



# UNIVERSIDAD DE INVESTIGACIÓN DE TECNOLOGÍA EXPERIMENTAL YACHAY

Escuela de Ciencias Matemáticas y Computacionales

## **TÍTULO: Computer-Assisted Mispronunciation Detection System for L2 Kichwa Speech**

Trabajo de integración curricular presentado como requisito para la  
obtención del título de Ingeniero en Tecnologías de la Información

### **Autor:**

Velasco Silva Ricardo Isaías

### **Tutor:**

PhD - Fonseca Delgado Rigoberto Salomón

### **Co-Tutor:**

PhD - Morales Navarrete Diego Fabián

Urcuquí, Marzo 2024

# Autoría

Yo, **VELASCO SILVA RICARDO ISAÍAS**, con cédula de identidad 175027344-1, declaro que las ideas, juicios, valoraciones, interpretaciones, consultas bibliográficas, definiciones y conceptualizaciones expuestas en el presente trabajo; así cómo, los procedimientos y herramientas utilizadas en la investigación, son de absoluta responsabilidad de el/la autor/a del trabajo de integración curricular. Así mismo, me acojo a los reglamentos internos de la Universidad de Investigación de Tecnología Experimental Yachay.

Urququí, Marzo 2024.

---

Velasco Silva Ricardo Isaías

CI: 175027344-1

# Autorización de publicación

Yo, **VELASCO SILVA RICARDO ISAÍAS**, con cédula de identidad 175027344-1, cedo a la Universidad de Investigación de Tecnología Experimental Yachay, los derechos de publicación de la presente obra, sin que deba haber un reconocimiento económico por este concepto. Declaro además que el texto del presente trabajo de titulación no podrá ser cedido a ninguna empresa editorial para su publicación u otros fines, sin contar previamente con la autorización escrita de la Universidad.

Asimismo, autorizo a la Universidad que realice la digitalización y publicación de este trabajo de integración curricular en el repositorio virtual, de conformidad a lo dispuesto en el Art. 144 de la Ley Orgánica de Educación

Urcuquí, Marzo 2024.

---

Velasco Silva Ricardo Isaías

CI: 175027344-1

# Resumen

Se ha realizado una evaluación inicial y experimental de un sistema de detección de errores de pronunciación para la lengua kichwa. Se empleó arquitecturas de redes neuronales convolucionales preentrenadas para clasificar espectrogramas de palabras pronunciadas con exactitud y palabras pronunciadas con inexactitud. El modelo inicial, conocido como un modelo basado en características de redes convolucionales, extrae características de las capas totalmente conectadas. A continuación, emplea una técnica de selección de características para separar las características discriminativas de las no discriminativas. Por último, estas características se clasifican mediante un clasificador KNN. El segundo modelo, basado en el aprendizaje por transferencia con redes neuronales convolucionales (CNN), utiliza el conocimiento de las capas convolucionales y adapta la capa clasificadora para la clasificación binaria, distinguiendo entre audios bien pronunciados y mal pronunciados. En cuanto al conjunto de datos utilizado, se construyeron dos conjuntos de datos que se utilizaron en este estudio: un conjunto de datos con palabras en Kichwa y palabras sintéticas, y el mismo pero con palabras sintéticas para el entrenamiento. En conclusión, el método basado en el aprendizaje por transferencia es superior al método basado en características en ambos conjuntos de datos. Concretamente, AlexNet con ajuste de hiperparámetros alcanza 0,90 y 0,92 en la métrica de valor predictivo equilibrado en ambos conjuntos de datos, respectivamente.

**Palabras Clave:**

modelo basado en características, modelo basado en aprendizaje por transferencia, redes neuronales convolucionales, ajuste de hiperparámetros, recopilación de conjuntos de datos, lengua Kichwa

# Abstract

An initial and experimental evaluation of a mispronunciation detection system was developed for the Kichwa language. The study implemented pretrained convolutional neural network architectures to classify spectrograms of accurately pronounced and inaccurately pronounced words. The initial model, known as the CNN feature-based model, extracts features from the fully connected layers. It then employs a feature selection technique to separate discriminative features from non-discriminative ones. Finally, these features are classified using a KNN classifier. The second model, which is based on transfer learning with convolutional neural networks (CNNs), uses the knowledge from convolutional layers and adapts the classifier layer for binary classification, distinguishing between well-pronounced and mispronounced audios. When referring to the used dataset, two datasets were constructed and used in this study: a dataset with Kichwa words and synthetic words, and the same but with synthetic words for training. In conclusion, the CNN transfer learning-based method is superior to the CNN feature-based method in both datasets. Concretely, AlexNet with hyperparameter tuning achieves 0.90 and 0.92 in the balanced predictive value metric in both datasets, respectively.

**Keywords:**

features-based model, transfer learning-based model, convolutional neural networks, hyperparameter tuning, dataset collection, Kichwa language

# Contents

<b>Resumen</b>	<b>iii</b>
<b>Abstract</b>	<b>iv</b>
<b>Contents</b>	<b>v</b>
<b>List of Tables</b>	<b>viii</b>
<b>List of Figures</b>	<b>x</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	1
1.2 Problem statement . . . . .	4
1.3 Objectives . . . . .	4
<b>2 Theoretical Framework</b>	<b>5</b>
2.1 Artificial Intelligence . . . . .	5
2.2 Machine Learning . . . . .	7
2.2.1 Data . . . . .	8
2.2.2 Types . . . . .	10
2.2.3 Tasks . . . . .	10
<b>3 State of the Art</b>	<b>17</b>
<b>4 Methodology</b>	<b>20</b>
4.1 Available Datasets . . . . .	20
4.2 Dataset Construction . . . . .	22
4.3 Dataset Labelling . . . . .	23

4.4	Dataset Splitting and Audio Distribution . . . . .	24
4.5	Preprocessing and Spectrogram Conversion . . . . .	25
4.6	Selected CNN Architectures . . . . .	25
4.6.1	AlexNet . . . . .	25
4.6.2	VGG19 . . . . .	27
4.6.3	ResNet50 . . . . .	27
4.7	Data Augmentation . . . . .	28
4.8	CNN Features-Based Model . . . . .	29
4.8.1	Feature Extraction . . . . .	29
4.8.2	Dataset Discretization . . . . .	30
4.8.3	Statistical-based Feature Selection Methods . . . . .	30
4.8.4	Classification and Hyperparameter Tuning . . . . .	31
4.9	CNN Transfer learning-based Model . . . . .	31
4.9.1	Importing Pre-trained CNN Network . . . . .	31
4.9.2	Modifying Final Layers . . . . .	32
4.9.3	Hyperparameter Tuning . . . . .	32
4.10	Model Performance Metrics . . . . .	32
4.10.1	PPV (Positive Predictive Value) . . . . .	33
4.10.2	NPV (Negative Predictive Value) . . . . .	33
4.10.3	Custom Metric: BPV (Balanced Predictive Value) . . . . .	34
<b>5</b>	<b>Results and Discussion</b>	<b>35</b>
5.1	Dataset Construction . . . . .	35
5.2	Syntetic Data . . . . .	36
5.3	Dataset Labelling . . . . .	37
5.4	Model's Testing . . . . .	37
5.4.1	CNN Features-based Model . . . . .	38
5.4.2	CNN Transfer Learning-based Model . . . . .	39
5.4.3	Observations from Table 5.4 . . . . .	39
<b>6</b>	<b>Conclusions</b>	<b>42</b>
6.1	Observations and Limitations . . . . .	42

<b>Bibliography</b>	<b>45</b>
<b>Appendices</b>	<b>54</b>
<b>A CNN Features-based Model Hyperparameter Tuning and Testing Results: Diagrams and Tables</b>	<b>55</b>
<b>B CNN Transfer Learning-based Model Hyperparameter Tuning and Testing Results: Diagrams and Tables</b>	<b>61</b>
<b>C Exemplary Code</b>	<b>77</b>
C.1 Dataset Splitting . . . . .	77
C.2 Audio Dataset Augmentation, Resampling, and Conversion . . . . .	78
C.3 CNN Features-based Model . . . . .	82
C.4 CNN Transfer Learning-based AlexNet Hyperparameter Tuning . . . . .	91
C.5 CNN Transfer Learning-based AlexNet Testing . . . . .	98



# List of Tables

4.1	Range of Parameters to be Chosen for Hyperparameter Tuning . . . . .	33
5.1	Number of participants in the dataset . . . . .	36
5.2	Distribution of words in Kichwa dataset . . . . .	36
5.3	Syntetic Participants Details . . . . .	37
5.4	Best Models' Performance on Testing Datasets. Each model has been selected according to validation metric. . . . .	41
A.1	CNN Features-based Hyperparameter Tuning on k and Testing in Original Dataset . . . . .	55
A.2	CNN Features-based Hyperparameter Tuning on k and Testing in Synthetic Dataset . . . . .	56
A.3	CNN Features-based Hyperparameter Tuning with Feature Selection and Testing in Original Dataset . . . . .	56
A.4	CNN Features-based Hyperparameter Tuning with Feature Selection and Testing in Synthetic Dataset . . . . .	56
B.1	AlexNet Hyperparameter Tuning on Original Dataset: best results per trial using validation loss as metric. . . . .	62
B.2	AlexNet Hyperparameter Tuning on Original Dataset: best results per trial using BPV as metric. . . . .	63
B.3	AlexNet Hyperparameter Tuning on Synthetic Dataset: best results per trial using validation loss as metric. . . . .	64
B.4	AlexNet Hyperparameter Tuning on Synthetic Dataset: best results per trial using BPV as metric. . . . .	65

B.5	VGG19 Hyperparameter Tuning on Original Dataset: best results per trial using validation loss as metric. . . . .	66
B.6	VGG19 Hyperparameter Tuning on Original Dataset: best results per trial using BPV as metric. . . . .	67
B.7	VGG19 Hyperparameter Tuning on Synthetic Dataset: best results per trial using validation loss as metric. . . . .	68
B.8	VGG19 Hyperparameter Tuning on Synthetic Dataset: best results per trial using BPV as metric. . . . .	69
B.9	ResNet50 Hyperparameter Tuning on Original Dataset: best results per trial using validation loss as metric. . . . .	70
B.10	ResNet50 Hyperparameter Tuning on Original Dataset: best results per trial using BPV as metric. . . . .	71
B.11	ResNet50 Hyperparameter Tuning on Synthetic Dataset: best results per trial using validation loss as metric. . . . .	72
B.12	ResNet50 Hyperparameter Tuning on Synthetic Dataset: best results per trial using BPV as metric. . . . .	73
B.13	AlexNet Mispronunciation Results in Original Dataset . . . . .	74
B.14	AlexNet Mispronunciation Results in Synthetic Dataset . . . . .	74
B.15	VGG19 Mispronunciation Results in Original Dataset . . . . .	75
B.16	VGG19 Mispronunciation Results in Synthetic Dataset . . . . .	75
B.17	ResNet50 Mispronunciation Results in Original Dataset . . . . .	76
B.18	ResNet50 Mispronunciation Results in Synthetic Dataset . . . . .	76

# List of Figures

2.1	Artificial intelligence Diagram [1] . . . . .	6
2.2	Multi-layer Perceptron Diagram with 16 inputs, 2 hidden layers and 3 outputs	14
2.3	CNN Architecture Diagram . . . . .	15
2.4	The transfer learning process involves transferring knowledge from a pre-trained model to a new deep learning model [1] . . . . .	16
4.1	Textgrid with manual annotations from a L2-ARCTIC’s sample audio file [2]. Top to bottom: speech waveform, spectrogram, words, phonemes, error tags, and comments from the annotator. . . . .	20
4.2	Textgrid with manual annotations from a TIMIT’s sample audio file [3]. Top to bottom: speech waveform, spectrogram, phonemes, words . . . . .	21
4.3	Folder directory labelling . . . . .	24
4.4	Spectrogram representation of the word “allku” pronounced by a non-native female speaker . . . . .	25
4.5	AlexNet Architecture Diagram [4] . . . . .	26
4.6	VGG19 Architecture Diagram [5] . . . . .	27
4.7	ResNet50 Architecture Diagram . . . . .	28
A.1	KNN Classifier Hyperparameter Tuning in Original Dataset. . . . .	57
A.2	KNN Classifier Hyperparameter Tuning in Synthetic Dataset. . . . .	58
A.3	Feature Selection Hyperparameter Tuning in Original Dataset . . . . .	59
A.4	Feature Selection Hyperparameter Tuning in Synthetic Dataset . . . . .	60

# Chapter 1

## Introduction

### 1.1 Background

The global population of indigenous people is estimated to be approximately 370 million individuals, residing in more than 90 nations worldwide [6]. Despite constituting just 5% of the global population, this group owns, inhabits, and/or uses 22% of the world's land area. A significant portion of the world's cultural diversity is represented by indigenous people, who speak a total of 7000 languages [7]. However, the enduring effects of historical inequality and marginalization have resulted in indigenous populations representing approximately 19% of the global impoverished demographic. These individuals experience a significantly reduced life expectancy, which has been calculated to be 20 years less than their non-indigenous counterparts worldwide. They encounter challenges such as the absence of official recognition regarding their lands, territories, and natural resources, delayed access to public investments in essential services and infrastructure, and numerous obstacles hindering their full engagement in the formal economy, access to justice, and participation in political processes and decision-making [8].

Alongside the challenges faced by indigenous communities, their languages suffer from similar issues. Many languages are currently facing the imminent threat of extinction, and various factors contribute to this pressing concern [9]:

- **Implementation of state policy.** Several governmental institutions have implemented initiatives with the objective of eradicating indigenous languages with legal regulations that prohibit their usage. As a result, numerous nations deny the exis-

tence of indigenous populations and classify their languages as mere dialects, thus undervaluing them, devaluing them compared to national languages.

- **Lack of support and attention.** Because of their minority status, indigenous populations frequently encounter a lack of recognition and support from governmental initiatives aimed at preserving languages, resulting in the marginalization of their linguistic heritage. Furthermore, even though there may be initiatives, they may not be well executed and consequently be overlooked. In the Philippines, the government has started to implement mother-tongue-based instruction in the classrooms [9]. Nonetheless, due to the lack of available resources, including qualified teachers and appropriate learning materials, the effective instruction of indigenous children in their respective mother tongues may be hindered, and consequently, individuals acquire proficiency in another language rather than their own, ultimately leading to the deterioration of their native language skills.
- **Obsolescence.** Due to a long history of discriminatory policies and practices, a significant number of indigenous parents have made the decision to teach and speak in current languages with their children [9]. This is done with the intention of fostering an environment that is beneficial to their children's development in a globalized society. In most cases, indigenous children are currently unable to take part in conversations with their grandparents, who commonly use their mother tongue to communicate with each other, inevitably leading to the loss of oral traditions and expressions.

Nevertheless, the increasing global recognition of indigenous knowledge systems has restored optimism concerning the preservation and growth of indigenous languages in both oral and written formats. Numerous indigenous communities have already implemented autonomous mechanisms aimed at restoring their respective languages. Moreover, legislative provisions and measures have been proposed. For instance, in Ecuador, the following articles support the right to education for indigenous people:

- "Article 29 of the Constitution declares that 'the State shall guarantee the freedom of teaching, academic freedom in higher education, and the right of individuals to learn in their own language and cultural environment...'" [10]

- "Article 57, number 14, of the Constitution agrees that the State must 'develop, strengthen and enhance the system of bilingual intercultural education, on the basis of quality criteria, from early childhood stimulation to the highest level, in accordance with cultural diversity, for the safeguarding and preserving identities in harmony with teaching and learning methodologies....'" [10]

These articles guarantee the rights of communities, peoples, and nationalities, emphasizing the state's obligation to integrate education with an intercultural environment [10]. By doing so, Ecuador acknowledges the significance of fostering inclusivity and demonstrating respect for the diverse cultural backgrounds of its communities.

In the end, this legislative provision resulted in Ecuador having 1,834 bilingual intercultural educational institutions that educated 156,027 students in 2015–2016, as well as 13 educational institutions designated as guardians of the language and staffed by teachers who speak one of the indigenous languages [6]. Similarly, there are other countries that have implemented similar legislation, such as Argentina, Australia, Bolivia, Brazil, Chile, Guatemala, Honduras, Iran, Mexico, New Zealand, Norway, Sweden, and Venezuela [6].

Unfortunately, indigenous peoples must still overcome several challenges before breaking free from marginalization. One-third of the world's impoverished population are indigenous, and in many countries, indigenous rights laws conflict with laws pertaining to agriculture, land, conservation, forestry, mining and other sectors [9].

