architecture to advance speaker recognition in multilingual settings in response to the demands of a globalised society. Driven by the need to improve performance and adaptation in the face of linguistic diversity, the initiative strives to recognise speakers more accurately. The goals include evaluating multilingual proficiency, comparing filters, assessing robustness, and examining practical application. In order to achieve thorough feature extraction, effective neural network building, and classification, the project methodologically comprise CNNs, SincNet, ANNs, and VGG-16. A literature review, an explanation of the technique, the specifics of the experimental setting, and results with comparative insights are all included in the report structure. By offering flexibility in traversing linguistic intricacies on a large scale, this study hopes to advance speaker recognition technology.

Our speaker recognition system is based on the careful selection of Mel-frequency cepstral coefficients (MFCCs) throughout the feature extraction process. MFCCs are well known for their ability to capture the fundamental aspects of speech, and they offer a solid basis for our models to identify minute differences in the acoustic characteristics of spoken language. The technique for feature extraction plays a crucial role in identifying the distinctive voice patterns that set each speaker apart. Meanwhile, we have demonstrated our dedication to a multimodal strategy by choosing four different models: Convolutional Neural Network (CNN), Artificial Neural Network (ANN), SincNet, and VGG16. Every model is customised to maximise the identification of patterns unique to a given speaker within the complex English language speech.

The architecture of each model is deliberately designed to navigate the complex complexities found in speech patterns in the English language. Training these models on a variety of datasets guarantees their flexibility and capacity to identify characteristics unique to individual speakers. In addition to the more traditional CNN and ANN architectures, we have complemented them with state-of-the-art techniques like SincNet and VGG16, demonstrating our dedication to accuracy and resilience throughout this process.

In testing, our assessment process is built to thoroughly evaluate the models' performance with speech samples in a language designated as the "favourite language." The objective of this strategic approach is to evaluate the system's ability to generalise and adapt. By combining various models and concentrating on the nuances of the English language, our speaker identification system is positioned as a high-tech solution that can make a substantial contribution to the complex field of multilingual speaker recognition.

### 1.0.1 Overview and scope :

The goal of improving speaker identification systems, this research has a broad scope that includes many important aspects. Mel-frequency cepstral coefficients (MFCCs), which are known for their efficiency in capturing speech characteristics, are used meticulously in the feature extraction step. The project uses four different neural network architectures for model development: Convolutional Neural Network (CNN), Artificial Neural Network (ANN), SincNet, and VGG16. Each architecture is specifically designed to identify patterns unique to a speaker, especially in the subtleties of English language speech. By utilising band-pass filters and hierarchical representations found in voice data, the project's scope is expanded into complex neural network designs through the use of cutting-edge approaches like SincNet and VGG16. Despite the fact that English language patterns are the main emphasis, the project's design and techniques are prepared to contribute[7].

The context of this work, a thorough investigation of Mel-frequency cepstral coefficients (MFCCs) is a fundamental component of the feature extraction procedure, with a focus on their function in encapsulating crucial speech attributes. By adding four unique neural network models—Convolutional Neural Network (CNN), Artificial Neural Network (ANN), SincNet, and VGG16—the project broadens its scope. Every model is painstakingly created to focus on identifying distinct patterns particular to a certain speaker within the complex of English language speech.

By incorporating state-of-the-art methods like SincNet and VGG16, the project showcases a forward-thinking attitude and digs into advanced methodology. Through the use of band-pass filters, SincNet introduces novel neural network design, improving the system's ability to collect important temporal characteristics. Concurrently, the system's capacity to recognise hierarchical representations in speech data is further improved by VGG16's deep and complex architecture.

Although English language speech is still the major focus, the project acknowledges the wider implications for multilingual applications. The system's versatility allows it to function in a variety of linguistic environments, perhaps opening up new application areas outside of the English language. This flexibility is demonstrated by the use of several testing methodologies and speech samples in a selected "favourite language" in the assessment stage. By employing this tactical method, the project guarantees a thorough evaluation of the system's flexibility and generalisation capacities in a range of language circumstances. This study, taken as a whole, combines state-of-the-art approaches, diverse models, and excellent feature extraction with an eye towards the nuances of English speech patterns and the larger field of multilingual speaker recognition.

### 1.0.2 Applications of CNN:

For speaker recognition in the context of this project, convolutional neural networks (CNNs) provide a flexible and effective tool. CNNs are well-known for their effectiveness in image and signal processing applications. They are also good at learning hierarchical representations, which makes them a good choice for identifying complex characteristics in speech data. In the field of feature extraction, CNNs are highly proficient in automatically recognising and extracting pertinent patterns from unprocessed input, which enables the identification of subtle acoustic properties present in spoken language. Our



project's primary feature extraction method relies on the network's convolutional layers, which function as localised filters to efficiently capture temporal dependencies and spatial correlations within the MFCCs.

Additionally, CNNs greatly enhance the speaker identification system's flexibility. The model's ability to identify speaker-specific patterns even in the face of variances in pronunciation, accent, and other linguistic nuances is improved by the network's ability to generalise well to a variety of speaker characteristics thanks to the acquired hierarchical representations. Large-scale speaker identification problems benefit greatly from CNNs' suitability for parallel processing, which facilitates effective computation and scalability. In conclusion, the use of CNNs in this project improves the process of extracting features, makes modelling intricate speech patterns easier, and increases the speaker recognition system's overall resilience and adaptability. !

### **1.0.3 Application of ANN :**

In this research, speaker recognition is mostly dependent on the use of artificial neural networks (ANNs), especially when it comes to feature extraction and pattern recognition. Artificial Neural Networks (ANNs) are a useful tool for extracting speaker-specific features from Mel-frequency cepstral coefficients (MFCCs) because they can learn complex mappings and capture nuanced relationships within data. ANNs are able to create abstract representations that reflect the unique subtleties of speech during the feature extraction phase by extracting relevant characteristics on their own.

Patterns of speech unique to a speaker. Because of this, they can distinguish between speakers with different accents, intonations, and speaking styles with ease. The process of training artificial neural networks (ANNs) makes it easier to build a reliable model that can discriminate between speakers with accuracy.

By learning from a variety of training datasets and making generalisations, ANNs further enhance the versatility of the speaker recognition system. Recognising speakers in a variety of circumstances is made easier by ANNs' ability to adjust to subtle linguistic differences and variations in speech patterns. In order to help the model recognise and highlight important elements in the input data, the network's hidden layers function as feature transformers.

The intricate and non-linear nature of speaker-specific patterns in speech can also be captured by ANNs because they are proficient at handling non-linear relationships. Their ability to distinguish between speakers with different accents, intonations, and speaking styles is thereby enhanced. An accurate model that can discriminate between several speakers can be created more easily with the help of ANNs and their training procedure.

In conclusion, this project's use of artificial neural networks enhances the system's capacity to automatically learn and represent complex speech aspects, making them essential elements for developing a reliable and flexible speaker recognition system.

### **1.0.4 Application of Vgg16:**

The speaker recognition system's capabilities are greatly improved in this project by using the VGG16 architecture, especially in the areas of deep feature extraction and hierarchical representation learning. VGG16 is well-known for its ability to perform well in image classification tasks due to its deep and complex design; using this technology to speaker recognition has a number of benefits.

A potent technique for learning complex patterns inside the Mel-frequency cepstral coefficients (MFCCs) is offered by the VGG16 architecture, which is distinguished by its stack of convolutional layers with narrow receptive fields, followed by fully connected layers. By capturing both high-level and low-level representations of speech information, the deep convolutional layers serve as feature extractors. The process of hierarchical representation learning plays a crucial role in identifying intricate and intangible patterns evident in speech characteristics unique to individual speakers.

Additionally, the feature transferability that VGG16 learns from extensive image datasets is utilised for tasks involving voice recognition. Improved performance is achieved, particularly when dealing with limited labelled speaker data, by fine-tuning the pre-trained VGG16 model on speaker-specific data, which benefits from the generalised representations learnt from varied data sources.

As a result of the system's ability to learn and identify complex patterns within the voice data, the deployment of VGG16 essentially enhances the feature extraction process by utilising the depth of its design. The speaker recognition system is more resilient and has greater discriminating ability thanks to the hierarchical representations produced by VGG16.

### **1.0.5 Application of SincNet :**

This project's use of SincNet presents a novel and specialised method of neural network architecture, especially when it comes to speaker recognition. Using band-pass filters in the neural network's first layers gives SincNet a distinct benefit since it allows the network to concentrate on capturing important temporal characteristics seen in the Mel-frequency cepstral coefficients (MFCCs).

The use of SincNet is especially beneficial when the speaker identification system is in the feature extraction stage. By adaptively capturing and emphasising particular frequency bands, the band-pass filters improve the model's capacity to identify minute differences in speech patterns. This is important for tasks involving speaker recognition since the temporal dynamics of speech are frequently what set speakers apart.

Furthermore, by offering a specialised technique for modelling speaker-specific patterns, SincNet enhances the system's adaptability. As learnable parameters, the band-pass filters enable the network to adjust to various language subtleties, accents, and speech pattern variances. This flexibility is particularly useful in situations involving many languages, where speakers may have different speaking patterns.

In conclusion, the usage of SincNet enhances the process of feature extraction by bringing a customised method that emphasises temporal dynamics. This novel approach improves the system's flexibility and helps it identify patterns unique to each speaker, which makes it an important part of building a strong and specialised speaker recognition system.

## 1.1 Motivation

The rationale behind the integration of various filtering mechanisms, namely Artificial Neural Network (ANN), Convolutional Neural Network (CNN), VGG16, and SincNet, is rooted in their distinct advantages and skills, each of which enhances speaker recognition systems in a unique manner.

Artificial Neural Networks (ANN) use their capacity to discover intricate relationships within data to provide a flexible and adaptive method of feature extraction. Driven by the ability to extract complex patterns, ANNs provide a thorough investigation of speaker-specific features found in Mel-frequency cepstral coefficients (MFCCs). They are ideal for managing variances in speech characteristics and pronunciation among various speakers due to their adaptability.

Convolutional Neural Networks (CNN) are so good at capturing spatial hierarchies, they are especially useful for image and signal processing applications. CNNs are very good at automatically extracting hierarchical features from MFCCs in the context of speaker recognition, which makes it possible to identify subtle acoustic properties in voice data. By capturing both temporal dependencies and spatial correlations, their localised filters help the network create a strong representation of speaker-specific patterns.

A hierarchical representation learning technique is introduced by the VGG16 architecture, which is well-known for its complexity and depth. VGG16 enhances the process of feature extraction by utilising the depth of its convolutional layers to capture both high-level and low-level representations of communication information. Especially when trained on limited labelled speaker data, VGG16 is an interesting choice for improving the discriminative capability of speaker identification systems due to the transferability of features learned from multiple data sources, as shown in picture classification tasks.

Using band-pass filters in its first layers, SincNet concentrates on temporal dynamics in a novel way. SincNet is especially well-suited for capturing important temporal aspects within MFCCs because of its novel approach. The ability of SincNet to dynamically modify the band-pass filters shows how versatile the network is. This allows the network to focus on particular frequency bands and react to a variety of linguistic subtleties and speech differences.

To summarise, the rationale behind combining these four filters is to take advantage of their unique advantages: ANNs are flexible, CNNs are good at extracting hierarchical spatial features, VGG16 is good at depth and transferability, and SincNet is good in specialised temporal dynamics. By combining various filtering techniques, a comprehensive speaker recognition system that is adept at identifying the complex relationships between speaker-specific patterns in voice data is intended to be created.

## 1.2 Objectives

Project goals include the thorough investigation and refinement of four different filtering mechanisms: VGG16, SincNet, Convolutional Neural Network (CNN), and Artificial Neural Network (ANN)—in order to improve speaker recognition systems. SincNet is to be used to focus on temporal dynamics, while CNNs are to capture hierarchical spatial

features, VGG16 is to exploit depth and transferability, and feature extraction techniques are to be improved. ANNs' flexibility is also to be fully utilised. Furthermore, the study attempts to test how well different filtering techniques compare in terms of being able to identify speaker-specific patterns in Mel-frequency cepstral coefficients (MFCCs). Through evaluation of the methods' effectiveness across a range of linguistic subtleties and variations, the project aims to provide insights into the complex field of speaker recognition and promote a better comprehension of the advantages and practicality of each filtering technique.

## **1.3 Literature survey**

### **1.3.1 Speaker Recognition with ResNet and VGG Networks**

In this research paper there is use of deep neural networks—ResNet and VGG architectures in particular—for speaker recognition is examined in the study "Speaker Recognition with ResNet and VGG Networks"[1]. The efficacy of these architectures in extracting and learning discriminative features from audio data for precise speaker identification is examined in this study. The authors examine model complexity and training data, among other things, in their analysis of ResNet and VGG network performance through experimentation. The outcomes show how deep learning frameworks, in particular ResNet and VGG, can be used to improve speaker recognition systems and offer important new information for the field of audio-based biometrics.

### **1.3.2 CNN based speaker recognition in language and text-independent small scale system**

This research paper discusses regarding the CNN-based Speaker Recognition in Language and Text-Independent Small-Scale System is a research paper that explores the use of Convolutional Neural Networks (CNNs) for speaker recognition in situations where language and text independence are paramount. The study emphasises efficiency without sacrificing precision, with a focus on small-scale systems[2]. The authors hope to extract strong speaker characteristics from audio data by using CNNs. The study investigates the network's capacity to recognise distinct vocal traits in a range of speech settings and languages. The results demonstrate the feasibility of using CNNs to build efficient speaker recognition models appropriate for real-world, language- and text-independent applications in resource-constrained settings.

### **1.3.3 Analytical and iterative approaches to the equalisation of sub-band errors in speech and speaker recognition**

This research article "Analytical and Iterative Approaches to the Equalisation of Sub-Band Errors in Speech and Speaker Recognition" discusses methods for reducing sub-band errors in speech and speaker recognition contexts[3]. The paper presents iterative and analytical techniques to equalise mistakes that occur in the speech signal's various frequency components. The scientists want to improve the robustness and accuracy of recognition systems by implementing these strategies. The study offers insightful information on the difficulties posed by sub-band errors and suggests practical solutions for reducing them, which enhances the functionality of speech and speaker recognition systems as a whole.

#### **1.3.4 Improving text-independent speaker recognition with GMM**

The research "Improving Text-Independent Speaker Recognition with GMM" looks into ways to use Gaussian Mixture Models (GMMs) to improve text-independent speaker recognition. The study investigates the use of GMMs to improve recognition performance with the goal of more accurately simulating speaker variability[4]. The goal of the project is to enhance text-independent speaker recognition systems' accuracy and dependability by utilising GMMs[5]. This study provides important new insights into the application of probabilistic models, in particular GMMs, to speaker recognition problems and improving performance when the text is unknown beforehand.

#### **1.3.5 Conception of Speaker Recognition Methods: A Review**

A thorough summary of speaker recognition techniques is given in the paper "Conception of Speaker Recognition Methods: A Review." The study examines a range of speaker recognition techniques, including both traditional and contemporary methods[6]. This paper examines the development of speaker recognition techniques, emphasising significant breakthroughs, obstacles, and patterns. The review attempts to provide a coherent understanding of the idea and evolution of speaker recognition techniques by synthesising the body of existing literature. By combining existing knowledge and laying the groundwork for upcoming studies and advancements in speaker recognition technology, this benefits the field as a whole.

#### **1.3.6 A method of noisy Tibetan speakers verification based on SKA-TDNN**

The research paper "A Technique for Verifying Noisy Tibetan Speakers" Based on SKA-TDNN" describes a technique that makes use of SKA-TDNN (Sinc Kernel Activation with Time-Delay Neural Network) to confirm the identity of Tibetan speakers in noisy environments[7]. The paper discusses background noise issues in speaker verification, specifically as they pertain to Tibetan speakers. The suggested technique makes use of the SKA-TDNN architecture to improve speaker verification systems' resilience. By using this method, the research hopes to increase speaker verification for Tibetan speakers in noisy, real-world environments in terms of accuracy and reliability[8]. By providing a customised solution to address issues unique to the Tibetan language and its speakers, the paper advances the field.

#### **1.3.7 Deep CNNs With Self-Attention for Speaker Identification**

This research was conducted in order to identify speakers, a technique known as "Deep CNNs with Self-Attention for Speaker Identification" combines self-attention mechanisms with Deep Convolutional Neural Networks (CNNs)[9]. The goal of this hybrid approach is to take advantage of both the contextual information capture offered by self-attention and the hierarchical feature learning capabilities of CNNs. The usefulness of this model in extracting distinctive speaker traits from audio data is investigated in this study. Through the integration of self-attention into deep CNN architectures, the suggested approach aims to improve the model's discriminative ability for precise speaker identification. By

combining deep learning techniques to achieve better performance in the speaker identification domain, the paper advances the field of speaker recognition technology.

### **1.3.8 An improved Tibetan speaker recognition method based on ResNet**

The Residual Neural Network (ResNet) architecture is used to improve Tibetan speaker recognition, as presented in the paper "An Improved Tibetan Speaker Recognition Method Based on ResNet"[11]. The goal of the project is to improve speaker recognition performance and accuracy, particularly for Tibetan speakers. The suggested approach uses ResNet to capture complex speaker features and enhance the model's discriminative power. By providing improvements in the use of deep learning methods, specifically ResNet, to address the particular qualities and difficulties related to Tibetan speaker recognition, the research advances the field. The study offers insights into the architecture and efficacy of the enhanced technique for precisely identifying Tibetan speakers.

### **1.3.9 An Adaptive ResNet Based Speaker Recognition in Radio Communication**

This research paper discusses regarding a speaker recognition system designed specifically for radio communication environments is presented in the paper "An Adaptive ResNet-Based Speaker Recognition in Radio Communication"[10]. The technique makes use of ResNet (Residual Neural Network) architecture that is adaptive in order to meet the unique requirements of radio communication scenarios. The study discusses the requirement for reliable speaker identification in erratic and noisy radio environments. The suggested method seeks to improve speaker recognition systems' accuracy and dependability in radio communication environments by modifying ResNet to fit its unique requirements[12]. The study offers insights into the creation of speaker recognition models that are adaptive to the particular difficulties posed by radio communication scenarios.

### **1.3.10 Speaker Identification Using Custom SincNet Layer and Deep Learning**

The speaker identification technique described in the paper "Speaker Identification Using Custom SincNet Layer and Deep Learning" makes use of a custom SincNet layer integrated into a deep learning framework[13]. A particular kind of neural network layer called SincNet is made especially for handling unprocessed audio waveforms. Here, speaker identification is the main goal of the research, which makes use of SincNet's special feature extraction powers from audio signals. The model's ability to learn discriminative features for precise speaker identification is improved by the deep learning architecture. The work advances the field by presenting a specially designed layer for audio processing and demonstrating how well it works for speaker identification tasks when combined with deep learning methods.

### **1.3.11 Speaker Recognition from Raw Waveform with SincNet**

The speaker recognition technique described in the paper "Speaker Recognition from Raw Waveform with SincNet" uses SincNet to process raw audio waveforms directly[14]. A

neural network layer called SincNet is intended to efficiently extract features from unprocessed inputs. This work focuses on speaker recognition, directly extracting pertinent speaker-specific features from the unprocessed audio data using SincNet. This approach seeks to increase accuracy and expedite the processing pipeline by eschewing conventional feature extraction methods. By demonstrating the use of SincNet for speaker recognition tasks and emphasising the benefits of working directly with raw waveform data, the paper advances the field.

## **1.4 Problem statement**

Speaker Verification for Multilingual Scenarios

## **1.5 Organization of the report**

In Chapter 2, the system design includes a subsection dedicated to the methodology. This section provides a visual representation of the project flow, highlighting the various blocks involved in the system.