7% of the indigenous people in Ecuador use Kichwa and quite frequently it has been reduced to a peasant language [11]. In order to contribute to the preservation of the Kichwa language, Yachay Tech University proposed a community service project called "Kichwa Yachay App" whose head chairman is Sophia L. Cadoux from the Yachay Tech English Department. Formerly, this was a software engineering project aimed at developing a mobile application for learning and practicing the Kichwa language. As a result, all the students gave a presentation presented prototypes of their applications. For one of the projects, there was an open thesis statement to develop a computer-assisted mispronunciation system. The development of such system will provide valuable insights into the capabilities and limitations of the developed models and the developed data, allowing for further improvements and advancements in the field of language learning for Kichwa.

## 1.2 Problem statement

It is important to know that recently, the adaptability of Computer-Aided Language Learning (CALL) has drawn a lot of attention because it lets students improve their language skills at their own pace. To achieve this cause, it is necessary to use CALL technologies. By definition, CALL involves using computers and computer-based resources like the Internet, to present, reinforce, and assess learning. Specifically, for this work, we focus on Computer-Aided Pronunciation Training (CAPT), a subset of CALL that focuses on finding and correcting pronunciation mistakes made by non-native Kichwa language learners [12].

Deep learning technology is considered one of the hot topics within the areas of machine learning, artificial intelligence, data science, and analytics due to its learning capabilities from the given data [1]. The problem begins with the scarcity of data for the Kichwa language (see Section 4.1), to which two models for mispronunciation detection based on convolutional neural networks and spectrograms were developed. In addition, two Kichwa datasets were collected and proposed for the sake of the project and future research. These models will be subjected to testing and measurement using various evaluation metrics, such as positive and negative predictive values, to evaluate their performance.

## 1.3 Objectives

- Generate two datasets of Kichwa words from both native and non-native speakers.
- Implement two models for Kichwa mispronunciation detection by means of convolutional neural network architectures (AlexNet, VGG19, ResNet): CNN features-based and transfer learning-based models.
- Compare the performance of the proposed models with adequate metrics: positive, negative and balanced predictive value.
- Evaluate how the model performs whenever a word is well-pronounced or mispronounced, without restricting only to evaluating the accuracy of the model.

# Chapter 2

## Theoretical Framework

### 2.1 Artificial Intelligence

The renowned author Karel Capek first used the term “robot” in his play “R.U.R.” (Rossum’s Universal Robots) in 1921, and it originally meant a factory of biomechanical machines used as forced labor [13]. During the mid-20th century, Isaac Asimov popularized the term “robot” through a compilation of contemporary science-fiction short stories [13]. Even though the word was slightly new, during the Renaissance period, Leonardo da Vinci performed a comprehensive study of human anatomy with the purpose of guiding the design of his humanoid robot, Leonardo’s robot. It was a mechanized knight automaton capable of assuming many postures, including standing, sitting, and gesturing with its arms, as well as articulating its head and jaw by using a combination of pulleys and cables [14]. Subsequently, the term “Artificial Intelligence”, derived from the idea of a robot, appeared years later. John McCarthy coined the term “artificial intelligence” (AI) in 1955, defining it as “the science and engineering of making intelligent machines”. He was very influential in the early development of AI. With his colleagues, he founded the field of AI in 1956 at a Dartmouth College conference on artificial intelligence. The conference gave birth to what developed into a new interdisciplinary research area. It provided an intellectual framework for all subsequent computer research and development efforts [15].

Given the general historical background, Artificial Intelligence (AI) focuses on the replication of cognitive abilities to address real-world challenges and the development of systems capable of learning and reasoning in a manner akin to human cognition [16]. More and



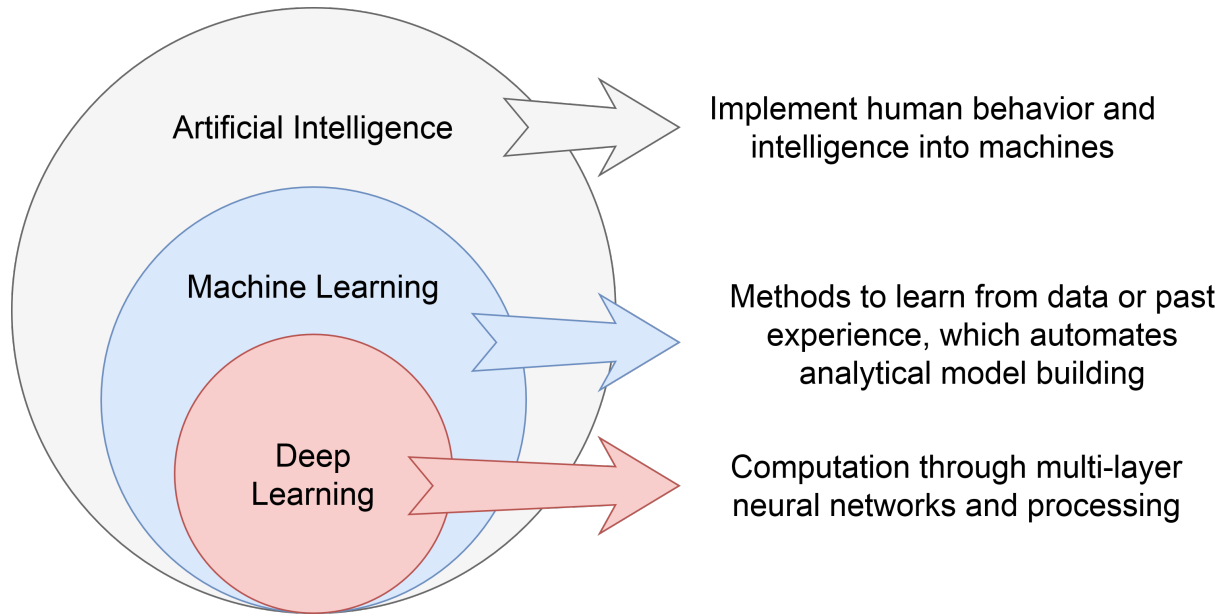


Figure 2.1: Artificial intelligence Diagram [1]

more studies show how important it is to give artificial intelligence systems the ability to build causal models of the world, giving explanation and understanding more weight than just solving pattern recognition problems [17]. Finally, the branches of artificial intelligence can be directly and informally represented in Figure 2.1.

With the advancement of technical elements such as data, algorithms, and computing capabilities reaching a relatively mature stage, artificial intelligence (AI) has begun to efficiently address problems and generate real economic advantages [18]. Thus, it is not surprising that AI is currently being applied in [19]:

- Automotive Industry.
- Financial Markets.
- Health Industry.
- Retailing Industry.
- Media Industry.
- Smart Payment Systems.
- Smart Homes

## 2.2 Machine Learning

The formulation of a singular, comprehensive definition for machine learning (ML) is challenging due to its broad scope, which covers several methodologies derived from computer science and multivariate statistics [20]. However, machine learning is precisely defined as "a set of methods that can automatically detect patterns in data and then use the uncovered patterns to predict future data or to perform other kinds of decision-making under uncertainty (such as planning how to collect more data!)" [21]. Given this definition, machine learning is defined as a subfield of artificial intelligence and a group of methods and algorithms with the capacity to learn continuously from accumulated or past experience and hone their skills through knowledge extraction from data without being specifically programmed. On top of that, it is the most famous and modern technology in the fourth industrial revolution [22], [23]. For example, in a real-world scenario where ML is used, in the context of a sophisticated system like a self-driving vehicle, machine learning algorithms must process extensive sensor data and utilize their acquired knowledge to identify potential hazards and subsequently choose appropriate control actions [24].

At the moment, the discipline of machine learning is structured around three main areas of research [25]:

- **Task-Oriented Studies.** It employs the engineering approach, which refers to the process of designing and evaluating learning systems with the aim of enhancing performance in a certain set of tasks. This involves the development and analysis of these systems.
- **Cognitive Simulation.** The study and computational modeling of human learning processes.
- **Theoretical Analysis.** The present study focuses on the theoretical investigation of the expansive range of potential learning methods and algorithms, absent any specific application domain.

When considering applications, machine learning is utilized across a diverse range of domains, including robotics, virtual personal assistants such as Google, computer games, pattern recognition, natural language processing, data mining, traffic prediction, online

transportation networks, product recommendation systems, share market prediction, medical diagnosis, online fraud detection, agriculture advisory services, search engine result optimization like Google's search engine, chatbots for online customer support, e-mail spam filtering, crime prediction through video surveillance systems, and face recognition in social media [26].

In order to properly summarize the previously mentioned field of inquiry, it is necessary to follow the appropriate structure suggested by [27], [28], [20] and [29]. Thus, in general, the following sections are going to consist of data, types of algorithms, and tasks related to machine learning.

### 2.2.1 Data

Because the machine learns by using specific data sets for training and then applying algorithms to those data sets, the availability of data is the main factor in the development of machine learning models in practical applications.

In the real world, data may appear in diverse formats, including structured, semi-structured, and unstructured forms [30]:

- **Structured.** The data exhibits a clearly defined structure and follows a data model that belongs to a standardized sequence, thereby enabling efficient organization and accessibility machine learning algorithms. Structured data is commonly stored in well-defined schemes, such as relational databases, where it is organized in a tabular fashion.
- **Semi-Structured.** Semi-structured data, in contrast to the structured data discussed earlier, is not often stored within a relational database. However, it possesses particular organizational characteristics that make its analysis feasible. Semi-structured data incorporates several forms, such as HTML, XML, JSON documents, and NoSQL databases.
- **Unstructured.** Unstructured data lacks a predetermined format or organizational structure, thereby making it considerably more challenging to gather, handle, and analyze. This type of data mostly includes textual and multimedia content. There

are various forms of unstructured data, such as sensor data, emails, blog entries, wikis, word processing documents, photographs, presentations, web pages, and other business documents, can be classified as examples of unstructured data.

Furthermore, since datasets are made up of data objects (e.g. rows of the dataset, observations) and attributes/features (e.g. columns of the dataset), it is important to be clear with which types of attributes we are dealing with. Overall, there are four types of attributes [31]:

- **Binary Attributes.** A binary attribute represents a nominal attribute that includes merely two categories or states, namely 0 or 1. In this context, 0 signifies the absence of the attribute, while 1 denotes its presence. Binary attributes are commonly known as booleans when the two states they represent correspond to true and false.
- **Nominal Attributes.** The values associated with a nominal attribute consist of symbolic representations or designations of entities. Values that represent nominal attributes, also referred to as categorical attributes, correspond to various categories, codes, or states. The values lack a recognizable or significant order. If *haircolor* is an attribute, its possible nominal values may be black, brown, yellow, or blonde. Besides, we can assign a code to such values, for instance, 0 for black, 1 for brown, and so on.
- **Ordinal Attributes.** An ordinal attribute refers to an attribute that contains potential values that indicate some sort of order or ranking, but the exact magnitude between consecutive values remains unknown. If *drinksize* is an attribute, its possible ordinal values may be small, medium, and big. These values follow a sequence by nature, but we cannot tell how much bigger “small” is than “medium.”
- **Numerical Attributes.** It is possible to classify a numerical attribute as quantitative, indicating that it is a measurable quantity that has an integer or real value as its representation.

Machine learning classification tasks often categorize attributes as discrete or continuous. In this context, discrete attributes have a finite or countably infinite set of values;

otherwise, continuous attributes can be represented as floating-point numbers [31]. In machine learning, it is essential to determine whether the attributes are discrete or continuous, as it affects the choice of appropriate algorithms and techniques.

### 2.2.2 Types

Algorithms for machine learning fall into the following categories [32]:

- **Supervised learning.** In a commonplace supervised situation, a teacher or supervisor gives the agent a precise error measure that is equivalent to output values. When referring to algorithms, actual methods use a training set of input-output pairs to provide this function. Based on this knowledge, the agent may adjust its settings to lower its minimize loss function. If the method is flexible and the data components are coherent, accuracy improves, and the difference between predicted and expected values approaches zero after each iteration. The applications of supervised learning can be in the fields of predictive analysis, spam detection, pattern detection, natural language processing, and sentiment analysis, among others.
- **Unsupervised learning.** This approach relies on the lack of a supervisor and, consequently, the absence of absolute error measures. Cluster analysis is a valuable tool for grouping elements based on their similarity or distance measure. In fact, unsupervised learning provides an implicit descriptive analysis of the data in order to obtain a complete understanding of the dataset. This approach can be applied to object segmentation, similarity detection, automatic labeling, and recommendation engines.

### 2.2.3 Tasks

#### Classification

Classification is a machine learning task that is employed to classify data into two or more distinct classes. The classification of potato diseases based on deep learning and computer vision methods has the potential to serve as an interesting instance of multi-class classification [33]. From the last example, it can be inferred that there are commonplace types of classification problems based on the label assigned to each instance.

- **Binary Classification.** The term "binary classification" refers to classification tasks in which there are two class labels that can be expressed by either numerical values (0 or 1) or textual descriptors (true or false).
- **Multi-Class Classification.** Multi-class classification involves the classification of more than two classes. In the context of a given set of classes, examples are categorized as members of a specific class.

For instance, some popular classification algorithms used are Naive Bayes, logistic regression, K-nearest neighbors (KNN), support vector machines, decision trees, random forests, adaptive boosting (AdaBoost), extreme gradient boosting (XGBoost), etc [27].

## Dimensionality Reduction

The objective of dimensionality reduction is to address the issue of high dimensionality in data, also referred to as the curse of dimensionality. The primary purpose is to enhance the quality of data by lowering its complexity [34]. In the fields of machine learning and data science, processing high-dimensional data poses a significant challenge that both researchers and application developers must address. Dimensionality reduction plays a crucial role as it encourages improved human interpretations, reduces computational costs, and decreases problems such as overfitting and redundancy by simplifying models [27]. Dimensionality reduction can be categorized into feature selection and feature extraction. The main difference between feature selection and feature extraction lies in their respective approaches: feature selection involves keeping a subset of the original features [35], whereas feature extraction involves creating whole new features [36]. More precisely, in the former, discriminative features are chosen to create a new subset of them, which allows us to delete irrelevant or redundant features present in the data. The selection of an appropriate and effective subset of features inside a given problem domain has the potential to reduce the issue of overfitting by simplifying and generalizing the model, hence enhancing its accuracy [35]. In the latter, a better understanding of the data is provided, which results in enhanced predictive accuracy and a reduction in computing costs or training duration.

## Deep Learning

The Fourth Industrial Revolution, also known as Industry 4.0, primarily centers around technology-driven automation and intelligent systems, which are applied in various areas, such as smart healthcare, business intelligence, smart cities, and cybersecurity intelligence [37]. Due to its excellent learning capabilities from historical data, DL techniques can play a key role in building intelligent data-driven systems that meet the needs of today and potentially alter the world and the way people live in it through their automation potential and their ability to learn from experience [1].

By definition, deep learning is a branch of machine learning that uses hierarchical structures to try to find complex abstractions or architectures in high-dimensional data [1], [38]. This branch of machine learning has become more popular and has been used in many areas, such as semantic parsing, transfer learning, natural language processing, computer vision, and many more. Furthermore, the current growth of DL can be attributed to three key factors: the notable advancements in chip processing capabilities, including GPU processing and Moore's Law; the substantial reduction in hardware construction costs; and the enormous progress made in machine learning algorithms [29]. Surprisingly, there are already many areas where deep learning has exceeded human-level capability and performance, e.g., predicting movie ratings, decisions to approve loan applications, and estimating car delivery times [39].

In order to employ deep learning techniques, it is crucial to consider the following properties and potential requirements [1]:

- **Amount of Data.** Deep learning relies on a substantial volume of data to construct a data-centric model for a specific problem domain. Deep learning algorithms tend to exhibit suboptimal performance when the dataset size is small [40].
- **Hardware Dependency.** Deep learning algorithms require significant computational resources when training models with extensive datasets. Due to the increased computational demands, GPUs are commonly employed to boost running time efficiency, as they offer significant advantages over CPUs. GPU hardware is essential for effective deep learning training. Accordingly, deep learning algorithms require

high-performance computing machines, such as the HPC CEDIA Cluster, to be run optimally [41].

- **Feature Engineering Efficiency.** Feature engineering involves extracting features from raw data by applying domain knowledge. One key difference between DL and other machine learning methods lies in the focus on extracting high-level features directly from data [27], [42]. Fortunately, deep learning reduces the time and effort needed to create a feature extractor for each problem.
- **Training and Execution Time.** Deep learning algorithms typically require a substantial amount of time for training due to the high number of parameters involved. Deep learning models can require several weeks to complete a training session, while training with machine learning algorithms typically takes significantly less time, ranging from seconds to hours [43], [41].
- **Blackbox Model.** Interpretability is a crucial consideration when comparing DL with ML. The process of obtaining a DL result can be challenging to explain due to its "black-box" nature.

We can generally specify several technologies used in deep learning and intended to be employed in this work [44], [1]:

- **Multi-layer Perceptron (MLP).** Multi-layer Perceptron (MLP) is a supervised learning method [45] and the foundational structure of deep neural networks (DNN). This type of architecture falls under the category of feedforward artificial neural networks (ANNs). A standard multi-layer perceptron (MLP) is a fully connected neural network consisting of an input layer for getting input data, an output layer for making decisions or predictions, and one or more hidden layers that serve as the network's computational engine (See Figure 2.2) [31], [46]. The MLP network's output is determined by employing various activation functions, also referred to as transfer functions, including ReLU (Rectified Linear Unit), Tanh, Sigmoid, and Softmax [45], [47]. The most commonly used algorithm for training MLP is "backpropagation", which is a supervised learning technique and considered the fundamental component of a neural network [31]. Different optimization approaches, including Stochastic



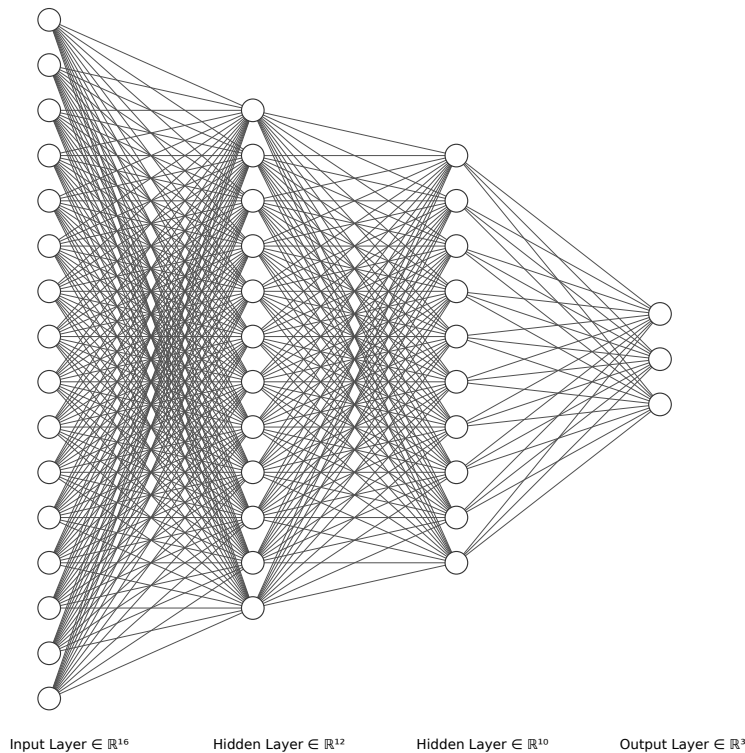


Figure 2.2: Multi-layer Perceptron Diagram with 16 inputs, 2 hidden layers and 3 outputs

Gradient Descent (SGD), Limited Memory BFGS (L-BFGS), and Adaptive Moment Estimation (Adam), are utilized in the training process. The MLP model needs hyperparameter tuning or adjusting the parameters, including the number of hidden layers, neurons, and iterations, which can result in a computationally demanding solution for complex models. MLP has the advantage of learning non-linear models in real-time or online through partial fit [45].

- **Convolutional Neural Networks (CNN).** CNN are widely used deep learning architectures that learn directly from input data without requiring human feature extraction. Each layer in a CNN considers optimal parameters to produce a meaningful output and also reduces the complexity of the model, thus improving the architecture of conventional ANNs, such as regularized MLP networks. Moreover, on top of that improvement, CNNs employ a 'dropout' technique [48], which involves randomly dropping out a certain percentage of nodes during training, which helps prevent the model from memorizing the training data and generalizing better, to address the issue of overfitting. CNNs are commonly used in various fields such as

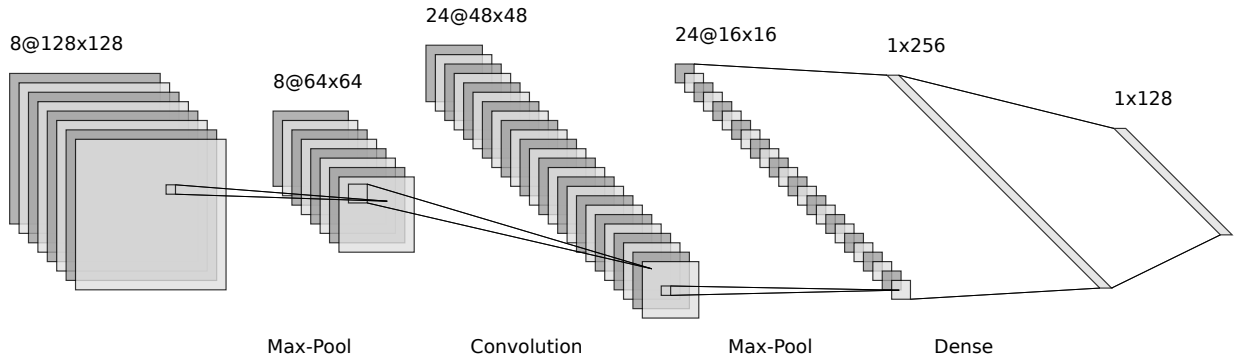


Figure 2.3: CNN Architecture Diagram

visual recognition, medical image analysis, image segmentation, and natural language processing since they are intended to deal with 2D shapes such as images (See Figure 2.3) [49], [47]. The autonomous identification of crucial input features without human intervention enhances the potency of this network in comparison to traditional networks. Various CNN variants have been developed in the field, such as VGG [50], AlexNet [51] and ResNet [52].

- **Deep Transfer Learning (DTL).** Transfer learning is a method that leverages pre-existing model knowledge to efficiently address a new task, requiring minimal training or fine-tuning. Training DL models requires a large amount of data. However, this can present a challenge for specific applications, such as medical image analysis. The collection and annotation of large datasets of patient data are not only challenging and expensive but also require significant computational resources for model training. Hopefully, transfer learning could help address this issue by transferring knowledge from a pretrained DL model to a specific or desired model to be trained. Because of its ability to successfully train deep neural networks with minimal input data, transfer learning has become a hot topic in the field of DL [53]. The process can be visualized in Figure 2.4.

There are many other deep learning methods that deserve mention, such as Recurrent Neural Networks (RNNs), Generative Adversarial Networks (GANs), Auto-Encoders (AEs), Self-Organizing Maps (SOMs), Deep Belief Networks (DBNs), Hybrid Deep Neural Networks, and Deep Reinforcement Learning, all of which are reviewed in detail in [1], but not in this work since those architectures aren't used at all.

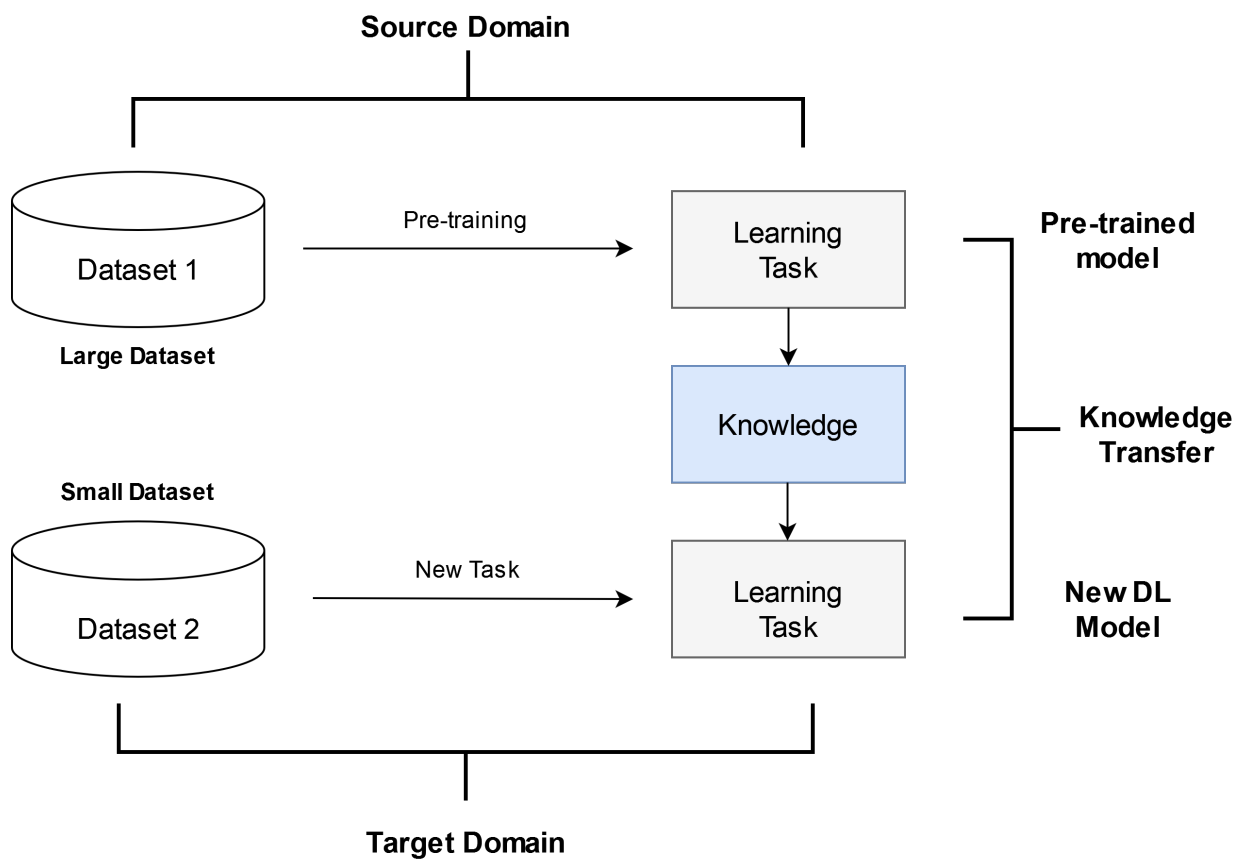


Figure 2.4: The transfer learning process involves transferring knowledge from a pre-trained model to a new deep learning model [1]

# Chapter 3

## State of the Art

There is no official state-of-art document for mispronunciation detection system. However, in [54], a literature review detailed three main techniques: posterior probability-based methods, classifier-based methods, and deep belief network-based methods. Probability-based techniques have the capability to identify the quality of pronunciation; yet, these scoring computations alone are insufficient to differentiate the type of error and accurately pinpoint its position. To solve this, classifier-based methods are used [54]. These methods can be summarized as follows:

- **Posterior-Probability-based Methods**

- The Goodness of Pronunciation (GOP) score was introduced by Witt and Young [55], who based it on a standardization of log-probability values by the duration of individual phone sections. A criterion was established for evaluating the accuracy of pronunciation for each phone.
- In their study, Franco et al. [56] developed three recognition models that used data from different levels of nativeness. The researchers also took into account the log-likelihood-based scores in order to determine the proportionality of these scores.
- The application of acoustic-phonetic parameters, such as log root mean square, energy, and zero-crossing rate, for the recognition of velar fricative and velar plosive was examined by Strik et al. [57].

- Minematsu et al. [58] introduced a widely recognized acoustic framework in speech that effectively eliminates non-semantic information.
- Wang and Lee [59] additionally incorporated GOP scores with error pattern identifiers to enhance the performance of identifying errors within the gathering of students from 36 nations learning Mandarin Chinese.
- Zhang et al. [60] presented scaled log-posterior probability (SLPP) and weighted phone SLPP to enhance measures of pronunciation quality.

- **Classifier-based Methods**

- The clusters of pronunciation rules were produced by Ito et al. [61], who also established a threshold for each cluster. The researchers proposed a clustering technique that is based on decision-making in order to improve the accuracy of error detection.
- The features obtained by Georgoulas et al. [62] by the utilization of a discrete wavelet transform. Wavelet analysis is a distinctive signal processing technique that uses a support vector machine (SVM) classification algorithm to identify instances of mispronunciation.
- In their work, Strik et al. [57] conducted a comparative analysis of four different methods, specifically GOP, decision tree, LDA-APF, and LDA-MFCC. The results of this study show that the LDA-APF and LDA-MFCC techniques did a better job of classifying the velar fricative /x/ and the velar plosive /k/ than the GOP scores and decision tree.
- Amdal et al. [63] conducted a study in which they made distinctions between short and long vowels in speech. The researchers employed acoustic-phonetic characteristics and conducted training on a linear discriminant analysis (LDA) classifier.
- Li et al. [64] introduced a detection strategy that is based on GLDS-SVM (Generalized Linear Discriminant Sequence-Support Vector Machine) and has been proven to be effective. The researchers integrated the GLDS-SVM method with

the UBM-GMM (universal background model) framework in order to improve performance.

- In their study, Wei et al. [65] employed log-likelihood ratios obtained from acoustic models as features for a support vector machine (SVM) classifier. Multiple acoustic models were used for each phone in order to determine the variety of pronunciation variations of that phone across various degrees of capability. This method facilitated the framework’s achievement of superior performance compared to conventional likelihood-based strategies.
- In their study, Maqsood et al. [66] employed the sequential floating-forward selection method to identify the discriminative features for the purpose of detecting Arabic mispronunciations. Additionally, a technique based on grouping was presented by the researchers to identify distinctive properties of Arabic phonemes.

- **Deep Learning**

- Nazir et al. and Akhtar et al. in [54] and [67] propose similar mispronunciation detection systems based on deep learning and machine learning. Most of them work by extracting features from spectrograms. In the former, it proposes hand-crafted features-based, CNN features-based, and transfer learning-based models subjected to feature selection and feature extraction (depending on the case). In the latter, it employs the same structure. The main differences between these articles are the different locations from which features are extracted.

# Chapter 4

## Methodology

### 4.1 Available Datasets

There are several classical datasets used for state-of-the-art mispronunciation system detection in Chapter 3. An example of big speech corpora designed for research purposes are L2-ARCTIC and TIMIT for English. Every speech corpora implements Textgrid as a tool for segmentation and labeling of the phonemes, words, and errors in an audio file as in Figure 4.1 and 4.2.

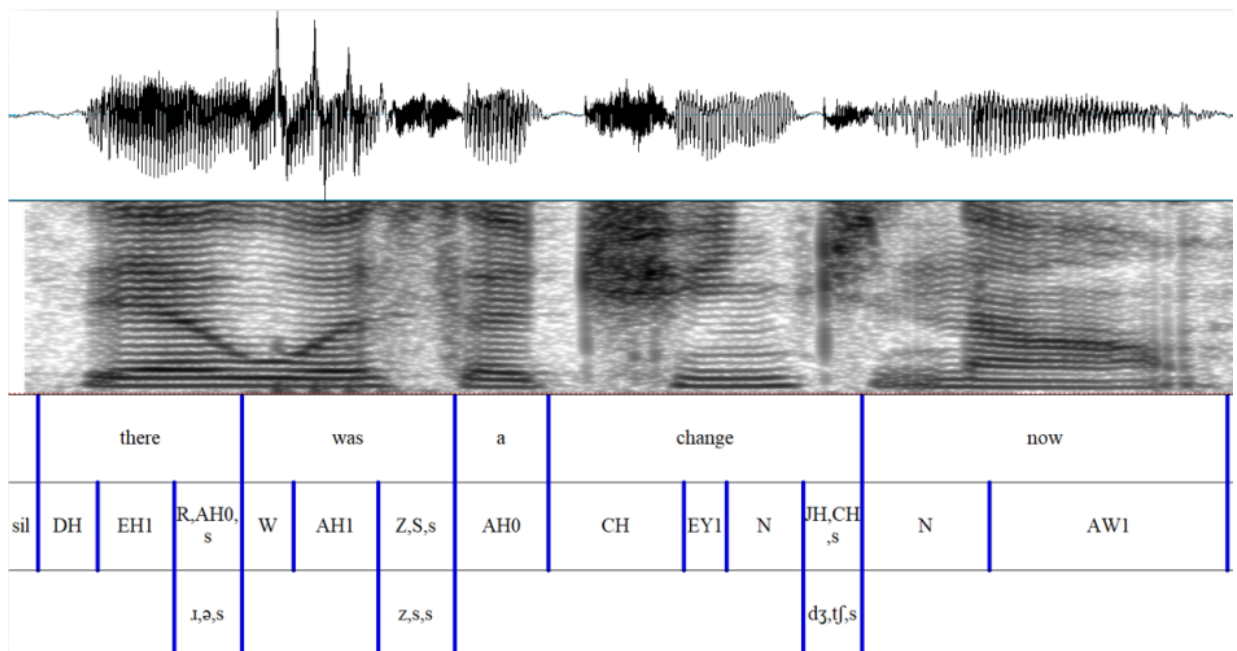


Figure 4.1: Textgrid with manual annotations from a L2-ARCTIC's sample audio file [2]. Top to bottom: speech waveform, spectrogram, words, phonemes, error tags, and comments from the annotator.

L2-ARCTIC [2] targets non-native English speakers. It consists of recordings from a total of 24 individuals whose first languages (L1s) are Hindi, Korean, Mandarin, Spanish, Arabic, and Vietnamese. Each L1 group includes two male and two female speakers, and each speaker contributed approximately one hour of recorded speech, specifically reading from CMU’s ARCTIC prompts. From these recordings, orthographic and forced-aligned phonetic transcriptions were made, and each speaker’s 150 utterances were manually annotated to look for substitutions, deletions, and additions, which are three different types of mispronunciation errors.