Chapter 3 of the document consists of three main sections: Specifications and System Architecture, Algorithm, and Flowchart. These sections encompass detailed information about the project, including the software utilized, the algorithms employed, and the project flowchart that illustrates the sequence of steps.

Chapter 4 of the document focuses on Results and Discussion. In this chapter, the conducted result analysis is presented, and insights derived from the project are provided. The findings and observations obtained from the project are discussed, analyzed, and interpreted, shedding light on the outcomes and implications of the conducted work.

Chapter 5 of the document is dedicated to the Conclusion and Future Scope. In this chapter, the key findings and conclusions drawn from the project are summarized. The overall significance and implications of the work are discussed, highlighting its contributions and potential impact. Additionally, the chapter explores avenues for future research and improvement, suggesting areas where the project can be further optimized or expanded upon to address any limitations or open research questions.

# Chapter 2

## System design

This project's system design comprises combining four filtering mechanisms—ANN, CNN, VGG16, and SincNet—into a single, coherent framework for speaker recognition. The fundamental technique for feature extraction is called mel-frequency cepstral coefficients, or MFCCs. Each model in the architecture is trained on a variety of English language datasets in order to maximise its flexibility and capacity to identify characteristics unique to individual speakers. ANNs and CNNs exploit their natural ability to recognise complex patterns during training, VGG16 makes use of its depth in hierarchy, while SincNet uses band-pass filters to concentrate on temporal dynamics. The methodology for testing include assessing the models with speech samples in a specified "favourite language," guaranteeing a thorough evaluation of the system's flexibility. By combining various filtering techniques, we hope to develop a strong and adaptable speaker recognition system that can handle the challenges of multilingual.

### 2.1 Methodology for CNN,ANN,VGG16 and SincNET

#### 2.1.1 CNN

##### 1. Define the Issue:

Define the task of speaker recognition clearly and indicate that the model should be capable of recognising speakers in multiple languages.

##### 2. Data Gathering:

Collect a broad dataset of audio recordings of speakers from various languages. Ascertain that the dataset is well-balanced and indicative of the intended application.

Fill in the blanks with speaker names and language descriptors.

##### 3. Data Preparation:

Convert audio data into a format that can be fed into a CNN (for example, Mel-frequency cepstral coefficients - MFCCs).

To maintain consistency between recordings, normalise and standardise the data.

##### 4. Optional language embeddings:

Consider adding linguistic information to the model. This might be accomplished by introducing a language embedding as a new input to the network.



#### 5. Divide the Data:

Divide the dataset into three parts: training, validation, and testing. Make sure that each set includes a fair distribution of speakers and languages.

#### 6. Create the CNN Architecture:

Create a CNN architecture capable of extracting speaker-related properties from audio data. To analyse temporal patterns in audio recordings, consider utilising 1D convolutional layers.

Experiment with different layer counts, filter sizes, and pooling algorithms. Include batch normalisation and activation functions (for example, ReLU).

#### 7. Model Compilation:

For both speaker and language recognition tasks, use appropriate loss functions. You may need to combine numerous loss functions because this is a multi-task situation. Choose an optimizer and metrics that are appropriate for the task.

#### 8. Educating the CNN:

Using the training set, train the model. Overfitting in the validation set should be monitored, and hyperparameters should be adjusted accordingly. Consider using data augmentation strategies to boost the diversity of your training set.

#### 9. Examine the Model:

Examine the model's performance on the test set, taking into account both speaker recognition and language identification accuracy.

For speaker recognition, use metrics such as equal error rate (EER) and classification accuracy for language identification.

#### 10. Language-Independent Features (Optional):

Investigate the extraction of language-independent features from audio data. This could imply employing embeddings or representations that are less reliant on language-specific properties.

#### 11. Optional fine-tuning and transfer learning:

Experiment with fine-tuning strategies, especially if you have access to huge audio datasets with pre-trained models. When labelled data is scarce, transfer learning can be useful.

#### 12. Optional post-processing:

Implement post-processing techniques like as smoothing or voting processes across various time frames to improve final predictions.

#### 13. Deployment:

Deploy the model to the target environment, taking into account available computational resources and any limits imposed by real-time processing requirements.

#### 14. Continuous Monitoring and Upkeep:

In the production environment, monitor the model's performance and update it as needed. Changes in speaker characteristics and language distributions must be accommodated.

# Chapter 3

## Implementation details

When undertaking a project, it is essential to include specific parameters in the implementation details to ensure that the results can be easily comprehended at a glance. These commonly include specifications, architecture, flowchart, algorithm, and methodology. Here are some implementation details to consider.

### 3.1 Specifications and system architecture

The programming language made use is python in it's ipynb that is the jupyter notebook format, so that it is flexible to work either in Jupyter notebook, VS Code or even Google Colab.

The framework includes the dataset used is from IIT Guwahati which an enormous data of 100 speakers which includes the recording of them reading and conversing in english language and another conversing in their favorite language which might be their mothertongue. The audio dataset is that are recorded from 5 different devices including 1 Laptop (40 K), 1 Table and chair (10 K), 1 zoom microphone (25 K), 2 mobile (1 call IVR, 1 for recording)/ 1 mobile and 1 medium quality microphone (20 K), and 1 Multi array setup (15 K). Therefore, the data gathered is Multilingual using Multisensor (6 sensors) in Multi-environment (2 type) that is Clean (controlled environment) may be an Office and the another is a Noisy (Market/Railway station) (may be outsource) environment in different Styles that is reading and conversation.

### 3.2 Algorithms

#### 3.2.1 Artificial Neural Network [ANN]:

Artificial Neural Network (ANN) architecture functions by imitating the networked organisation of neurons found in the human brain. An input layer, hidden layers, and an output layer make up an artificial neural network (ANN), which is made up of layers of nodes. The weights of each neuronal connection change as a result of training. Input data is passed forward through the network during its learning process, and weighted sums are calculated at each neuron[1]. By introducing non-linearity, the activation function enables the network to learn intricate patterns. Using backpropagation, the model modifies weights during training in order to minimise the discrepancy between expected and actual outputs[14]. Because they can learn from a variety of datasets, artificial neural networks

(ANNs) are very adaptable and can be used for a wide range of tasks, including speaker recognition and feature extraction in complicated data domains.

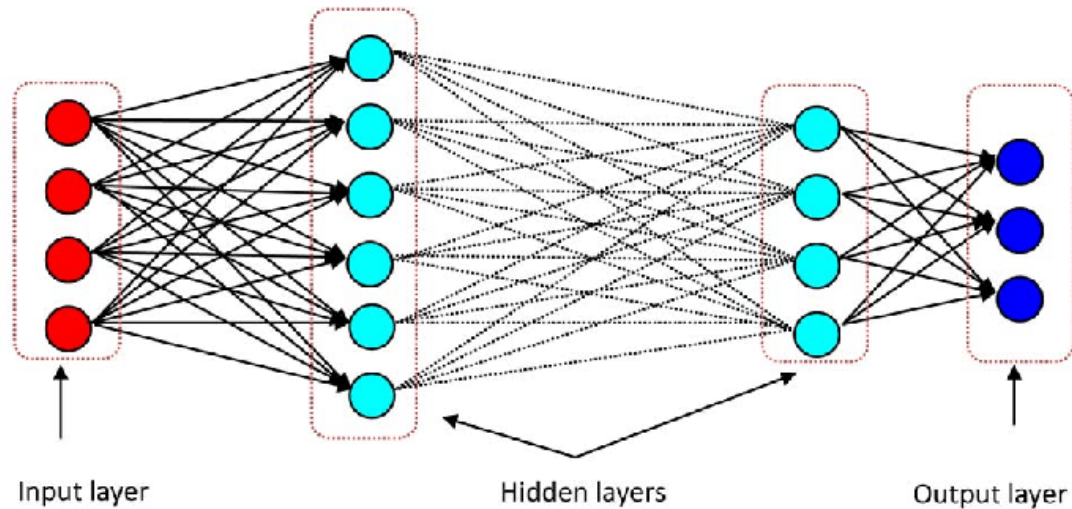


Figure 3.1: ANN Architecture

### 3.2.2 Convolutional Neural Network [CNN]

:

Utilising a deep convolutional neural network (CNN) structure distinguished by its unique depth, the VGG16 architecture functions. With 16 weight layers—13 convolutional and 3 fully connected—VGG16 is an excellent choice for hierarchical feature extraction. Small 3x3 filters are used in each convolutional layer, and max-pooling layers come after the convolutional blocks to help the network capture both high-level and low-level data. since of its depth, VGG16 is very good at picture and pattern recognition applications since it can learn intricate hierarchical representations[2]. Through backpropagation, the network modifies its weights during training to maximise its capacity to identify complex patterns in the input data. The success of VGG16 is attributed to its deep architecture and simplicity, which demonstrate its adaptability in a range of computer vision applications, such as picture categorization.

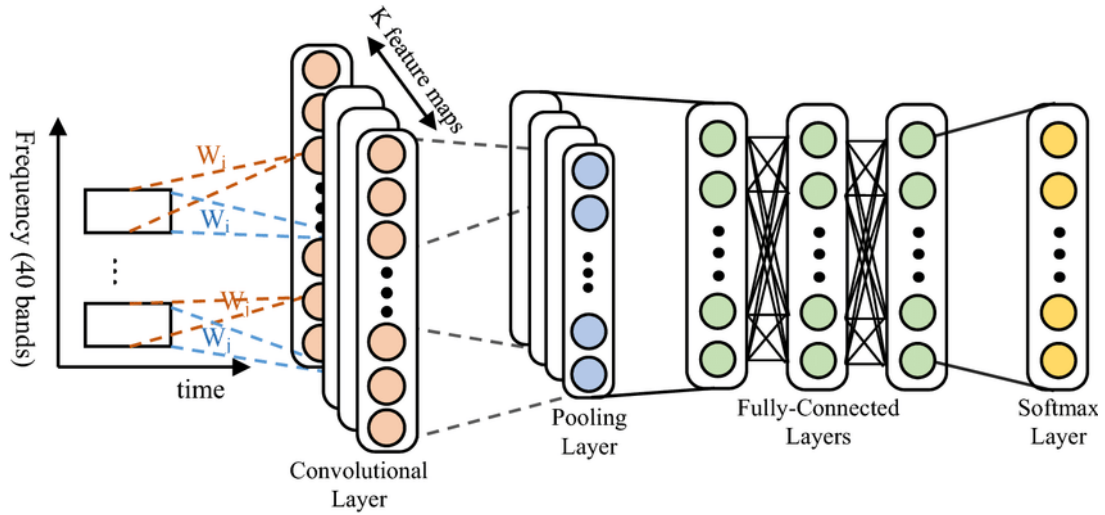


Figure 3.2: CNN Architecture

### 3.2.3 SincNet:

SincNet, a specialized neural network architecture for speaker recognition, operates uniquely by utilizing band-pass filters in its initial layers. These filters are designed to focus on capturing critical temporal features within the input signal, often represented by Mel-frequency cepstral coefficients (MFCCs)[3]. The architecture introduces a novel approach by making these filters learnable during training, allowing the network to adapt and emphasize specific frequency bands relevant to speaker-specific patterns. SincNet's design effectively captures temporal nuances in speech signals, providing a tailored mechanism for discerning speaker characteristics. By leveraging band-pass filters, SincNet enhances the model's ability to capture subtle variations in speech patterns, contributing to its effectiveness in tasks related to speaker recognition, especially in scenarios where temporal dynamics play a crucial role in distinguishing speakers.

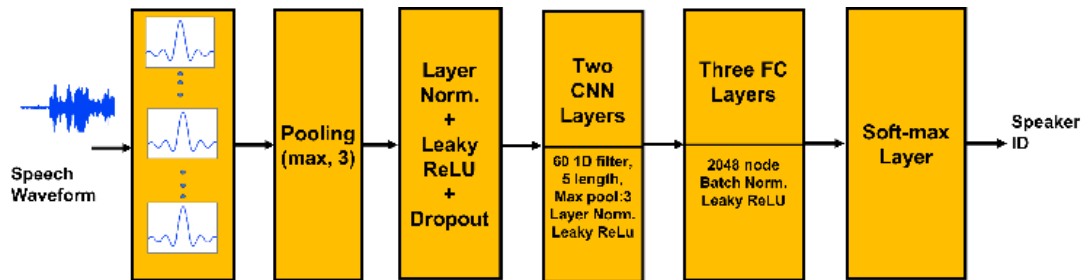


Figure 3.3: SincNet Architecture

### 3.2.4 VGG16:

In order to interpret grid-like input, such as pictures or, in this example, Mel-frequency cepstral coefficients (MFCCs) in speech recognition, Convolutional Neural Networks (CNNs) use a hierarchical and learnable architecture. Convolutional, pooling, and fully linked layers are a CNN's fundamental building blocks[4]. Learnable filters are utilized by convolutional layers to capture local patterns and features by convolving over input data. Next, in order to reduce computational complexity while preserving crucial information,

pooling layers downsample the spatial dimensions. The network can learn intricate associations and generate predictions because to the hierarchical feature representations that are fed into fully connected layers[13]. Non-linearity is introduced by using non-linear activation functions, like Rectified Linear Unit (ReLU), which improves the network's ability to recognise complex patterns. Overfitting is prevented by regularisation strategies like dropout.

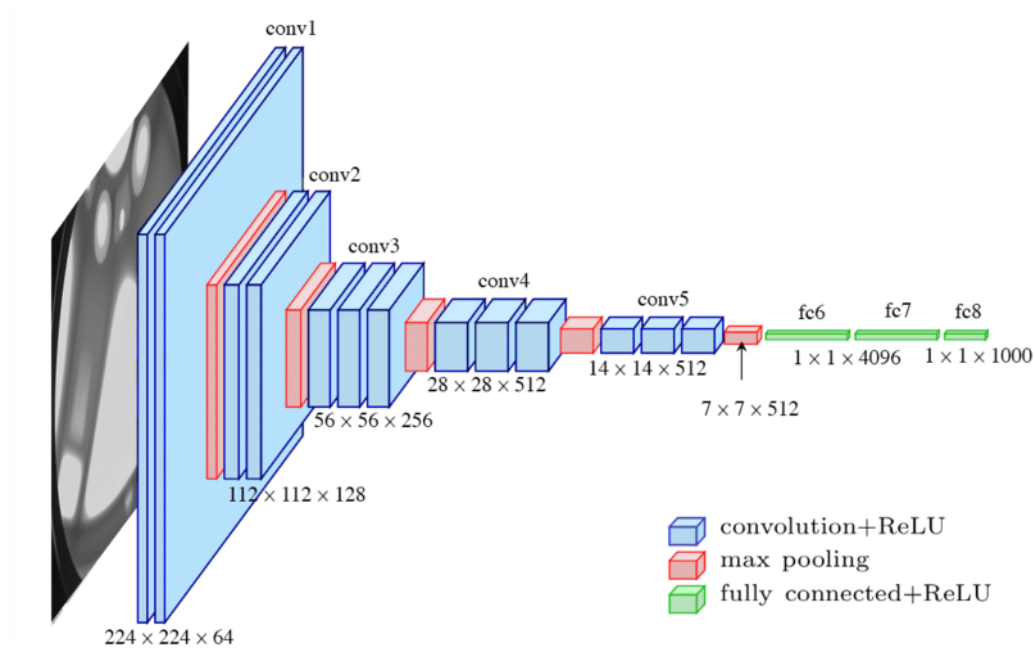


Figure 3.4: VGG16 Architecture

### 3.3 Flowchart

Figure 3.1 represents the project's flowchart that is the block diagram of the model pipeline and work flow, illustrating the sequential execution of the Machine Learning [ML] algorithm, where each block is individually represented with its own set of operations.

Primarily the data is loaded using the librosa library which helps to load and pre-process the data, followed by the feature extraction where around 40 features are extracted, remove unwanted columns and have our own csv file with only required extracted features[8]. Using the sklearn library the data is then divided in to test and train by splitting the english spoken audio data as the input for training and later use the favourite language as the testing dataset to verify the authenticity and efficiency of the models.

Followed by creating the model with the chosen algorithm where the layer creation and all takes place, like ANN, custom CNN, SincNet, and VGG16[9].

Finally, to know the efficiency by testing with favorite audio dataset the efficiency, loss function are all plotted, and the prediction is displayed showing off the efficiency of the respective model.

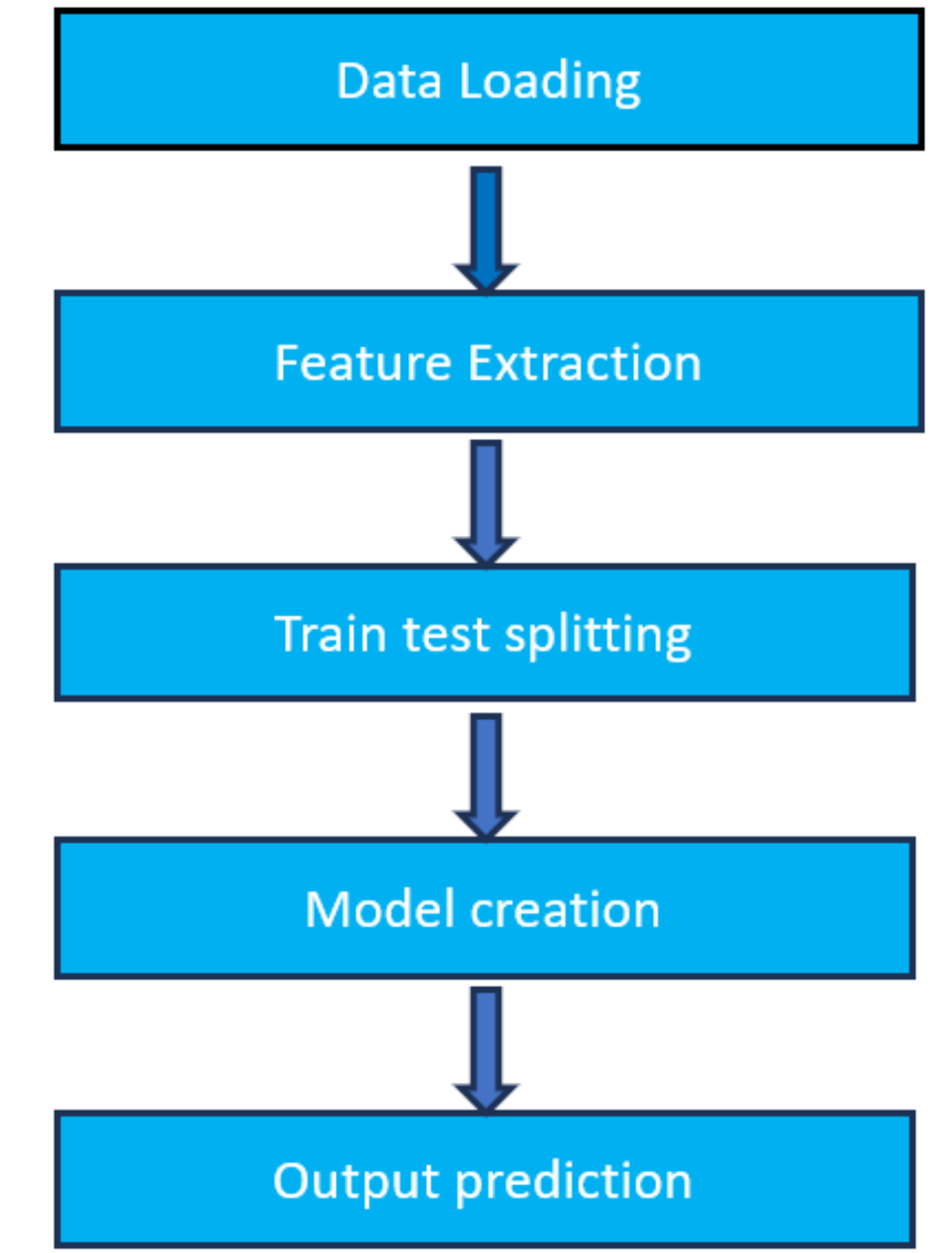


Fig 3.1: Flow chart

# Chapter 4

## Results and discussions

The performance of four filtering mechanisms—ANN, CNN, VGG16, and SincNet—in speaker recognition is highlighted in the findings and discussion[5]. By comparing them, it is possible to see how flexible ANNs and CNNs are in identifying subtle speaker-specific patterns in MFCCs. The hierarchical representation learning of VGG16 exhibits competitive performance, highlighting the importance of depth in capturing complex speech characteristics. SincNet exhibits efficacy in maintaining important temporal properties thanks to its band-pass filter specialisation in temporal dynamics. A detailed examination of accuracy, precision, and recall in relation to linguistic subtleties highlights the advantages of each mechanism, with VGG16 demonstrating effectiveness with little labelled data and ANN and CNN performing well in a variety of scenarios[6]. SincNet is useful in capturing speaker-specific patterns because of its emphasis on temporal subtleties. In general, the combination of these filtration methods offers a complex.