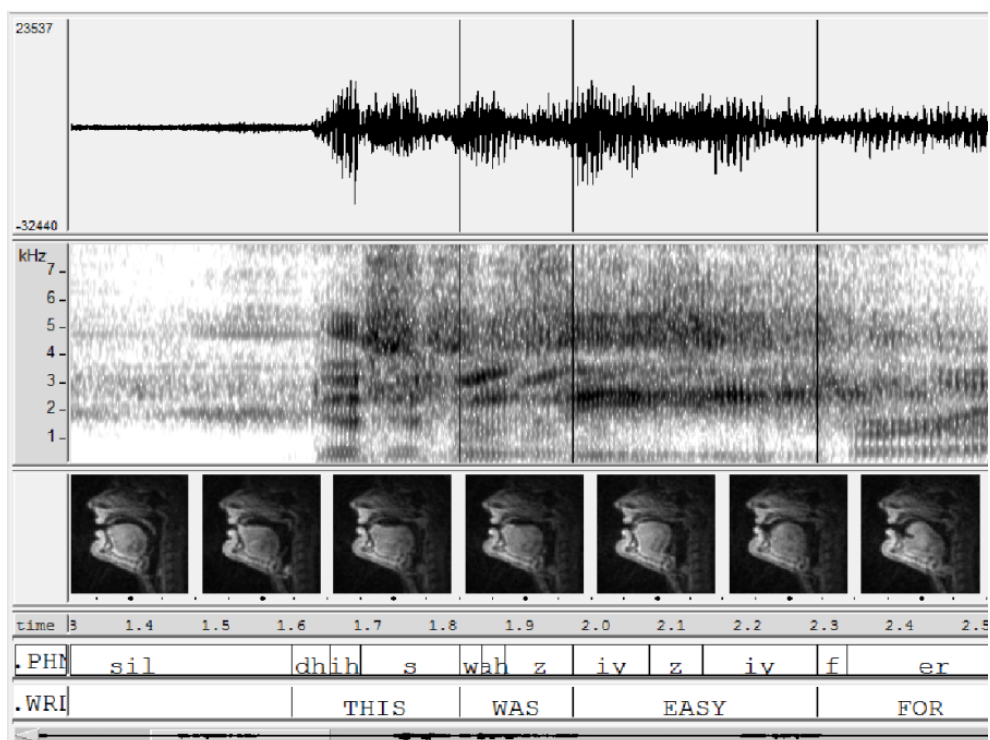


Figure 4.2: Textgrid with manual annotations from a TIMIT’s sample audio file [3]. Top to bottom: speech waveform, spectrogram, phonemes, words

TIMIT [3] is a speech dataset corpus that supplies speech data that can be used for acoustic-phonetic research as well as for the advancement and assessment of automatic speech recognition systems. It is also a collection of broadband recordings of 630 people who speak eight of the most common types of American English. Each person was asked to read ten phonetically dense sentences. The TIMIT corpus consists of orthographic, phonetic, and word transcriptions that are aligned with time, along with a speech waveform file for each utterance, which is encoded in 16-bit format and has a sampling rate of 16 kHz.



The development of the corpus was a collaborative endeavor between the Massachusetts Institute of Technology (MIT), SRI International (SRI), and Texas Instruments, Inc. (TI). The lecture was recorded at TI, transcribed at the MIT, and then checked and prepared for CD-ROM manufacturing by the National Institute of Standards and Technology (NIST).

There are other datasets like TIMIT and L2-ARCTIC that target other languages, such as CU-CHLOE [68] and iCALL [69] for Chinese.

## 4.2 Dataset Construction

It was developed a simple and effective way to encourage participation and collect audio files for the dataset. The desired user (native or non-native) is required to record short audio clips using their own devices and send them to our number as WhatsApp audio files. In this way, by implementing state-of-the-art deep learning methods, the audio file is only required to generate its corresponding spectrogram to create a dataset, thus dealing with data scarcity and a lack of resources. To expand the dataset, the following text was sent to several test subjects in order to obtain different audios:

Hello, I need you to do me a favor, help me recording some audios of some words in Kichwa, it is for a thesis. Send me as Whatsapp audios right here.

Here are the requirements for the recordings:

- Record at a safe distance from the microphone.
- Speak at normal speed and timbre.
- Avoid at all costs blowing into your phone's microphone to avoid generating noise, being in a drafty environment, or being in an environment with background noise.
- One audio per word. If there are 5 words, 5 audios.
- Record the word completely.
- Record the word as you think it is pronounced (depending on whether or not you speak Quechua).

These are the words: yana, yanami, yurak, yurakmi, killu, killumi, puka, pukami, misi, misika, allku, allkuka, wakra, wakra, atallpa, atallpaka, kuchi, kuchika, challwa, challwaka, amaru, amaruka, wiwa, Maykan, Kan.

Thank you for your cooperation and please take the requirements into account.

After that, they immediately send the audio files, and then the dataset collector can easily download the audio files in a single batch by using WhatsApp Web.

### 4.3 Dataset Labelling

For the labeling process, a procedure is followed in order to save the tags corresponding to the audio file. In this project, there are two labels: either well-pronounced or mispronounced, which can also be represented with the integers 1 or 0, respectively. Therefore, it is employed a folder structure in order to extract the labels from each audio file. This method consists of creating a root directory for the Kichwa dataset, creating subdirectories with the names of the labels, and saving the corresponding audio files to each label according to some criteria (see Figure 4.3). This type of labeling enables us to extract the labels and preprocess them later.

The criteria employed for labeling each audio file were partly based on the fluency in Kichwa of the participant in the non-synthetic case. This means that if the participant was a native Kichwa speaker, the label of its audio files is 1. Otherwise, if the participant is a non-native Kichwa speaker, the data collector, in this case the author himself, takes a native Kichwa speaker as a model and compare the spoken word. If the word pronounced by the non-native speaker matches up with the word pronounced by the other counterpart, it was categorized as 1; otherwise, it was categorized as 0. Finally, in the case of synthetic data (see Section 5.2), since all the speakers are pronouncing Kichwa words according their languages, their respective audio files will be labelled as 0.

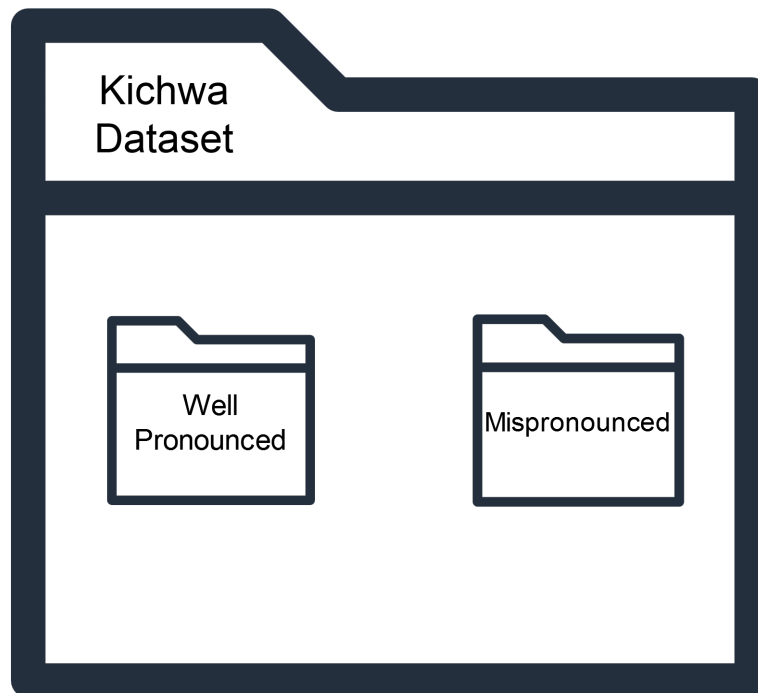


Figure 4.3: Folder directory labelling

## 4.4 Dataset Splitting and Audio Distribution

Given the non-synthetic and synthetic audio files of various speakers, it raises the question of how the data is distributed over the dataset. For this work, it was decided to split and distribute the data in two ways, therefore creating two datasets.

- **Original Dataset.** The synthetic data is considered only for the training set. Since the dataset is imbalanced at the beginning, stratified dataset splitting is used on the original dataset to obtain metrics such as PPV and NPV. In other words, the training, validation, and testing sets have the same proportion of well-pronounced and mispronounced words. After the splitting, all the synthetic audio samples were added to the mispronounced training set. Only non-synthetic voices were added to the validation and testing datasets.
- **Synthetic Dataset.** The synthetic data is considered for training, testing, and validation sets. In this way, a balanced dataset is obtained and, in consequence, split the data into 70%, 20%, and 10% for training, validation, and testing sets, respectively.

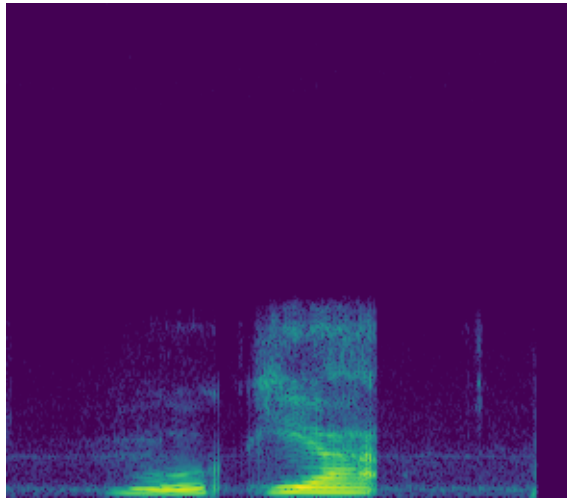


Figure 4.4: Spectrogram representation of the word “allku” pronounced by a non-native female speaker

## 4.5 Preprocessing and Spectrogram Conversion

Since a CNN-based deep learning model is used, it is necessary to decide how the audio dataset is represented in a 2D input way. The traditional spectrogram (see Figure 4.4), like other time-frequency representations, offers a visual representation of the temporal and frequency distribution of energy in a signal, specifically in relation to Fourier’s frequency [70].

After deciding the representation of the data, each audio file is resampled before the spectrogram conversion. In the case of non-synthetic data, audio files are sampled at 48000 Hz, and in the case of non-synthetic data, at 24000. Then both audios are resampled to 44100 Hz. Finally, the duration of all audio files is adjusted to 3 seconds.

## 4.6 Selected CNN Architectures

### 4.6.1 AlexNet

Alexnet [71], [72] is a convolutional neural network (CNN) that achieved success in the 2012 Imagenet Large Scale Visual Recognition Challenge, specifically in the tasks of object recognition and image categorization. The model performed training using ImageNet dataset, which includes over 1.2 million photos belonging to 1000 distinct categories. It includes a total of eight trainable layers. Those layers have been divided into a total of five

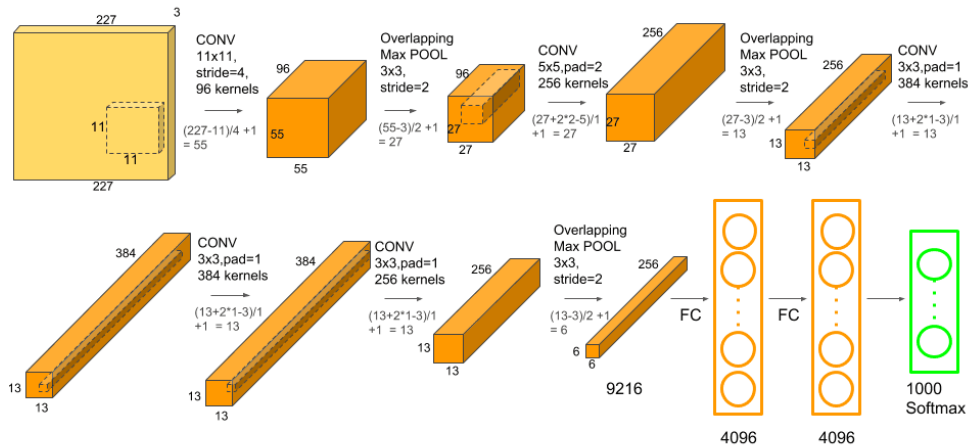


Figure 4.5: AlexNet Architecture Diagram [4]

convolutional layers and three fully connected layers. The Rectified Linear Unit (ReLU) activation function is employed in all trainable layers except for the last fully connected layer, which adopts a softmax function. The architectural design also has non-trainable layers, including three pooling layers, two normalization layers, and one dropout layer, which fulfill the purpose of preventing overfitting (See Figure 4.5). Researchers choose to employ the Rectified Linear Unit (ReLU) function. It was shown that deep convolutional neural networks implementing rectified linear units (ReLU) exhibited much faster training speeds compared to their counterparts employing hyperbolic tangent (tanh) units. When compared to the tanh activation function, the use of the ReLU activation function solely generated a substantially faster training process, achieving a 25% error rate on the training set six times more quickly. Stochastic Gradient Descent has been used with a learning rate of 0.01, a momentum of 0.9, and a weight decay of 0.0005. The training process takes place using two GPUs, specifically GTX 580, in order to achieve parallelism. The graphics processing units (GPUs) employed in the system have a memory capacity of 3 gigabytes (GB) each. The network is divided into two equal halves distributed across the two graphical processing units (GPUs).

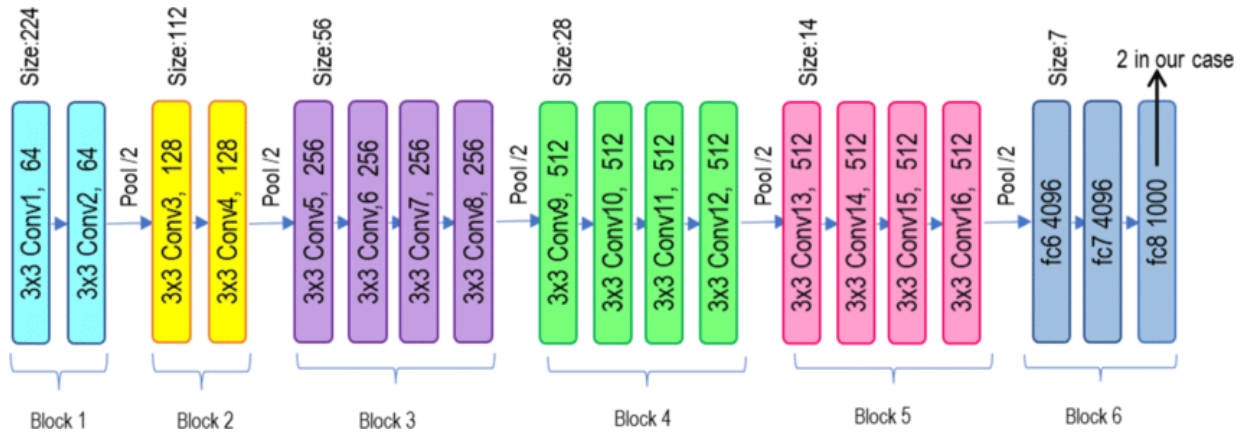


Figure 4.6: VGG19 Architecture Diagram [5]

## 4.6.2 VGG19

At the 2014 ILSVRC, A. Zisserman and K. Simonyan from the University of Oxford’s convolutional neural network model, known as the VGG model or VGGNet, took first and second place in the categories of object detection and image classification, respectively [73]. The proposed approach substitutes the use of big kernel-sized filters with a sequence of several  $3 \times 3$  kernel-sized filters, resulting in a notable improvement compared to the AlexNet architecture [73]. The VGG network is designed using convolutional filters of reduced dimensions. The VGG-16 architecture consists of a series of 13 convolutional layers, three fully connected layers, and a total of 16 layers in all. In contrast, VGG-19 shows an identical architecture to VGG-16, although with the inclusion of three more convolutional layers (See Figure 4.6) [73].