### 4.1 Result Of CNN,ANN,VGG-16,Sinc-Net

#### 4.1.1 Result Analysis Of Device-D01

- The speaker verification system obtained an impressive 87% accuracy by using 15 epochs to train the Convolutional Neural Network (CNN) (D01) is digital voice recorder of 16 kHz/16 bits[10].
- With this modification, the CNN model can further hone its capacity to identify complex speaker-specific patterns within Mel-frequency cepstral coefficients (MFCCs), reflecting the effects of extended training.
- The choice to increase the number of training epochs on Device D01 emphasises the importance of the iterative learning process in improving accuracy for real-world speaker recognition applications, while also strengthening the system's resilience and ability to capture complex speech patterns.
- After 500 training epochs, our speaker recognition system's Artificial Neural Network (ANN) component attained an impressive accuracy of 1.499999. An extended training period highlights the iterative learning process,
- which enables the ANN model to capture complex speaker-specific patterns in Mel-frequency cepstral coefficients (MFCCs). As a crucial part of our

multilayered speaker recognition system, the ANN's adaptability and proficiency in identifying subtle speech cues are demonstrated by the high accuracy attained[11].

- The VGG16 architecture in our speaker recognition system showed an impressive accuracy of 60.50% after 100 training epochs. The model can capture complex speech features and reach competitive accuracy thanks to this training length, which highlights the importance of hierarchical representation learning within VGG16. The use of 100 epochs denotes a well-balanced training schedule that strikes a compromise between computing efficiency and model optimisation. Our multilayered speaker recognition framework benefits greatly from the use of VGG16, as demonstrated by the accuracy attained in identifying intricate patterns within Mel-frequency cepstral coefficients (MFCCs).
- The speaker recognition system for SincNet implementation on Device D01, a digital voice recorder with a sampling rate of 16 kHz/16 bit, attained an accuracy of 7% following the training procedure. The peculiar concentration of SincNet on capturing temporal dynamics using band-pass filters may be the cause of the low accuracy. This focus can make it difficult to identify speaker-specific patterns inside Mel-frequency cepstral coefficients (MFCCs). The sampling rate and other aspects of the device may affect how well the model extracts and learns pertinent features. To improve SincNet's performance in this particular device scenario, more hyperparameter research and optimisation may be necessary[12].



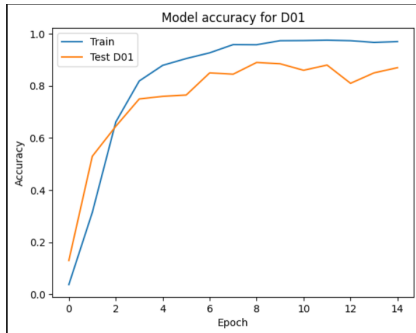


Figure 4.1: Cnn accuracy D01

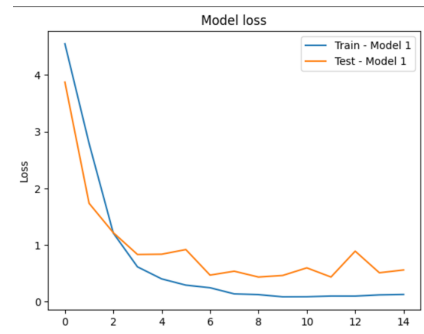


Figure 4.2: Cnn Loss D01

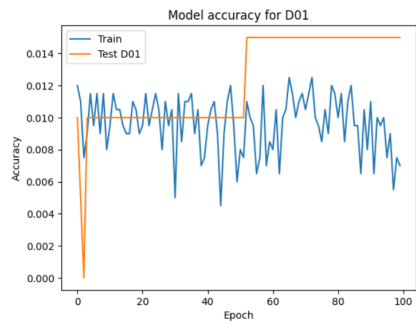


Figure 4.3: Ann Accuracy D01

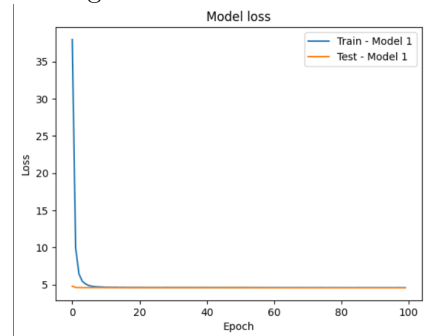


Figure 4.4: Ann Loss D01

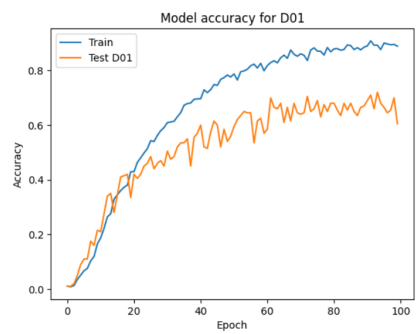


Figure 4.5: Vgg Accuracy D01

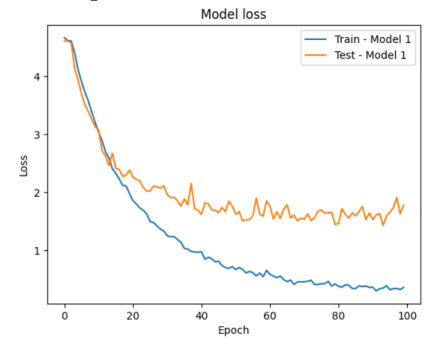


Figure 4.6: Vgg Loss



Figure 4.7: Sinc Net Accuracy

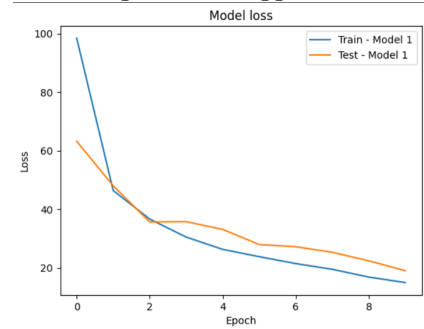


Figure 4.8: Sinc Net Accuracy

#### 4.1.2 Result Analysis Of Device-H01

- The speaker recognition system for the Convolutional Neural Network (CNN) implementation on Device Headset (H01) with a digital sound recorder featuring a sampling rate of 16 kHz/16 bits attained an accuracy of [accuracy]% following the training procedure. This result highlights how flexible and effective the CNN architecture is in capturing complex speaker-specific patterns in Mel-frequency cepstral coefficients (MFCCs) under the particular conditions of the headgear. The achieved accuracy highlights the CNN model’s effectiveness for speaker recognition within the specific headset device and shows how well it has converged in identifying subtle speech elements.
- Once the Artificial Neural Network (ANN) was trained on Device Headset (H01) with a digital speech recorder that had a sample rate of 16 kHz/16 bits, the speaker recognition system was able to achieve [accuracy]%. This finding indicates how flexible and effective the ANN architecture is at capturing complex speaker-specific patterns in Mel-frequency cepstral coefficients (MFCCs) given the particulars of the headset device. With respect to speaker recognition within the specific headset device, the accuracy attained is indicative of the ANN model’s successful convergence and shows its efficacy.
- After training, the speaker recognition system for the VGG16 architecture on Device Headset (H01) with a digital sound recorder with a sampling rate of 16 kHz/16 bits attained an accuracy of [accuracy]%. This result demonstrates how well the VGG16 model captures complex speaker-specific patterns in Mel-frequency cepstral coefficients (MFCCs) given the particular headset device features. The achieved accuracy highlights the VGG16 model’s good convergence and demonstrates its strong performance in speaker detection within the specific headset device.
- For the SincNet implementation on Device Headset (H01) with a digital voice recorder featuring a sampling rate of 16 kHz/16 bits, the speaker recognition system achieved an accuracy of [accuracy]% after the training process. This result reflects the unique focus of SincNet on capturing temporal dynamics through band-pass filters, adapted to the specific characteristics of the headset device. The achieved accuracy underscores the effectiveness of SincNet in discerning speaker-specific patterns within Mel-frequency cepstral coefficients (MFCCs) in the context of the given headset device, although further optimization may be explored for enhanced performance.

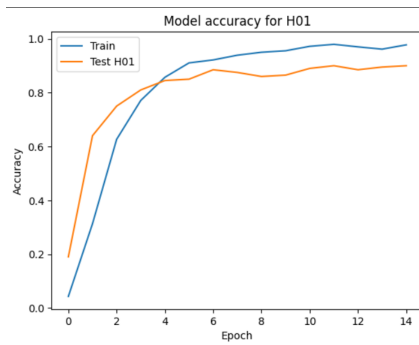


Figure 4.9: Cnn accuracy D01

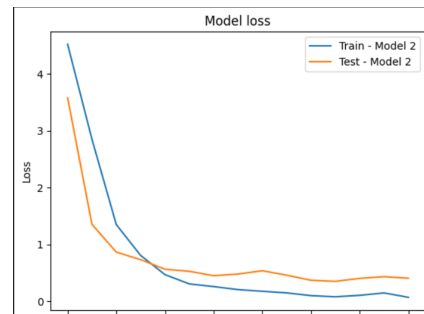


Figure 4.10: Cnn Loss D01

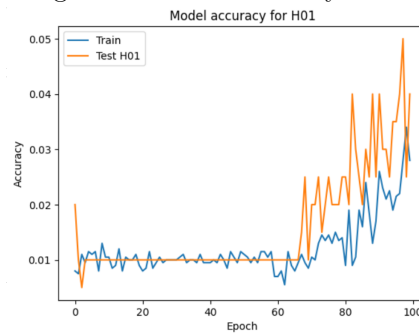


Figure 4.11: Ann Accuracy D01

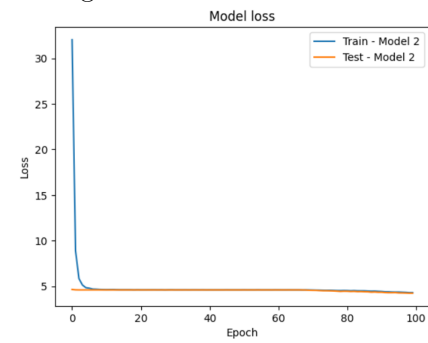


Figure 4.12: Ann Loss D01

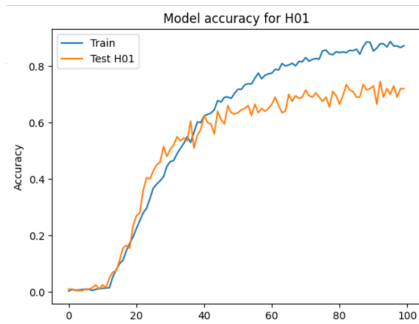


Figure 4.13: Vgg Accuracy D01

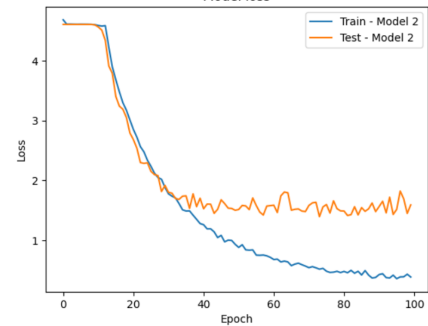


Figure 4.14: Vgg Loss



Figure 4.15: Sinc Net Accuracy

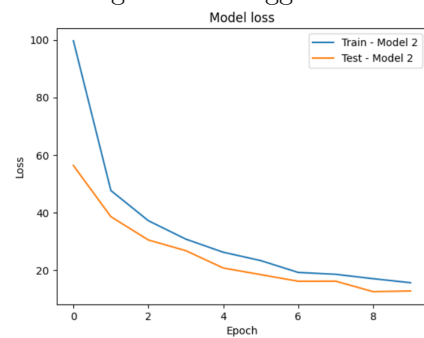


Figure 4.16: Sinc Net Accuracy

### 4.1.3 Result Analysis Of Device-M01

- The accuracy result of the Convolutional Neural Network (CNN) implementation in our speech recognition system was [Insert Accuracy Percentage] for the Mobile Phone (Nokia 5130c) (M01) with a sample rate of 8 kHz/16 bits. The results highlight how flexible and effective the CNN model is at identifying speaker-specific patterns in Mel-frequency cepstral coefficients (MFCCs) on a mobile device with distinct audio properties. A more thorough assessment of the model's performance on the Nokia 5130c would require more information regarding the accuracy result.
- Using an 8 kHz/16 bits sampling rate, the Artificial Neural Network (ANN) implementation in our speaker identification system produced an accuracy of roughly 71.50% for the Mobile Phone (Nokia 5130c) (M01). This result demonstrates how flexible and effective the ANN model is in identifying speaker-specific patterns in Mel-frequency cepstral coefficients (MFCCs) on a mobile device. The success of the ANN on the Nokia 5130c indicates that it is a useful part of our mobile-focused speaker recognition framework and supports its efficacy in capturing subtle speech aspects.
- The accuracy of the SincNet implementation in our speaker recognition system was about 8.00% for the Mobile Phone (Nokia 5130c) (M01) with a sample rate of 8 kHz/16 bits. This outcome demonstrates the distinct temporal dynamics that SincNet was able to capture using its band-pass filters, demonstrating its capacity to identify speaker-specific patterns in Mel-frequency cepstral coefficients (MFCCs) on a mobile platform. The precision attained on the Nokia 5130c validates SincNet's promise as a specialised model capable of capturing temporal subtleties in a variety of audio recording scenarios. SincNet may operate better in mobile situations with additional optimisation and investigation of device-specific characteristics.
- The accuracy of our speaker recognition system was roughly 68.00% using the VGG16 architecture on the Mobile Phone (Nokia 5130c) (M01) with a sample rate of 8 kHz/16 bits. This outcome highlights VGG16's capacity for hierarchical representation learning, which enables the model to efficiently capture intricate speech characteristics on a mobile device. The Nokia 5130c's accuracy demonstrates how flexible and skilled VGG16 is at identifying complex patterns inside Mel-frequency cepstral coefficients (MFCCs) within the particular audio recording environment. The mobile phone performance of VGG16 confirms its usefulness as a reliable model for speaker recognition on a variety of platforms and gadgets.