## 4.6.3 ResNet50

ResNets make it possible to train up to hundreds or even thousands of layers and still achieve outstanding performance, which is the case with ResNet-152, while still having lower complexity than VGGNet [74]. In general, training a deep neural network with lots of layers involves the vanishing gradient problem. The backpropagation algorithm is employed to iteratively adjust the weights of a neural network by using the chain rule of derivatives. As the process progresses, the repeated multiplication of gradients can result in the weights becoming exceedingly small as they propagate towards earlier layers [75]. To avoid this problem, ResNets use the concept of skip connections. In a conventional network

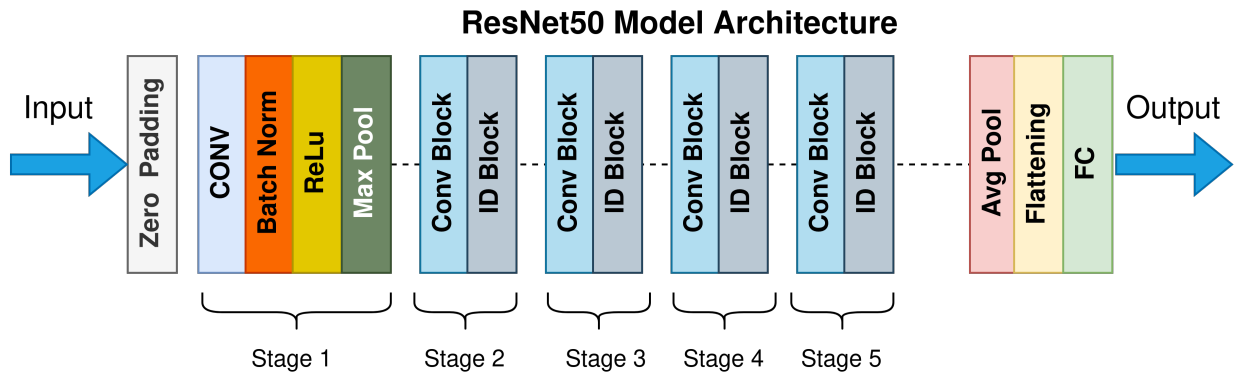


Figure 4.7: ResNet50 Architecture Diagram

architecture, the convolutional layers are sequentially arranged, but in skip connections, the conventional layers are sequentially arranged, with the original input being added into the output of the convolutional block [75].

Finally, ResNet50 can be described as a CNN architecture that consists of 50 layers that are divided into five residual blocks. Each residual block contains two convolutional layers, a batch normalization unit, and a ReLu unit after it. The output of the second convolutional layer is added to the input of the residual block and then passed to the next residual block [76]. The final layer is a fully connected layer that retrieves a flattened output from the last residual block and applies a softmax activation function to produce the final output probabilities.

## 4.7 Data Augmentation

First of all, in order to standardize all the audios in the dataset, the duration of all the audios was adjusted to 3 seconds. Then shifting time, shifting pitch, stretching time, adding Gaussian noise, and masking time were considered as audio data augmentation techniques. The following is a brief explanation of each audio augmentation technique [77]:

- **Noise Injection.** It inserts a random number into the data.
- **Time Shift.** It moves the audio to the left or right at an arbitrary second interval. When The audio is moved to the left with  $x$  seconds, the first  $x$  seconds will be

marked as 0 (which indicates quiet) or rollover. When the audio is moved to the right after  $x$  seconds, the final  $x$  seconds will be marked as 0 or rollover.

- **Pitch Shift.** It changes pitch randomly or within a range of semitones.
- **Time Stretch.** It stretches an audio file by a fixed rate, or in other words, changes the speed of the audio.
- **Time Masking.** It deletes a portion of audio.

For this work, each audio is augmented 100 times, and the probability that each type of audio augmentation is applied to such audio is 0.5.

## 4.8 CNN Features-Based Model

### 4.8.1 Feature Extraction

Abstract characteristics and concepts can be taught to CNN by feeding them with raw image pixels. Features like edges and basic textures are taught to the first convolutional layer(s); complex textures and patterns are examples of the types of characteristics that later convolutional layers can learn; and the final convolutional layers are responsible for learning features such as objects or parts of things [78]. Finally, the fully linked layers are trained to establish associations between feature activations and target classes.

Thus, for the current work, knowledge of low-level features and high-level features acquired from pretrained CNN architectures are used to extract the deep learning features from the fully connected layers inspired by S. Akhtar et al. [67].

After obtaining the audio spectrograms and resampling the images to pass them through the CNNs, pre-trained models from AlexNet, VGG19, and ResNet50 are imported, and the point of extraction for the features are defined. For AlexNet, the features are extracted from the output of the first and second fully connected layers; for VGG19, similarly, from the output of the first and second fully connected layers; and for ResNet50, from the input of the first and only fully connected layer. The dimensions of the array of features extracted (number of features) are 4096 in both layers for AlexNet and VGG19 and 2048 in ResNet50.



On top of that, in order to optimize the process of feature extraction, extraction was performed in batches. Particularly, the batch size used was 1024 for every model specified.

### 4.8.2 Dataset Discretization

The majority of classification tasks in the field of machine learning often involve the process of acquiring the ability to differentiate between different nominal class values. However, it is important to note that these tasks can include features that possess ordinal or continuous attributes in addition to their nominal characteristics. Discretization is a widely employed method for achieving this objective. Discretization refers to the procedure of converting attributes with continuous values into nominal values.

For the purpose of discretizing numerical features, Fayyad and Irani [79] proposed a discretization technique that uses the minimum entropy heuristic. The approach employs the concept of class entropy to determine the optimal cut point for discretization among the candidate divisions. The aforementioned approach can be iteratively implemented on the two intervals resulting from the previous split until certain halting criteria are met. This process generates several intervals for the feature. Since a number of studies have shown that this discretization technique is superior to any other technique overall [80], this technique is used in this work.

### 4.8.3 Statistical-based Feature Selection Methods

Statistical-based feature selection methods originate from different statistical measures. The majority of these methods are filter-based, utilizing statistical measures rather than learning algorithms to evaluate feature relevance. Furthermore, the majority of statistical algorithms tend to analyze features in isolation, and therefore, feature redundancy is neglected during the selection phase [81].

For our purpose, we are using the Chi-square score, proposed in [82] as a feature selection method, that utilizes the test of independence to assess whether the feature is independent of the class label. Since this feature selection method can only be used in discrete or ordinal values [81], dataset discretization must be used before as a preprocessing step.

#### 4.8.4 Classification and Hyperparameter Tuning

For simplicity, KNN is used. In particular, KNN is a non-generalizing learning algorithm, also referred to as an instance-based or lazy learning algorithm. The approach does not prioritize the construction of a comprehensive internal model. Instead, it retains all instances associated with the training data in an n-dimensional space. KNN uses data and classifies new data points based on similarity measures (e.g., the Euclidean distance function) [45]. The classification process involves determining the class label of a point by taking into account the majority vote of its k nearest neighbors. The model exhibits resilience to noisy training data, and its accuracy is dependent upon the quality of the data [27]. A major challenge in K-nearest neighbors (KNN) is determining the optimal number of neighbors to consider. KNN is applicable for both classification and regression tasks [27].

In order to evaluate some possible effects between selecting features and not selecting features, the best parameter for the number of neighbors (k) is searched in the dataset without any preprocessing or feature selection methods. Then, using the same obtained hyperparameter, preprocessing, feature selection and hyperparameter tuning in the best parameter for the number of selected features is performed using validation datasets. Finally, KNN is tested on feature-selected data and non-feature-selected data using the testing dataset and searched hyperparameters. For hyperparameter tuning the HPC CEDIA Cluster was used. More precisely, a "gpu" partition with 64 CPU cores, 128000 MB of memory, and 1 Nvidia A100 SXM4 with 40 GB of memory was reserved for 48 hours.

### 4.9 CNN Transfer learning-based Model

#### 4.9.1 Importing Pre-trained CNN Network

The transfer learning model incorporates pre-trained versions of the AlexNet, VGG19, and ResNet networks. After preprocessing the audio datasets, the models trained on the ImageNet database were loaded, which comprises a large number of images. PyTorch provides pre-trained models within its API.

### 4.9.2 Modifying Final Layers

The final three fully connected layers of both the pre-trained AlexNet and VGG models are designed to classify 1000 different classes. ResNet employs a single, fully connected layer consisting of 1000 classes. The layers in question acquire high-level features for data classification, while others acquire low-level features. Consequently, these layers were optimized to conduct classification on our dataset. In this scenario, it was opted to maintain consistency in the number of fully connected layers across different architectures. However, the output of the last layer were modified to two outputs, enabling binary classification. Following the reset and modification of the final layers in each model, the parameters for a fully connected layer were defined using new data. Cross-entropy was employed as the loss function.

### 4.9.3 Hyperparameter Tuning

In order to perform hyperparameter tuning and training, it is need to define hardware specifications and parameters to be tuned.

The HPC CEDIA Cluster for hyperparameter tuning was used. More precisely, it was reserved a “gpu” partition for 48 hours, 64 CPU cores, 128000 MB of memory, and 1 Nvidia A100 SXM4 with 40 GB of memory. Additionally, the learning rate and training batch size were first set as hyperparameters to be tuned during the hyperparameter process in all models. The learning rate determines the step size at each iteration, while the batch size defines the number of samples processed in each training iteration. The optimizers employed were Stochastic Gradient Descent (SGD) with 0.09 momentum for Alexnet, Adam for VGG19, and ResNet. In the case of Adam, we also set the weight decay parameter to be tuned during the hyperparameter process. In the case of Alexnet and VGG19, the number of output features of the first two fully connected layers were set to be tuned. Table 4.1 summarizes the range of values to be chosen for hyperparameter tuning.

## 4.10 Model Performance Metrics

In order to access the mispronunciation system, distinctive metrics are introduced that evaluate performance based on a particular class or label, such as PPV, NPV, and BPV.

Parameter/Model	AlexNet	VGG19	ResNet50
Learning Rate	Logarithmic Random Value from 1e-4 to 1e-1	Logarithmic Random Value from 1e-4 to 1e-1	Logarithmic Random Value from 1e-4 to 1e-1
Training Batch Size	64, 128, 256, 512	64, 128, 256, 512	64, 128, 256, 512
First Fully Connected Output Features	32, 64, 128, 256, 512, 1024, 2048, 4096	32, 64, 128, 256, 512, 1024, 2048, 4096, 8192, 16384	
Second Fully Connected Output Features	32, 64, 128, 256, 512, 1024, 2048, 4096	32, 64, 128, 256, 512, 1024, 2048, 4096, 8192, 16384	
Weight Decay		Uniform Float Value from 1e4 to 1e-2	Uniform Float Value from 1e4 to 1e-2

Table 4.1: Range of Parameters to be Chosen for Hyperparameter Tuning

#### 4.10.1 PPV (Positive Predictive Value)

It estimates the ratio of accurately predicted positive occurrences (true positives) to the total number of cases that the model identified as positive (See Formula 4.1). Intuitively, PPV is the ability of the classifier or model not to label as well-pronounced an audio that is mispronounced [83].

$$PPV = \frac{TP}{TP + FP} \quad (4.1)$$

where  $TP$  and  $FP$  stands for true positives and false positives.

#### 4.10.2 NPV (Negative Predictive Value)

It estimates the ratio of accurately predicted negative occurrences (true negatives) to the total number of cases that the model identified as negative (See Formula 4.2). Intuitively, NPV is the ability of the classifier or model not to label as mispronounced an audio that is well-pronounced [83].

$$NPV = \frac{TN}{TN + FN} \quad (4.2)$$

where  $TN$  and  $FN$  stands for true negatives and false negatives.

### 4.10.3 Custom Metric: BPV (Balanced Predictive Value)

It is the arithmetic mean of the positive and negative predictive values (See Formula 4.3). This metric aims to assess how well-balanced the model's predictive power is. A high value for BPV indicates that the model does a good job of not confusing well-pronounced words with mispronounced words and vice versa. If not, the model does confuse a lot when trying to distinguish well-pronounced words from mispronounced ones, and vice versa.

$$BPV = \frac{PPV + NPV}{2} \quad (4.3)$$

where  $PPV$  and  $NPV$  stands for positive and negative predictive values, respectively.

# Chapter 5

## Results and Discussion

### 5.1 Dataset Construction

When considering building a mispronunciation detection system taking into account the state-of-the-art used datasets (See Section 4.1) for Kichwa, there were several problems.

- There are no Kichwa language speech corpora for mispronunciation system detection.
- In order to test and evaluate a mispronunciation system's detection, it is necessary to have a dataset like L2-ARCTIC and TIMIT, but for the Kichwa language.
- Since L2-ARCTIC, TIMIT, and other similar datasets use Textgrid and audio files to align the transcriptions and label the phonemes and mispronunciation errors, it would be ideal to have a similar format for the Kichwa language dataset, but this is clearly outrageously difficult due to the lack of personnel.
- The process of text-gridding audio files in the Kichwa language requires extensive resources, expertise, and lots of time.
- It is extremely hard to find participants who are fluent in Kichwa and willing to contribute to the dataset. This case also applies even when finding participants who are not fluent in Kichwa, which is completely disheartening.
- The process of collecting data can discourage Kichwa native speakers and non-native speakers from participating in the process due to the duration of the process, lack of incentives, potential privacy concerns, transportation issues, lack of time, etc.

For this work, to overcome these problems, the previously described method in Section was employed. Then a Kichwa dataset that consists of a total of 500 recorded words from various Kichwa speakers and non-speakers was constructed. Each participant in the collection recorded and sent 25 audio files corresponding to the words in Kichwa: yana, yanami, yurak, yurakmi, killu, killumi, puka, pukami, misi, misika, allku, allkuka, wakra, wakraka, atallpa, atallpaka, kuchi, kuchika, challwa, challwaka, amaru, amaruka, wiwa, Maykan, Kan. We summarize the dataset in Tables 5.1 and 5.2.

Number of Participants			
Fluency	Adult Male	Adult Female	Total
Native	3	6	9
Non-Native	6	5	11
Total	9	11	20

Table 5.1: Number of participants in the dataset

Words	
Label	Total
Well-pronounced	100
Mispronounced	400
Total	500

Table 5.2: Distribution of words in Kichwa dataset

## 5.2 Syntetic Data

As in Table 5.2, there is an imbalanced proportion of well-pronounced words. and mispronounced words. When developing transfer learning models, or, to be more precise, hyperparameter tuning, this imbalance can pose a challenge to accurately predict the mispronounced words. In fact, the value of the negative predictive value was always 0, and the positive predictive value was greater than it. This means that the model does not have enough data to learn how to distinguish mispronounced words from pronounced words.

As a solution, it was decided to use synthetic data in the dataset. To generate such data, a website called "Coqui Studio" was used. Initially, the website allows users to input text and have it converted into audio, providing a large amount of data for the model to learn from. Thus, by inputting Kichwa words into the generator, a person's voice and

<b>Name</b>	<b>Gender</b>	<b>Language</b>
Aaron Dreschner	Male	English
Szofi Granger	Female	English
Ige Behringer	Male	English
Claribel Dervla	Female	German
Ilkin Urbano	Male	French
Camilla Holmstrom	Female	Portuguese
Baldur Sanjin	Male	Portuguese
Tanja Adelina	Female	French
Wulf Carlevaro	Male	German
Annmarie Nele	Female	German
Marcos Rudaski	Male	Portuguese
Vjollca Johnnie	Female	English

Table 5.3: Syntetic Participants Details

language can be selected to generate mispronounced Kichwa audio outputs. The synthetic participants’ names and their respective languages used in this work can be detailed in Table 5.3.

In total, 300 Kichwa audios were collected from each artificial participant. It is important to clarify that the website’s model used in this work was XTTS, which is a newer and more expensive model. Therefore, it was expected that the voices generated by XTTS would be of higher quality and more natural-sounding compared to other models, which, in fact, they were.

### 5.3 Dataset Labelling

In fact, it should be ideal to have a voting mechanism in which if two of some experts in Kichwa language agree on the same label, then that label is assigned to that data as in [67], but that is not clearly possible given the circumstances.

### 5.4 Model’s Testing

The best models’ testing metrics can be shown in Table 5.4. Therefore, it might be found appropriate to expose the results and discussion for this section by exposing how the aforementioned table was constructed in addition to relevant observations in their respective stages of construction.



### 5.4.1 CNN Features-based Model

- From Figure A.1 and A.2, it can be observed the hyperparameter best selection given the number of neighbours  $k$ . The highest peak in the figure indicates the best  $k$  selection. It is important to note that the original dataset was transformed in another dataset given the feature extraction location in AlexNet, ResNet50 and VGG19. Moreover, the hyperparameter tuning was done given the training set and validation set.
- From Tables A.1 and A.2, it can be detailed better the model, its feature extraction point, the best  $k$  hyperparameter selection and its corresponding validation balanced predictive value metric given the dataset. Thus, from A.1 and A.2, it is seen that "ResNet50: AvgPool" and "AlexNet: FC1" are the best models without feature selection taking into account the validation BPV for the original and synthetic dataset, respectively. Hereby, those models are also tested using the testing dataset and measured using PPV, NPV and BPV, and added to Table 5.4.
- After establishing the baseline models without feature selection, it is fixed  $k$  for each model and then perform dataset discretization and use Chi-square statistical-based feature selection method to obtain the best features. In Figures A.3 and A.4 we plot the possible improvement of performance based on BPV metric against the number of best important features taken into account. In Chi-square feature selection, the greater the Chi-square score is, the better the feature is [81]. The red line is the best BPV score without feature selection. Thus, the purpose of feature selection here, similarly, is to search for the optimal number of features to be taken into account. It can be observed that in some models, the feature selection method presuppose a greater performance in the validation dataset, while in others, the performance is not evident.
- In Tables A.3 and A.4, the best selection for the number of features according to the validation BPV is extracted for each model, and then obtain the PPV, NPV and BPV metrics using the testing dataset. Then, the best model given its validation BPV is added to Table 5.4.

- When comparing Tables A.1, A.2, A.3 and A.4, it is noted that: in the original dataset, feature selection improves the NPV value for the first three models, on the other hand, in the synthetic dataset, the NPV performance is worsened in some models;

#### 5.4.2 CNN Transfer Learning-based Model

- For the hyperparameter tuning, it is selected a metric in order to optimize it, in this case, validation loss. However, Ray Tune library allows to eventually save the best models per trial given validation loss and validation BPV. This permits to compare the performance between different metric targets when searching a metric for optimization.
- In Tables B.1 to B.12, it is shown the best models given validation loss or balanced predictive value. Moreover, we record some other valuable metrics, for instance, PPV, NPV and BPV for training and validation sets. In the columns named “config”, it can be seen the hyperparameters used in that trial.
- In Tables B.13 to B.18, it is extracted the TOP-5 best models by its corresponding best validation metric (loss or BPV), and tested them by using the testing set. In these tables, it is important to note that there are NPV perfect values, i.e. with a value of 1, which can be a symptom of an inappropriate size of the dataset. By the size of the dataset or the quality of the dataset, they may be deceptively values. Nevertheless, finally, we extract the best models and add them to the general table.

#### 5.4.3 Observations from Table 5.4

- In general, transfer learning-based models achieve better performance in BPV metric.
- Feature selection technique did not show any improvement performance against models without feature selection. On the contrary, the testing metric shows that the best models with feature selection are worse than the ones without it.
- It is not necessary to use deeper convolutional neural networks such as VGG19 or Resnet50 to achieve the best results. For instance, AlexNet achieves the best results

in CNN Transfer Learning-based Model for both original and synthetic datasets.

- In terms of running time, in CNN Transfer Learning-based model, VGG19 took roughly 24 hours to complete hyperparameter search using original dataset and synthetic datasets. Meanwhile, the other two architectures took, in average, 4 hours to complete hyperparameter search for both datasets. On the other hand, in CNN Features-based model, running time was negligible since we performed feature extraction in batches, and the KNN Classifier used 64 cores.

Best Models' Performance on Testing Set			
Original Dataset			
CNN Features-Based Model			
Model	Test PPV	Test NPV	Test BPV
<b>ResNet50: AvgPool Without Feature Selection</b>	<b>0.84</b>	<b>0.50</b>	<b>0.67</b>
VGG19: FC2 with Feature Selection	0.85	0.37	0.61
CNN Transfer Learning Based Model			
Model	Test PPV	Test NPV	Test BPV
AlexNet Loss	0.87	0.61	0.74
<b>AlexNet BPV</b>	<b>0.80</b>	<b>1.00</b>	<b>0.90</b>
VGG19 Loss	0.83	0.64	0.73
VGG19 BPV	0.81	0.40	0.60
ResNet50 Loss	0.80	0.00	0.40
ResNet50 BPV	0.80	0.00	0.40
Synthetic Dataset			
CNN Features-Based Model			
Model	Test PPV	Test NPV	Test BPV
<b>AlexNet: FC1 Without Feature Selection</b>	<b>0.88</b>	<b>0.90</b>	<b>0.89</b>
AlexNet: FC1 with Feature Selection	0.81	0.86	0.84
CNN Transfer Learning-Based Model			
Model	Test PPV	Test NPV	Test BPV
<b>AlexNet Loss</b>	<b>0.83</b>	<b>1.00</b>	<b>0.92</b>
<b>AlexNet BPV</b>	<b>0.83</b>	<b>1.00</b>	<b>0.92</b>
VGG19 Loss	0.83	0.94	0.88
VGG19 BPV	0.79	0.91	0.85
ResNet50 Loss	0.52	1.00	0.76
ResNet50 BPV	0.65	1.00	0.82

Table 5.4: Best Models' Performance on Testing Datasets. Each model has been selected according to validation metric.

# Chapter 6

## Conclusions

This work represents an initial attempt at developing a mispronunciation system for the Kichwa language. In this study, pretrained CNN architectures were employed to classify spectrograms of both well-pronounced and mispronounced words. The first model, the CNN features-based model, extracts the features from the fully connected layers, applies a feature selection method technique in order to extract the discriminative features from the non-discriminative ones, and classifies them with a KNN classifier. The second model, CNN transfer-learning-based, uses knowledge from convolutional layers and modifies the classifier layer for binary classification (well-pronounced or mispronounced). Overall, we can conclude that the CNN transfer learning-based method is better than the CNN feature-based method, but nevertheless, the CNN feature-based method achieves similar performance to the CNN transfer learning-based method in the synthetic dataset, With a larger dataset, the transfer learning-based method may outperform the CNN feature-based method in the original and synthetic datasets.

### 6.1 Observations and Limitations

- Drawing significant conclusions from a small dataset is nearly impossible, especially from a small and imbalanced dataset. However, a larger and more balanced dataset would allow for more meaningful conclusions, and it would facilitate a comprehensive study of feature selection techniques and deep learning architectures with increased reliability.

- Since the majority of Kichwa words proposed were actually similar in pronunciation to Spanish words, most of the audios went to the majority class of well-pronounced words, leaving the minority class of poorly pronounced words underrepresented in the recordings. Thus, it is probable that the selection of words may be unfavorable for our work. In other words, it should have been chosen words with much more difficult pronunciation.
- The lack of participants and an official project for Kichwa preservation makes it difficult to gather a diverse range of audio samples. Additionally, the absence of standardized pronunciation guidelines further complicates the accuracy and consistency of the recordings.
- The use of synthetic data may ultimately distort the results and limit the reliability of any conclusions drawn from the study, particularly in the use of the synthetic dataset.
- The way words were collected may not be appropriate due to their individual recordings. In fact, it could be better if a person recorded different sentences from different sources. Hence, it should have been better to follow recording and organizational guidelines from official mispronunciation datasets such as TIMIT or L2-Arctic.
- It might be a risky decision to construct a mispronunciation model based on Kichwa first rather than a model based on official mispronunciation datasets such as TIMIT or L2-Arctic. In this way, it could have drawn better conclusions from the model, given the official datasets and learning how such datasets are constructed, in order to implement a similar one for Kichwa.
- CNN architectures may not be suitable for phrases or sentences. Instead, recurrent neural networks or transformer models are often employed due to their ability to capture sequential dependencies.
- The work only uses a KNN classifier and a feature selection method technique to evaluate performance, thus probably ignoring better alternatives.

- It is necessary to implement and train the models in a high-performance computing cluster.

# Bibliography

- [1] I. H. Sarker, “Deep learning: a comprehensive overview on techniques, taxonomy, applications and research directions,” *SN Computer Science*, vol. 2, no. 6, p. 420, 2021.
- [2] G. Zhao, S. Somsat, A. Silpachai, I. Lucic, E. Chukharev-Hudilainen, J. Levis, and R. Gutierrez-Osuna, “L2-arctic: A non-native english speech corpus,” in *Proc. Interspeech*, 2018, p. 2783–2787. [Online]. Available: <http://dx.doi.org/10.21437/Interspeech.2018-1110>
- [3] J. S. Garofolo, “Timit acoustic phonetic continuous speech corpus,” *Linguistic Data Consortium, 1993*, 1993.
- [4] A. Pujara, “Concept of alexnet:- convolutional neural network,” Jun 2021. [Online]. Available: <https://medium.com/analytics-vidhya/concept-of-alexnet-convolutional-neural-network-6e73b4f9ee30>
- [5] A. Khattar and S. Quadri, “Generalization of convolutional network to domain adaptation network for classification of disaster images on twitter,” *Multimedia Tools and Applications*, vol. 81, no. 21, pp. 30 437–30 464, 2022.
- [6] UNESCO, “Indigenous peoples’ right to education: overview of the measures supporting the right to education for indigenous peoples reported by member states in the context of the ninth consultation on the 1960 convention and recommendation against discrimination in education,” 2019.
- [7] “Indigenous peoples - UNESCO,” <https://www.unesco.org/en/indigenous-peoples>, (accessed Aug 28, 2023).



- [8] “Indigenous peoples overview,” <https://www.worldbank.org/en/topic/indigenouspeoples#1>, (accessed Aug 28, 2023).
- [9] UNESCO, “The UNESCO Courier,” in *2019: the International Year of Indigenous Languages*, 2019.
- [10] M. de Educación, “Informe en relación a la “convención relativa a la lucha contra las discriminaciones en la esfera de la enseñanza (la convención) y la recomendación relativa a la lucha contra las discriminaciones en la esfera de la enseñanza (la recomendación), aprobadas por la conferencia general de la unesco en 1960.”
- [11] J. G. Armijos Monar, B. N. Fuertes Lopez, J. E. Delgado Altamirano, and V. M. V. Villa, “University indigenous students’ perceptions towards kichwa, spanish and english.” *English Language Teaching*, vol. 11, no. 2, pp. 131–148, 2018.
- [12] M. A. Peabody, “Methods for pronunciation assessment in computer aided language learning,” Ph.D. dissertation, Massachusetts Institute of Technology, 2011.
- [13] T. R. Kurfess *et al.*, *Robotics and automation handbook*. CRC press Boca Raton, FL, 2005, vol. 414.
- [14] M. Rosheim, *Leonardo’s Lost Robots*. Springer Science & Business Media, 2006.
- [15] P. Hamet and J. Tremblay, “Artificial intelligence in medicine,” *Metabolism*, vol. 69, pp. S36–S40, 2017.
- [16] A. Holzinger, G. Langs, H. Denk, K. Zatloukal, and H. Müller, “Causability and explainability of artificial intelligence in medicine,” *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 9, no. 4, p. e1312, 2019.
- [17] B. M. Lake, T. D. Ullman, J. B. Tenenbaum, and S. J. Gershman, “Building machines that learn and think like people,” *Behavioral and brain sciences*, vol. 40, p. e253, 2017.
- [18] Z. Zhang, P. Cui, and W. Zhu, “Deep learning on graphs: A survey,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 1, pp. 249–270, 2020.

- [19] C. Zhang and Y. Lu, “Study on artificial intelligence: The state of the art and future prospects,” *Journal of Industrial Information Integration*, vol. 23, p. 100224, 2021.
- [20] M. Bertolini, D. Mezzogori, M. Neroni, and F. Zammori, “Machine learning for industrial applications: A comprehensive literature review,” *Expert Systems with Applications*, vol. 175, p. 114820, 2021.
- [21] C. Robert, “Machine learning, a probabilistic perspective,” 2014.
- [22] I. H. Sarker, A. Kayes, S. Badsha, H. Alqahtani, P. Watters, and A. Ng, “Cybersecurity data science: an overview from machine learning perspective,” *Journal of Big data*, vol. 7, pp. 1–29, 2020.
- [23] S. Das, A. Dey, A. Pal, and N. Roy, “Applications of artificial intelligence in machine learning: review and prospect,” *International Journal of Computer Applications*, vol. 115, no. 9, 2015.
- [24] G. Carleo, I. Cirac, K. Cranmer, L. Daudet, M. Schuld, N. Tishby, L. Vogt-Maranto, and L. Zdeborová, “Machine learning and the physical sciences,” *Reviews of Modern Physics*, vol. 91, no. 4, p. 045002, 2019.
- [25] R. S. Michalski, J. G. Carbonell, and T. M. Mitchell, *Machine learning: An artificial intelligence approach*. Springer Science & Business Media, 2013.
- [26] S. Ray, “A quick review of machine learning algorithms,” in *2019 International conference on machine learning, big data, cloud and parallel computing (COMITCon)*. IEEE, 2019, pp. 35–39.
- [27] I. H. Sarker, “Machine learning: Algorithms, real-world applications and research directions,” *SN computer science*, vol. 2, no. 3, p. 160, 2021.
- [28] D. Dhall, R. Kaur, and M. Juneja, “Machine learning: a review of the algorithms and its applications,” *Proceedings of ICRIC 2019: Recent Innovations in Computing*, pp. 47–63, 2020.

- [29] L. Deng, “A tutorial survey of architectures, algorithms, and applications for deep learning,” *APSIPA transactions on Signal and Information Processing*, vol. 3, p. e2, 2014.
- [30] A. Arya and S. Sridhar, “Overview of big data analytics technologies in smart grid,” 2023.
- [31] J. Han, J. Pei, and H. Tong, *Data mining: concepts and techniques*. Morgan kaufmann, 2022.
- [32] G. Bonaccorso, *Machine learning algorithms*. Packt Publishing Ltd, 2017.
- [33] A. Arshaghi, M. Ashourian, and L. Ghabeli, “Potato diseases detection and classification using deep learning methods,” *Multimedia Tools and Applications*, vol. 82, no. 4, pp. 5725–5742, 2023.
- [34] F. Anowar, S. Sadaoui, and B. Selim, “Conceptual and empirical comparison of dimensionality reduction algorithms (pca, kpca, lda, mds, svd, lle, isomap, le, ica, t-sne),” *Computer Science Review*, vol. 40, p. 100378, 2021.
- [35] I. H. Sarker, Y. B. Abushark, F. Alsolami, and A. I. Khan, “Intrudtree: a machine learning based cyber security intrusion detection model,” *Symmetry*, vol. 12, no. 5, p. 754, 2020.
- [36] I. H. Sarker, Y. B. Abushark, and A. I. Khan, “Contextpca: Predicting context-aware smartphone apps usage based on machine learning techniques,” *Symmetry*, vol. 12, no. 4, p. 499, 2020.
- [37] I. H. Sarker, “Data science and analytics: an overview from data-driven smart computing, decision-making and applications perspective,” *SN Computer Science*, vol. 2, no. 5, p. 377, 2021.
- [38] Y. Guo, Y. Liu, A. Oerlemans, S. Lao, S. Wu, and M. S. Lew, “Deep learning for visual understanding: A review,” *Neurocomputing*, vol. 187, pp. 27–48, 2016.
- [39] A. Ng, “Machine learning yearning: Technical strategy for ai engineers in the era of deep learning. retrieved online at h ttps,” 2019.

- [40] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [41] Y. Xin, L. Kong, Z. Liu, Y. Chen, Y. Li, H. Zhu, M. Gao, H. Hou, and C. Wang, “Machine learning and deep learning methods for cybersecurity,” *Ieee access*, vol. 6, pp. 35 365–35 381, 2018.
- [42] L. Deng, D. Yu *et al.*, “Deep learning: methods and applications,” *Foundations and trends® in signal processing*, vol. 7, no. 3–4, pp. 197–387, 2014.
- [43] I. H. Sarker, A. Kayes, and P. Watters, “Effectiveness analysis of machine learning classification models for predicting personalized context-aware smartphone usage,” *Journal of Big Data*, vol. 6, no. 1, pp. 1–28, 2019.
- [44] L. Deng, “Three classes of deep learning architectures and their applications: a tutorial survey,” *APSIPA transactions on signal and information processing*, vol. 57, p. 58, 2012.
- [45] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg *et al.*, “Scikit-learn: Machine learning in python,” *the Journal of machine Learning research*, vol. 12, pp. 2825–2830, 2011.
- [46] I. H. Sarker, M. H. Furhad, and R. Nowrozy, “Ai-driven cybersecurity: an overview, security intelligence modeling and research directions,” *SN Computer Science*, vol. 2, pp. 1–18, 2021.
- [47] I. H. Sarker, “Deep cybersecurity: a comprehensive overview from neural network and deep learning perspective,” *SN Computer Science*, vol. 2, no. 3, p. 154, 2021.
- [48] A. Géron, *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow.* ” O’Reilly Media, Inc.”, 2022.
- [49] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

- [50] K. He, X. Zhang, S. Ren, and J. Sun, “Spatial pyramid pooling in deep convolutional networks for visual recognition,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 9, pp. 1904–1916, 2015.
- [51] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Advances in neural information processing systems*, vol. 25, 2012.
- [52] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [53] K. Weiss, T. M. Khoshgoftaar, and D. Wang, “A survey of transfer learning,” *Journal of Big data*, vol. 3, no. 1, pp. 1–40, 2016.
- [54] F. Nazir, M. N. Majeed, M. A. Ghazanfar, and M. Maqsood, “Mispronunciation detection using deep convolutional neural network features and transfer learning-based model for arabic phonemes,” *IEEE Access*, vol. 7, pp. 52 589–52 608, 2019.
- [55] S. M. Witt and S. J. Young, “Phone-level pronunciation scoring and assessment for interactive language learning,” *Speech communication*, vol. 30, no. 2-3, pp. 95–108, 2000.
- [56] H. Franco, L. Neumeyer, M. Ramos, and H. Bratt, “Automatic detection of phone-level mispronunciation for language learning,” in *Sixth European Conference on Speech Communication and Technology*, 1999.
- [57] H. Strik, K. Truong, F. De Wet, and C. Cucchiarini, “Comparing different approaches for automatic pronunciation error detection,” *Speech communication*, vol. 51, no. 10, pp. 845–852, 2009.
- [58] N. Minematsu, S. Asakawa, and K. Hirose, “Structural representation of the pronunciation and its use for call,” in *2006 IEEE Spoken Language Technology Workshop*. IEEE, 2006, pp. 126–129.