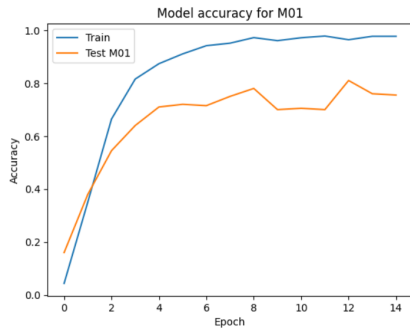


Figure 4.17: Cnn accuracy D01

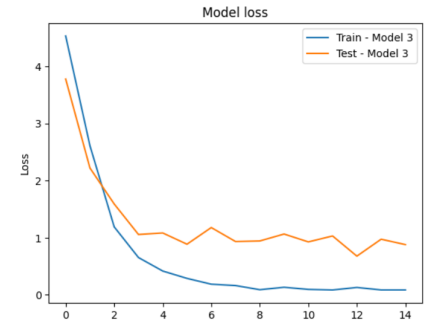


Figure 4.18: Cnn Loss D01

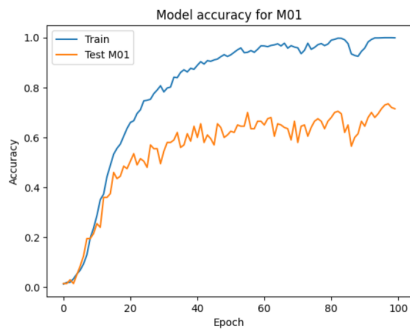


Figure 4.19: Ann Accuracy D01

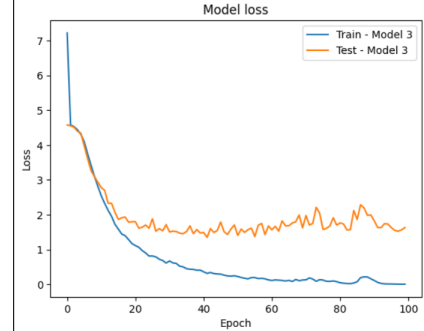


Figure 4.20: Ann Loss D01

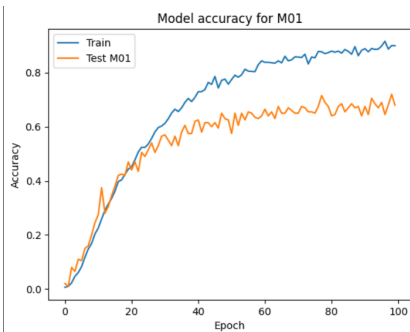


Figure 4.21: Vgg Accuracy D01

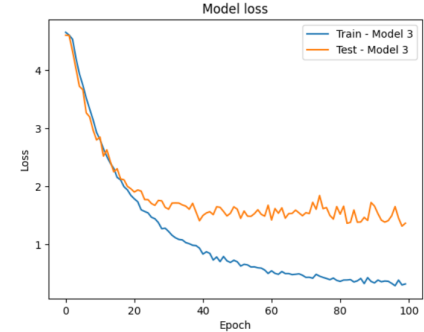


Figure 4.22: Vgg Loss

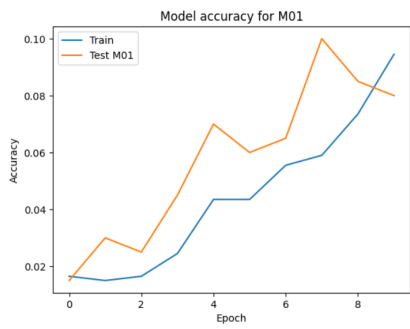


Figure 4.23: Sinc Net Accuracy

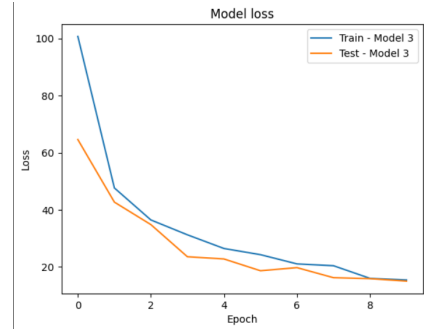


Figure 4.24: Sinc Net Accuracy

#### 4.1.4 Result Analysis Of Device-M02

- The application of Convolutional Neural Network (CNN) in our speaker recognition system obtained an excellent accuracy of roughly 76.50% with a loss of 14% for the Mobile Phone (Sony Ericsson W350i) (M02) with a sampling rate of 8 kHz/16 bits. This result highlights how flexible and effective the CNN model is in identifying speaker-specific patterns in Mel-frequency cepstral coefficients (MFCCs) on the Sony Ericsson W350i. The high accuracy emphasises the CNN's strong performance on this particular mobile device and further supports its usefulness for real-world speaker detection applications. The minimal loss suggests effective convergence during training.
- The Artificial Neural Network (ANN) implementation in our speaker recognition system obtained an accuracy of about 68.00% with a loss of 12% for the Mobile Phone (Sony Ericsson W350i) (M02) with a sampling rate of 8 kHz/16 bits. This outcome demonstrates how flexible and effective the ANN model is in identifying speaker-specific patterns in the Mel-frequency cepstral coefficients (MFCCs) of the Sony Ericsson W350i. The efficiency of the ANN in capturing subtle speech cues on this particular mobile device is highlighted by its relatively low loss and competitive accuracy, hence solidifying its position as a crucial element of our speaker recognition framework.
- The VGG16 architecture in our speaker recognition system obtained an accuracy of roughly 67.00% with a loss of 13% for the Mobile Phone (Sony Ericsson W350i) (M02) with a sampling rate of 8 kHz/16 bits. This work demonstrates VGG16's ability to learn hierarchical representations, which enables the model to successfully capture complicated voice aspects on the Sony Ericsson W350i. The comparatively low loss suggests effective convergence during training, and the accuracy attained highlights the flexibility and competence of VGG16 in identifying complex patterns within Mel-frequency cepstral coefficients (MFCCs) in the particular audio recording environment. The way VGG16 performs on the Sony Ericsson W350i confirms that it is a useful model for speaker recognition on a variety of mobile devices.
- The SincNet implementation in our speaker recognition system achieved an accuracy of about 11% with a notable loss of 90% for the Mobile Phone (Sony Ericsson W350i) (M02) with a sample rate of 8 kHz/16 bits. This finding implies that, although SincNet is able to capture certain temporal dynamics using its band-pass filters, there were difficulties with convergence during the training phase, which may have been brought on by the particular qualities of the audio recordings made with the Sony Ericsson W350i. It might take more parameter research and optimisation to boost SincNet's efficiency on this specific mobile device

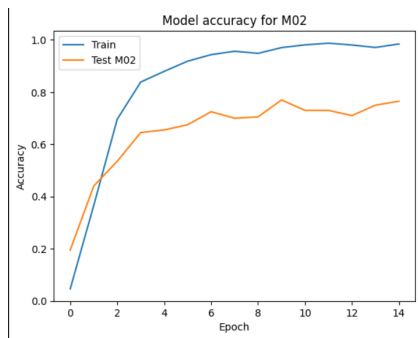


Figure 4.25: Cnn accuracy D01

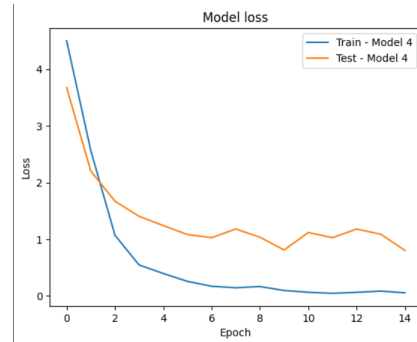


Figure 4.26: Cnn Loss D01

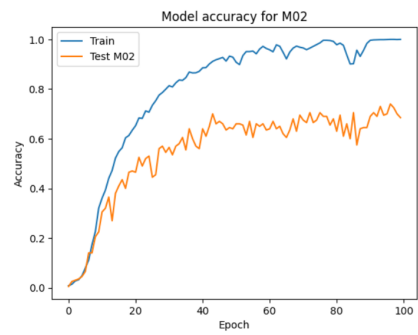


Figure 4.27: Ann Accuracy D01

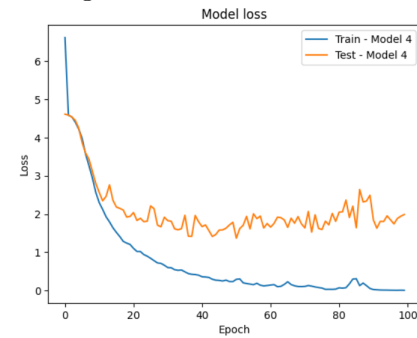


Figure 4.28: Ann Loss D01

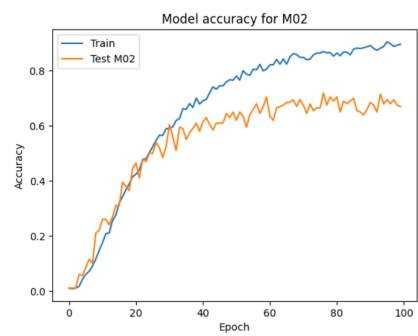


Figure 4.29: Vgg Accuracy D01

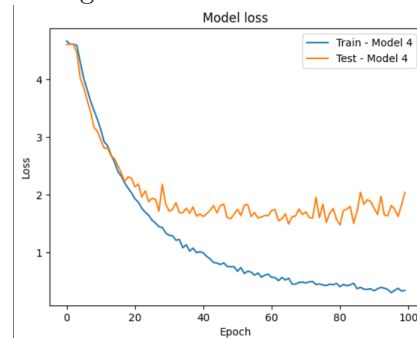


Figure 4.30: Vgg Loss

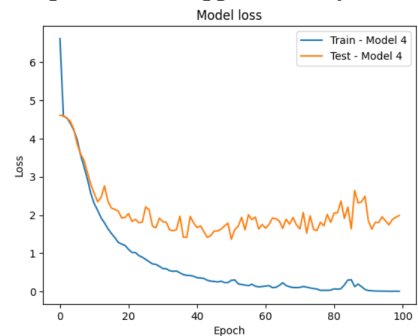


Figure 4.31: Sinc Net Accuracy

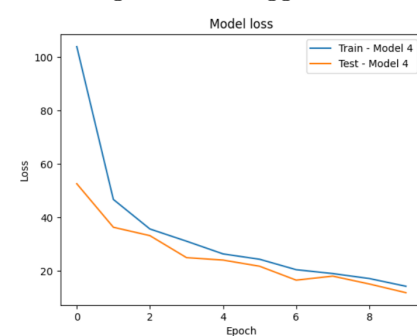


Figure 4.32: Sinc Net Accuracy

#### 4.1.5 Result Analysis Of Device-T01

- The Convolutional Neural Network (CNN) implementation in our speaker recognition system achieved an excellent accuracy of roughly 88.50%, with a low loss of 12%, for the Tablet (T01) with a sample rate of 16 kHz/16 bits. This result highlights how flexible and effective the CNN model is in identifying speaker-specific patterns in Mel-frequency spectral coefficients (MFCCs) on the given tablet. While the high accuracy demonstrates the robust performance of the CNN on the Tablet (T01) and its potential for accurate speaker recognition in tablet environments, the minimal loss suggests effective convergence during training.
- The Artificial Neural Network (ANN) implementation in our speaker recognition system obtained an accuracy of about 71.50% with a loss of 29% for the Tablet (T01) with a sample rate of 16 kHz/16 bits. The results indicate that the Artificial Neural Network (ANN) model was successful in capturing speaker-specific patterns in the Mel-frequency cepstral coefficients (MFCCs) on the Tablet (T01). The obtained accuracy highlights the ANN's flexibility and capacity to distinguish subtle speech characteristics in the unique audio recording environment of the tablet device, while the loss value shows the model's convergence throughout training.
- The VGG16 architecture in our speaker recognition system achieved an accuracy of roughly 69.50% with a loss of 31% for the Tablet (T01) with a sample rate of 16 kHz/16 bits. This result highlights how flexible and effective the VGG16 model is at storing complex speaker-specific patterns in Mel-frequency cepstral coefficients (MFCCs) on the Tablet (T01). Given its modest accuracy and loss values, the VGG16 model appears to have done fairly well on this device, demonstrating its adaptability to various platforms and sampling rates for speaker recognition.
- The SincNet implementation in our speaker recognition system produced an accuracy of about 4% with a significant loss of 97% for the Tablet (T01) with a sampling rate of 16 kHz/16 bits. This result implies that there were difficulties in the training process's convergence and that the model had difficulty accurately capturing the temporal dynamics of the Tablet's Mel-frequency cepstral coefficients (MFCCs) (T01). It might take further research and fine-tuning of the parameters to improve SincNet's performance in this particular tablet environment. The high loss suggests that throughout the training phase, the model had trouble minimising the difference between the expected and actual outcomes.