- [59] Y.-B. Wang and L.-S. Lee, “Improved approaches of modeling and detecting error patterns with empirical analysis for computer-aided pronunciation training,” in *2012 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2012, pp. 5049–5052.
- [60] F. Zhang, C. Huang, F. K. Soong, M. Chu, and R. Wang, “Automatic mispronunciation detection for mandarin,” in *2008 IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2008, pp. 5077–5080.
- [61] A. Ito, Y.-L. Lim, M. Suzuki, and S. Makino, “Pronunciation error detection method based on error rule clustering using a decision tree,” in *Ninth European Conference on Speech Communication and Technology*, 2005.
- [62] G. Georgoulas, V. C. Georgopoulos, and C. D. Stylios, “Speech sound classification and detection of articulation disorders with support vector machines and wavelets,” in *2006 International Conference of the IEEE Engineering in Medicine and Biology Society*. IEEE, 2006, pp. 2199–2202.
- [63] I. Amdal, M. H. Johnsen, and E. Versvik, “Automatic evaluation of quantity contrast in non-native norwegian speech,” in *International Workshop on Speech and Language Technology in Education*, 2009.
- [64] H. Li, J. Liang, S. Wang, and B. Xu, “An efficient mispronunciation detection method using glds-svm and formant enhanced features,” in *2009 IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2009, pp. 4845–4848.
- [65] S. Wei, G. Hu, Y. Hu, and R.-H. Wang, “A new method for mispronunciation detection using support vector machine based on pronunciation space models,” *Speech Communication*, vol. 51, no. 10, pp. 896–905, 2009.
- [66] M. Maqsood, “Feature selection for arabic mispronunciation detection based on sequential floating forward selection and data mining classifiers,” *Pakistan Journal of Science*, vol. 68, no. 4, 2016.
- [67] S. Akhtar, F. Hussain, F. R. Raja, M. Ehatisham-ul haq, N. K. Baloch, F. Ishmanov, and Y. B. Zikria, “Improving mispronunciation detection of arabic words for non-

- native learners using deep convolutional neural network features,” *Electronics*, vol. 9, no. 6, p. 963, 2020.
- [68] M. Wu, K. Li, W.-K. Leung, and H. Meng, “Transformer based end-to-end mispronunciation detection and diagnosis.” in *Interspeech*, 2021, pp. 3954–3958.
- [69] N. F. Chen, R. Tong, D. Wee, P. Lee, B. Ma, and H. Li, “icall corpus: Mandarin chinese spoken by non-native speakers of european descent,” in *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.
- [70] S. A. Fulop and K. Fitz, “A spectrogram for the twenty-first century,” *Acoustics today*, vol. 2, no. 3, pp. 26–33, 2006.
- [71] B. Akera, “Alexnet: A brief review,” May 2020. [Online]. Available: <https://medium.com/ai-research-lab-kampala/alexnet-a-brief-review-14979ce7cc84>
- [72] M. Z. Alom, T. M. Taha, C. Yakopcic, S. Westberg, P. Sidike, M. S. Nasrin, B. C. Van Esesn, A. A. S. Awwal, and V. K. Asari, “The history began from alexnet: A comprehensive survey on deep learning approaches,” *arXiv preprint arXiv:1803.01164*, 2018.
- [73] G. Boesch, “Vgg very deep convolutional networks (vggnet) - what you need to know - viso.ai,” Oct 2021. [Online]. Available: <https://viso.ai/deep-learning/vgg-very-deep-convolutional-networks/>
- [74] Jan 2019. [Online]. Available: <https://iq.opengenius.org/resnet/>
- [75] S. R. Sapiroddy, “Resnet-50: Introduction,” Jul 2023. [Online]. Available: <https://srsapiroddy.medium.com/resnet-50-introduction-b5435fdb66f>
- [76] wisdomml, “Understanding resnet-50 in depth: Architecture, skip connections, and advantages over other networks - wisdom ml,” Mar 2023. [Online]. Available: <https://wisdomml.in/understanding-resnet-50-in-depth-architecture-skip-connections-and-advantages-over-other-networks/>
- [77] E. Ma, “Data augmentation for audio,” Jun 2019. [Online]. Available: <https://medium.com/@makcedward/data-augmentation-for-audio-76912b01fdf6>

- [78] C. Molnar, “Interpretable machine learning,” Aug 2023. [Online]. Available: <https://christophm.github.io/interpretable-ml-book/cnn-features.html>
- [79] U. Fayyad and K. Irani, “Multi-interval discretisation of continuous-valued attributes,[in:] proceedings of the xiii international joint conference on artificial intelligence,” 1993.
- [80] M. A. Hall, “Correlation-based feature selection for machine learning,” Ph.D. dissertation, The University of Waikato, 1999.
- [81] J. Li, K. Cheng, S. Wang, F. Morstatter, R. P. Trevino, J. Tang, and H. Liu, “Feature selection: A data perspective,” *ACM computing surveys (CSUR)*, vol. 50, no. 6, pp. 1–45, 2017.
- [82] H. Liu and R. Setiono, “Chi2: Feature selection and discretization of numeric attributes,” in *Proceedings of 7th IEEE international conference on tools with artificial intelligence*. Ieee, 1995, pp. 388–391.
- [83] [Online]. Available: [https://scikit-learn.org/stable/modules/generated/sklearn.metrics.precision\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.precision_score.html)



# Appendices

# Appendix A

## CNN Features-based Model

### Hyperparameter Tuning and Testing Results: Diagrams and Tables

Hyperparameter Tuning (k) and Testing without Feature Selection in Original Dataset					
Best Results					
Model	k	Validation BPV	Test PPV	Test NPV	Test BPV
AlexNet: FC1	9	0.62	0.82	0.35	0.59
AlexNet: FC2	4	0.58	0.82	0.26	0.54
VGG19: FC1	9	0.60	0.84	0.38	0.61
VGG19: FC2	5	0.63	0.86	0.45	0.66
ResNet50: AvgPool	5	0.64	0.84	0.50	0.67

Table A.1: CNN Features-based Hyperparameter Tuning on k and Testing in Original Dataset

Hyperparameter Tuning (k) and Testing without Feature Selection in Synthetic Dataset					
Best Results					
Model	k	Validation BPV	Test PPV	Test NPV	Test BPV
AlexNet: FC1	9	0.90	0.88	0.90	0.89
AlexNet: FC2	9	0.88	0.83	0.94	0.88
VGG19: FC1	10	0.89	0.82	0.89	0.85
VGG19: FC2	1	0.90	0.82	0.91	0.87
ResNet50: AvgPool	9	0.89	0.84	0.92	0.88

Table A.2: CNN Features-based Hyperparameter Tuning on k and Testing in Synthetic Dataset

Hyperparameter Tuning (number of features) and Testing with Feature Selection in Original Dataset					
Best Results					
Model	Number of Features	Validation BPV	Test PPV	Test NPV	Test BPV
AlexNet: FC1	245	0.69	0.82	0.43	0.62
AlexNet: FC2	88	0.62	0.83	0.29	0.56
VGG19: FC1	188	0.60	0.84	0.45	0.65
VGG19: FC2	47	0.72	0.85	0.37	0.61
ResNet50: AvgPool	346	0.65	0.82	0.32	0.57

Table A.3: CNN Features-based Hyperparameter Tuning with Feature Selection and Testing in Original Dataset

Hyperparameter Tuning (number of features) and Testing with Feature Selection in Synthetic Dataset					
Best Results					
Model	Number of Features	Validation BPV	Test PPV	Test NPV	Test BPV
AlexNet: FC1	360	0.92	0.81	0.86	0.84
AlexNet: FC2	249	0.91	0.84	0.94	0.89
VGG19: FC1	249	0.91	0.80	0.82	0.81
VGG19: FC2	166	0.90	0.80	0.82	0.81
ResNet50: AvgPool	386	0.90	0.78	0.94	0.86

Table A.4: CNN Features-based Hyperparameter Tuning with Feature Selection and Testing in Synthetic Dataset

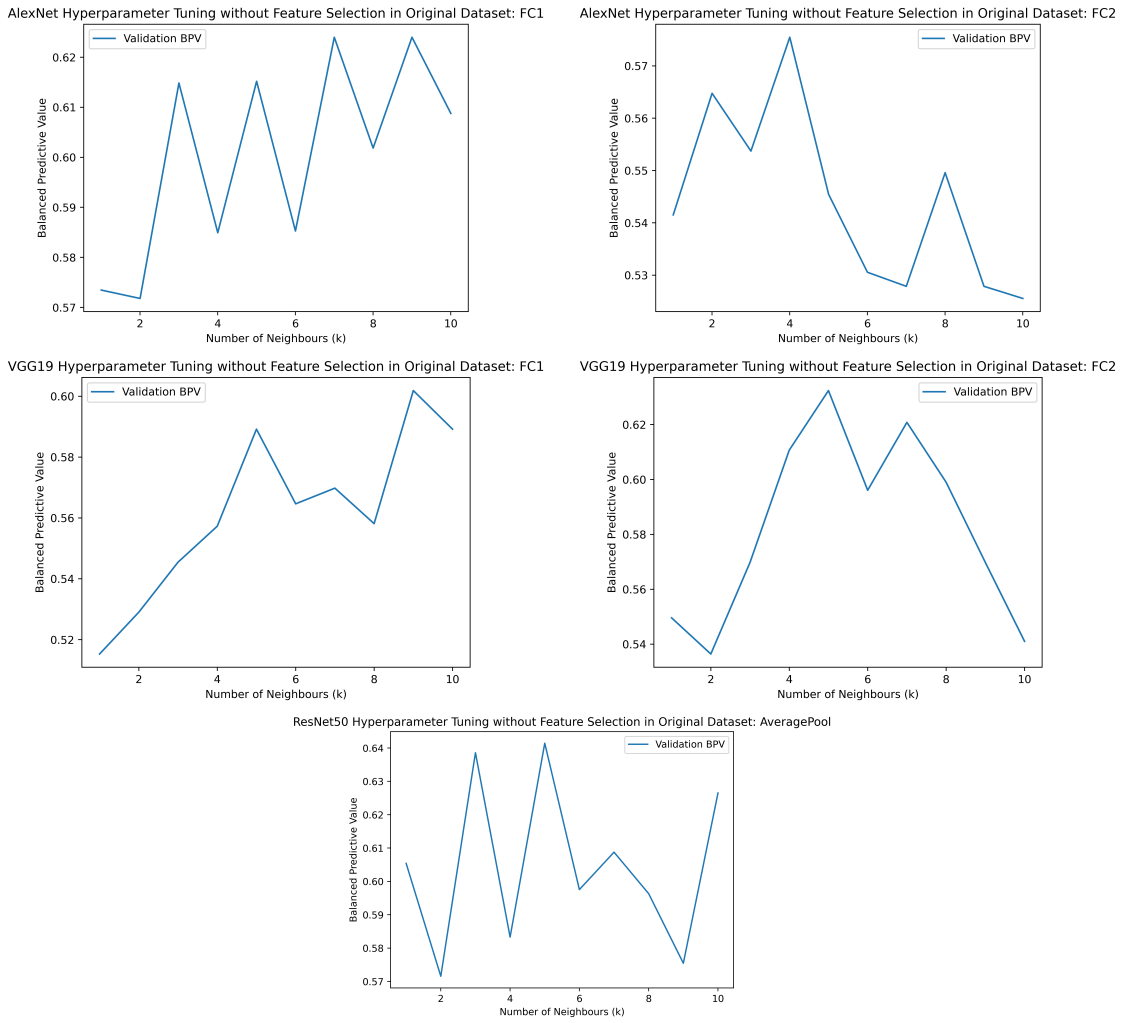
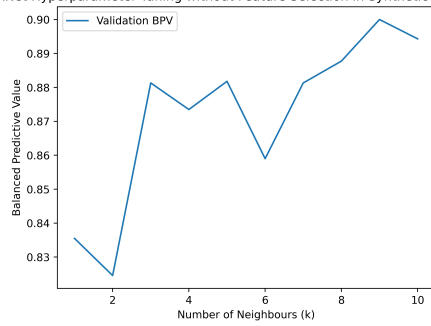
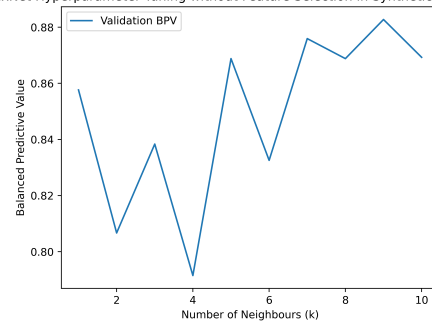


Figure A.1: KNN Classifier Hyperparameter Tuning in Original Dataset.

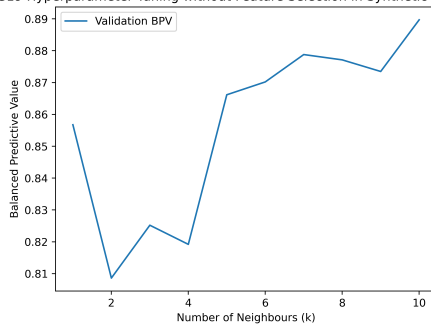
AlexNet Hyperparameter Tuning without Feature Selection in Synthetic Dataset: FC1



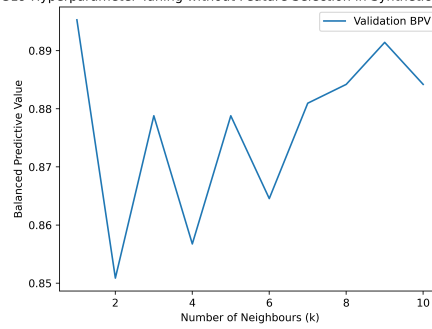
AlexNet Hyperparameter Tuning without Feature Selection in Synthetic Dataset: FC2



VGG19 Hyperparameter Tuning without Feature Selection in Synthetic Dataset: FC1



VGG19 Hyperparameter Tuning without Feature Selection in Synthetic Dataset: FC2



ResNet50 Hyperparameter Tuning without Feature Selection in Synthetic Dataset: AveragePool

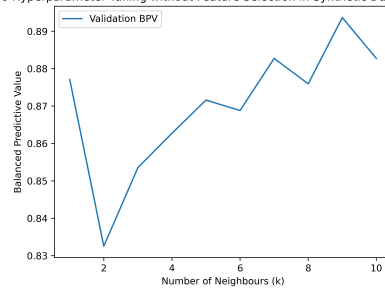


Figure A.2: KNN Classifier Hyperparameter Tuning in Synthetic Dataset.