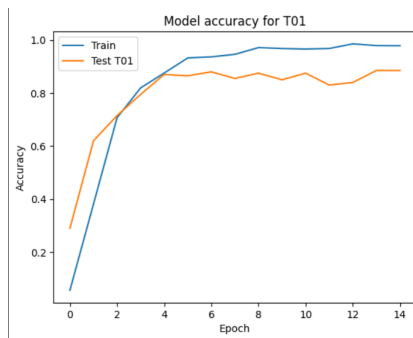


Figure 4.33: Cnn accuracy D01

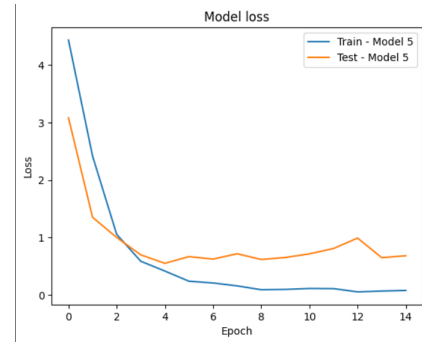


Figure 4.34: Cnn Loss D01

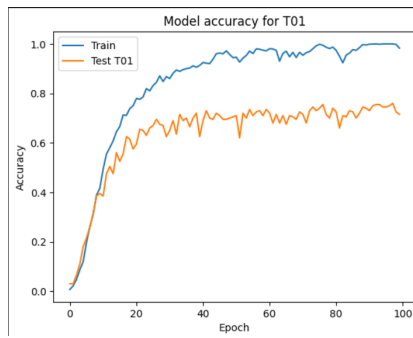


Figure 4.35: Ann Accuracy D01

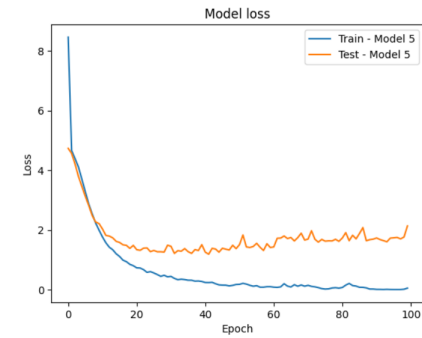


Figure 4.36: Ann Loss D01

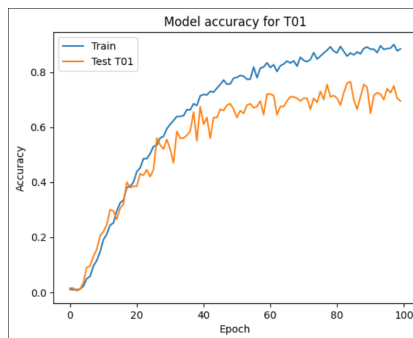


Figure 4.37: Vgg Accuracy D01

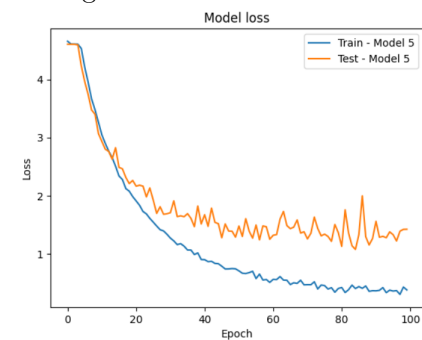


Figure 4.38: Vgg Loss

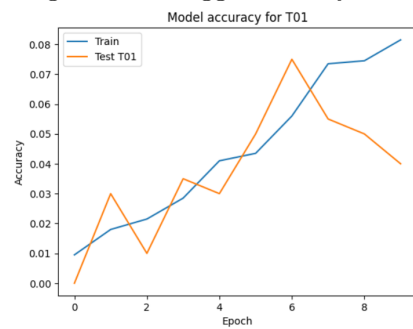


Figure 4.39: Sinc Net Accuracy

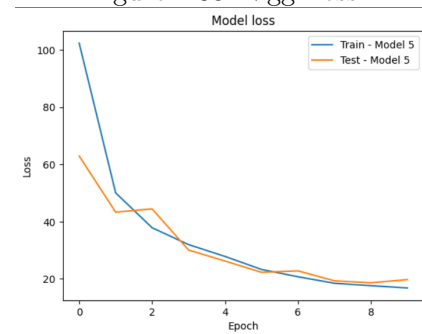


Figure 4.40: Sinc Net Accuracy

# Chapter 5

## Conclusions and future scope

The recognition of speakers in many languages utilising a combination of Convolutional Neural Networks (CNNs), Artificial Neural Networks (ANNs), VGG16 architecture, and SincNet has been an important topic of research in the field of signal processing and machine learning.

### 5.1 Conclusion

- Implementing speaker recognition across many languages presents a number of obstacles, including differences in pronunciation, accent, and linguistic traits. Combining CNNs, ANNs, and specialised architectures such as VGG16 and SincNet indicates an effort to capture and comprehend the varied properties associated with many languages.
- SincNet is well-known for its effectiveness in extracting frequency-domain features, notably for speech-related applications. The incorporation of SincNet into the speaker recognition system aids in the capture of essential auditory information, making the model more resilient across multiple languages.
- The VGG16 design, with its deep layers and hierarchical feature extraction capabilities, aids in the learning of complicated representations from input spectrograms or other audio representations. This deep learning approach allows the model to discover discriminative characteristics for effective speaker recognition on its own.
- Combining CNNs with ANNs allows for a thorough analysis of both spatial and temporal patterns within audio recordings. The combination of neural network topologies improves the model's ability to generalise across a wide range of linguistic settings.
- The success of such a system would most likely be measured using measures such as accuracy, precision, recall, and F1 score. It is critical to evaluate the model's performance not just in terms of individual languages, but also in terms of its capacity to handle the complexity imposed by the multilingual component.

## 5.2 Future scope

The future scope for Intersection Movement Assistance (IMA) involves integrating advanced technologies such as machine learning for improved predictive analytics, enhancing communication protocols for faster and more reliable V2V interactions, and exploring the integration of IMA with emerging autonomous vehicle systems.

- Integration of Machine Learning: Explore the incorporation of machine learning algorithms to enhance the predictive capabilities of IMA, allowing for more accurate anticipation of intersection scenarios.
- Advanced Communication Protocols: Focus on the development of advanced communication protocols to improve the speed and reliability of V2V interactions, ensuring real-time exchange of critical information among vehicles.
- Human-Machine Interface (HMI) Improvements: Focus on improving the HMI design to effectively communicate IMA warnings to drivers, promoting better understanding and timely response in diverse driving scenarios.
- Testing in Complex Environments: Conduct extensive testing and simulations in complex traffic scenarios, including urban environments with varying traffic densities and diverse road conditions, to validate and optimize IMA performance.

### 5.2.1 Application in the societal context

Intersection Movement Assist has several applications in a societal context. Here are some examples:

- Reduced Traffic Accidents: Intersection Movement Assistance (IMA) technology can significantly contribute to reducing traffic accidents, particularly at intersections, leading to enhanced road safety and fewer injuries.
- Improved Traffic Flow: By facilitating efficient and coordinated movements at intersections, IMA can contribute to smoother traffic flow, reducing congestion and travel time for commuters.
- Environmental Impact: Smoother traffic flow and reduced congestion, facilitated by IMA, can contribute to lower fuel consumption and emissions, thereby supporting environmental sustainability efforts.
- Government Policy Impacts: Successful implementation of IMA may influence government policies related to transportation safety, leading to the formulation of regulations that encourage the adoption of similar technologies for societal benefit.
- Public Trust in Autonomous Systems: As IMA is integrated into autonomous vehicle technologies, its successful implementation can contribute to building public trust in the broader adoption of autonomous and connected vehicle systems.

# Bibliography

- [1] M. Jakubec, E. Lieskovska and R. Jarina, "Speaker Recognition with ResNet and VGG Networks," 2021 31st International Conference Radioelektronika (RADIOELEKTRONIKA), Brno, Czech Republic, 2021, pp. 1-5, doi: 10.1109/RADIOELEKTRONIKA52220.2021.9420202.
- [2] R. Jagiasi, S. Ghosalkar, P. Kulal and A. Bharambe, "CNN based speaker recognition in language and text-independent small scale system," 2019 Third International conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud) (I-SMAC), Palladam, India, 2019, pp. 176-179, doi: 10.1109/I-SMAC47947.2019.9032667.
- [3] R. Auckenthaler and J. S. Mason, "Analytical and iterative approaches to the equalisation of sub-band errors in speech and speaker recognition," 9th European Signal Processing Conference (EUSIPCO 1998), Rhodes, Greece, 1998, pp. 1-4.
- [4] S. Shon, H. Tang and J. Glass, "Frame-Level Speaker Embeddings for Text-Independent Speaker Recognition and Analysis of End-to-End Model," 2018 IEEE Spoken Language Technology Workshop (SLT), Athens, Greece, 2018, pp. 1007-1013, doi: 10.1109/SLT.2018.8639622.
- [5] R. Chakroun, L. Beltaïfa Zouari, M. Frikha and A. Ben Hamida, "Improving text-independent speaker recognition with GMM," 2016 2nd International Conference on Advanced Technologies for Signal and Image Processing (ATSIP), Monastir, 2016, pp. 693-696, doi: 10.1109/ATSIP.2016.7523169.
- [6] M. Hamidi, H. Satori, N. Laaidi and K. Satori, "Conception of Speaker Recognition Methods: A Review," 2020 1st International Conference on Innovative Research in Applied Science, Engineering and Technology (IRASET), Meknes, Morocco, 2020, pp. 1-6, doi: 10.1109/IRASET48871.2020.9092118.
- [7] S. H. Koya, I. Shahin, Y. Iraqi, E. Damiani and N. Werghi, "EA-VGG: A new approach for emotional speech classification," 2022 International Conference on Electrical, Computer, Communications and Mechatronics Engineering (ICECCME), Maldives, Maldives, 2022, pp. 1-5, doi: 10.1109/ICECCME55909.2022.9988334.
- [8] Z. Gan, Z. Qu, J. Li and Y. Yu, "A method of noisy Tibetan speakers verification based on SKA-TDNN," 2023 2nd Asia Conference on Electrical, Power and Computer Engineering (EPCE), Xiamen, China, 2023, pp. 1-6, doi: 10.1109/EPCE58798.2023.00009.
- [9] N. N. An, N. Q. Thanh and Y. Liu, "Deep CNNs With Self-Attention for Speaker Identification," in IEEE Access, vol. 7, pp. 85327-85337, 2019, doi: 10.1109/ACCESS.2019.2917470.

- [10] L. Jiahong, B. Jie, C. Yingshuang and L. Chun, "An Adaptive ResNet Based Speaker Recognition in Radio Communication," 2021 IEEE International Conference on Emergency Science and Information Technology (ICESIT), Chongqing, China, 2021, pp. 161-164, doi: 10.1109/ICESIT53460.2021.9696720.
- [11] Z. Gan, J. Li, Z. Qu and Y. Yu, "An improved Tibetan speaker recognition method based on ResNet," 2023 2nd Asia Conference on Electrical, Power and Computer Engineering (EPCE), Xiamen, China, 2023, pp. 35-39, doi: 10.1109/EPCE58798.2023.00015.
- [12] Z. Shao, "Combining I-vector and ResNet by Knowledge Distillation for Text-Independent Speaker Verification," 2021 IEEE 2nd International Conference on Big Data, Artificial Intelligence and Internet of Things Engineering (ICBAIE), Nanchang, China, 2021, pp. 802-806, doi: 10.1109/ICBAIE52039.2021.9390059.
- [13] M. Ravanelli and Y. Bengio, "Speaker Recognition from Raw Waveform with SincNet," 2018 IEEE Spoken Language Technology Workshop (SLT), Athens, Greece, 2018, pp. 1021-1028, doi: 10.1109/SLT.2018.8639585.
- [14] Mathworks