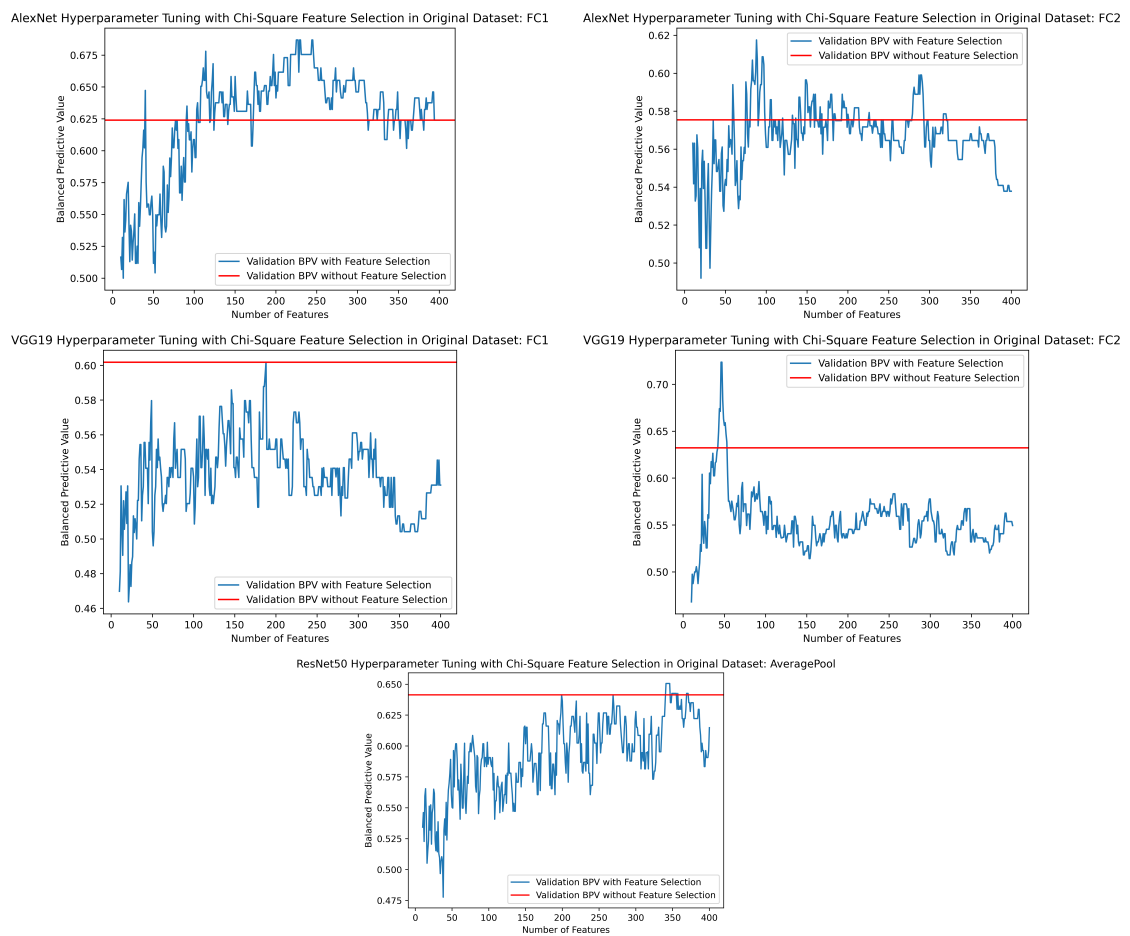


Figure A.3: Feature Selection Hyperparameter Tuning in Original Dataset

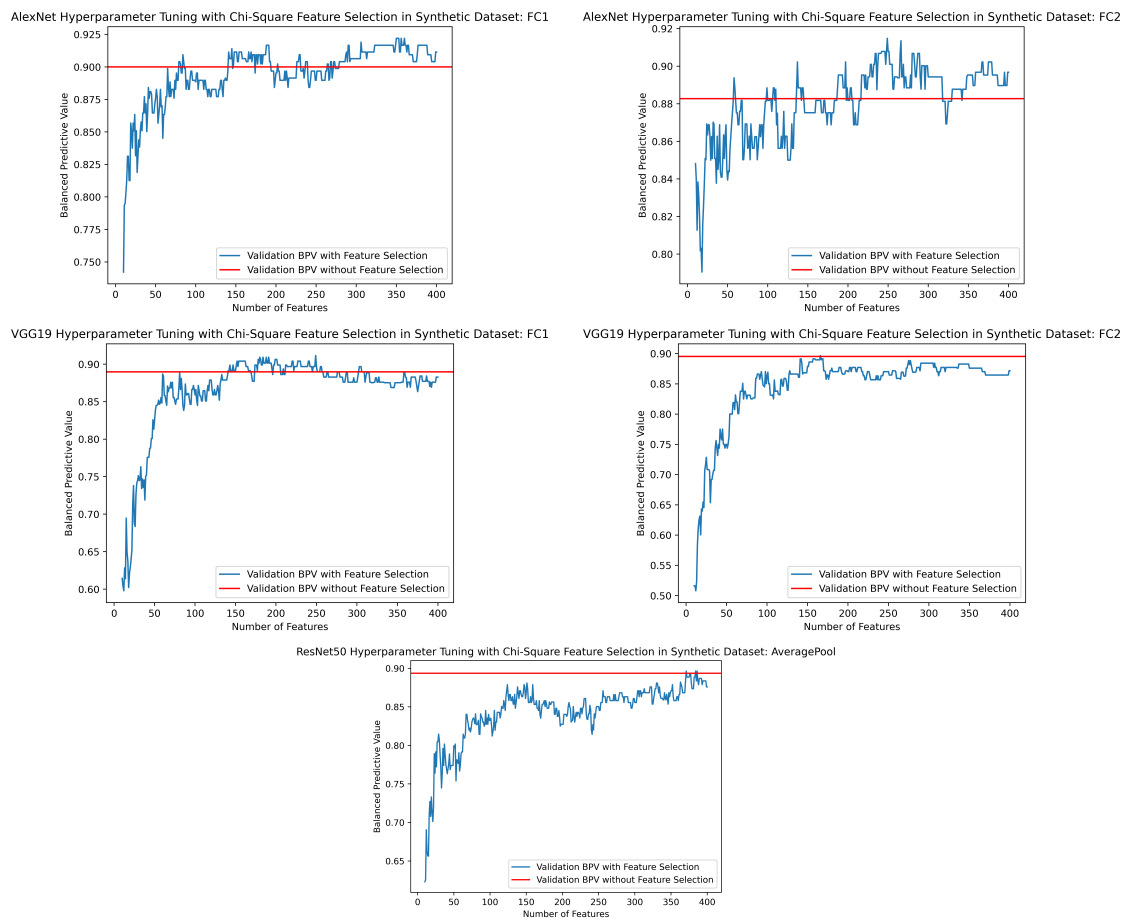


Figure A.4: Feature Selection Hyperparameter Tuning in Synthetic Dataset

## Appendix B

# CNN Transfer Learning-based Model Hyperparameter Tuning and Testing Results: Diagrams and Tables



trial_id	trainable_location	date	time_stamp	time_total_s	total_epochs	time_since_restore	config_batch_size	config_lr
0	FALSE	2023-06-18 16:35:49	165051000	1212.16877	2557764	43	512	0.002301758771
1	TRUE	2023-06-18 16:46:35	165051001	81.1230731	2557764	3	4096	0.003630359205
2	TRUE	2023-06-18 16:47:57	165051002	82.5685716	2557764	3	512	0.002217178746
3	TRUE	2023-06-18 16:49:13	165051003	76.10372654	2557764	3	256	0.001242772164
4	FALSE	2023-06-18 16:50:49	165051004	55.88678921	2557764	2	64	0.00076767619899
5	TRUE	2023-06-18 16:52:02	165051005	83.88647911	2557764	3	512	0.013365939793
6	TRUE	2023-06-18 17:01:57	165051006	371.283212	2557764	37	1024	0.00343302052
7	TRUE	2023-06-18 17:03:05	165051007	371.283212	2557764	37	512	0.00343302052
8	TRUE	2023-06-18 17:25:29	165051008	81.98223025	2557764	3	128	0.01533303318
9	FALSE	2023-06-18 17:30:56	165051009	327.2882318	2557764	11	512	0.0015811781800
10	TRUE	2023-06-18 17:38:37	165051010	74.19143736	2557764	3	256	0.00112546068511
11	FALSE	2023-06-18 17:49:53	165051011	675.6061273	2557764	22	2048	0.00011677424515
12	TRUE	2023-06-18 18:14:03	165060013	80.54112773	2557764	3	256	0.0004008885029
13	FALSE	2023-06-18 18:20:57	165060125	413.92855	2557764	13	256	0.0003014335096
14	TRUE	2023-06-18 18:48:40	165062101	30.83533897	2557764	3	64	0.0004002520534
15	TRUE	2023-06-18 18:49:06	165062102	30.83533897	2557764	3	64	0.0004002520534
16	FALSE	2023-06-18 18:49:46	165062103	128.60385	2557764	4	4096	0.00378332148
17	FALSE	2023-06-18 18:51:15	165063105	64.9788511	2557764	2	512	0.00378332148
18	TRUE	2023-06-18 18:53:15	165063195	90.93019009	2557764	3	128	0.0022767400059
19	FALSE	2023-06-18 18:55:46	165063226	31.12317562	2557764	1	512	0.002828669564
20	FALSE	2023-06-18 18:56:16	165063376	90.96127868	2557764	3	246	0.000618620733
21	FALSE	2023-06-18 19:15:28	165064528	1152.30871	2557764	33	4096	0.000210286851
22	TRUE	2023-06-18 19:25:46	165065116	89.98190122	2557764	3	256	0.000146888437
23	TRUE	2023-06-18 19:26:19	165065209	63.39333828	2557764	2	256	0.0003537943210
24	TRUE	2023-06-18 19:26:59	165065301	63.39333828	2557764	2	256	0.0003537943210
25	FALSE	2023-06-18 19:30:39	165065439	103.5358237	2557764	3	4096	0.004133837518
26	TRUE	2023-06-18 19:32:27	165065517	107.781488	2557764	3	2048	0.0023122316418
27	TRUE	2023-06-18 19:33:57	165065637	89.46723366	2557764	3	512	0.003816522605
28	TRUE	2023-06-18 19:40:06	165066006	279.7443788	2557764	10	64	0.000439123205
29	FALSE	2023-06-18 19:43:57	165066027	177.1907592	2557764	6	4096	0.0008639901321
30	TRUE	2023-06-18 19:45:19	165066319	81.6595919	2557764	3	128	0.00060665775196
31	FALSE	2023-06-18 19:48:49	165066629	209.9001175	2557764	7	256	0.000608477816

Table B.1: AlexNet Hyperparameter Tuning on Original Dataset: best results per trial using validation loss as metric.

























<b>AlexNet Mispronunciation Results: Original Dataset</b>							
<b>Top-5 Best Trials with Validation Loss</b>							
<b>Trials</b>	<b>Epoch</b>	<b>Validation Loss</b>	<b>hn1</b>	<b>hn2</b>	<b>PPV</b>	<b>NPV</b>	<b>BPV</b>
21	33	0.386	4096	512	0.87	0.61	0.74
11	22	0.391	2048	256	0.84	0.64	0.74
9	11	0.393	512	2048	0.81	0.67	0.74
0	43	0.398	512	1024	0.83	0.64	0.73
6	47	0.403	1024	2048	0.84	0.67	0.75
<b>Top-5 Best Trials with BPV</b>							
<b>Trial</b>	<b>Epoch</b>	<b>Validation BPV</b>	<b>hn1</b>	<b>hn2</b>	<b>PPV</b>	<b>NPV</b>	<b>BPV</b>
16	1	0.907	4096	2048	0.80	1.00	0.90
28	8	0.843	64	64	0.81	0.67	0.74
6	10	0.831	1024	2048	0.83	0.78	0.81
0	9	0.796	512	1024	0.82	0.55	0.69
11	6	0.782	2048	256	0.81	0.67	0.74

Table B.13: AlexNet Mispronunciation Results in Original Dataset

<b>AlexNet Mispronunciation Results: Synthetic Dataset</b>							
<b>Top-5 Best Trials with Validation Loss</b>							
<b>Trial</b>	<b>Epoch</b>	<b>Validation Loss</b>	<b>hn1</b>	<b>hn2</b>	<b>PPV</b>	<b>NPV</b>	<b>BPV</b>
12	13	0.22	256	2048	0.83	1.00	0.92
0	45	0.23	128	128	0.86	0.94	0.90
13	13	0.23	2048	128	0.83	0.97	0.90
8	3	0.23	4096	64	0.81	0.97	0.89
7	4	0.23	1024	64	0.82	1.00	0.91
<b>Top-5 Best Trials with BPV</b>							
<b>Trial</b>	<b>Epoch</b>	<b>Validation BPV</b>	<b>hn1</b>	<b>hn2</b>	<b>PPV</b>	<b>NPV</b>	<b>BPV</b>
13	12	0.92	2048	128	0.83	1.00	0.92
7	4	0.92	1024	64	0.82	1.00	0.91
14	4	0.92	1024	64	0.82	1.00	0.91
8	3	0.91	4096	64	0.81	0.97	0.89
24	6	0.91	4096	128	0.87	0.97	0.92

Table B.14: AlexNet Mispronunciation Results in Synthetic Dataset

<b>VGG19 Mispronunciation Results: Original Dataset</b>							
<b>Top-5 Best Trials with Validation Loss</b>							
<b>Trial</b>	<b>Epoch</b>	<b>Validation Loss</b>	<b>hn1</b>	<b>hn2</b>	<b>PPV</b>	<b>NPV</b>	<b>BPV</b>
29	3	0.42	1024	4096	0.83	0.64	0.73
2	64	0.42	128	64	0.83	0.53	0.68
20	4	0.43	1024	64	0.84	0.50	0.67
19	6	0.43	256	2048	0.85	0.52	0.69
10	6	0.43	8192	1024	0.84	0.50	0.67
<b>Top-5 Best Trials with BPV</b>							
<b>Trial</b>	<b>Epoch</b>	<b>Validation BPV</b>	<b>hn1</b>	<b>hn2</b>	<b>PPV</b>	<b>NPV</b>	<b>BPV</b>
5	18	0.91	512	1024	0.81	0.40	0.60
2	1	0.90	128	64	0.80	1.00	0.90
8	7	0.90	32	128	0.81	0.50	0.65
20	1	0.90	1024	64	0.80	1.00	0.90
18	62	0.90	4096	1024	0.81	1.00	0.90

Table B.15: VGG19 Mispronunciation Results in Original Dataset

<b>VGG19 Mispronunciation Results: Synthetic Dataset</b>							
<b>Top-5 Best Trials with Validation Loss</b>							
<b>Trial</b>	<b>Epoch</b>	<b>Validation Loss</b>	<b>hn1</b>	<b>hn2</b>	<b>PPV</b>	<b>NPV</b>	<b>BPV</b>
25	59	0.21	16384	8192	0.83	0.94	0.88
24	3	0.23	1024	128	0.82	0.89	0.85
6	62	0.24	4096	1024	0.82	0.91	0.87
18	12	0.24	1024	4096	0.87	0.83	0.85
2	7	0.24	16384	256	0.83	0.87	0.85
<b>Top-5 Best Trials with BPV</b>							
<b>Trial</b>	<b>Epoch</b>	<b>Validation BPV</b>	<b>hn1</b>	<b>hn2</b>	<b>PPV</b>	<b>NPV</b>	<b>BPV</b>
2	3	0.92	16384	256	0.79	0.91	0.85
25	2	0.92	16384	8192	0.80	1.00	0.90
8	60	0.92	256	256	0.81	0.94	0.87
11	1	0.92	512	8192	0.83	0.94	0.88
17	16	0.92	4096	16384	0.80	0.89	0.84

Table B.16: VGG19 Mispronunciation Results in Synthetic Dataset

<b>ResNet50 Mispronunciation Results: Original Dataset</b>					
<b>Top-5 Best Trials with Validation Loss</b>					
<b>Trial</b>	<b>Epoch</b>	<b>Validation Loss</b>	<b>PPV</b>	<b>NPV</b>	<b>BPV</b>
13	4	0.56	0.80	0.00	0.40
20	20	0.65	0.80	0.00	0.40
27	1	0.73	0.81	0.40	0.60
3	1	0.77	0.80	0.00	0.40
22	1	0.80	0.80	0.00	0.40
<b>Top-5 Best Trials with BPV</b>					
<b>Trial</b>	<b>Epoch</b>	<b>Validation BPV</b>	<b>PPV</b>	<b>NPV</b>	<b>BPV</b>
27	3	0.61	0.80	0.00	0.40
20	28	0.61	0.80	0.00	0.40
23	3	0.61	0.80	0.00	0.40
5	37	0.60	0.80	0.00	0.40
26	6	0.60	0.80	0.00	0.40

Table B.17: ResNet50 Mispronunciation Results in Original Dataset

<b>ResNet50 Mispronunciation Results: Synthetic Dataset</b>					
<b>Top-5 Best Trials with Validation Loss</b>					
<b>Trial</b>	<b>Epoch</b>	<b>Validation Loss</b>	<b>PPV</b>	<b>NPV</b>	<b>BPV</b>
3	4	0.55	0.52	1.00	0.76
16	2	0.55	0.61	1.00	0.80
9	2	0.56	0.66	1.00	0.83
5	60	0.56	0.68	1.00	0.84
0	19	0.56	0.63	1.00	0.82
<b>Top-5 Best Trials with BPV</b>					
<b>Trial</b>	<b>Epoch</b>	<b>Validation BPV</b>	<b>PPV</b>	<b>NPV</b>	<b>BPV</b>
22	2	0.78	0.65	1.00	0.82
28	1	0.78	0.51	1.00	0.76
0	19	0.78	0.63	1.00	0.82
5	42	0.77	0.66	1.00	0.83
3	4	0.77	0.52	1.00	0.76

Table B.18: ResNet50 Mispronunciation Results in Synthetic Dataset

# Appendix C

## Exemplary Code

This appendix is devoted to present the code used in the graduation project. In the case of the CNN Features-based model, feature extraction for all models is shown. However, in the case of the CNN Transfer Learning-based model, part of the code shown for AlexNet is already identical to the other models, but with some modifications. Therefore, we only show the CNN Transfer Learning-based model using AlexNet architecture.

### C.1 Dataset Splitting

```
1 import splitfolders
2
3 #Split original dataset
4 input_folder = '/home/ricardo.velasco/TT_kichwa/resources/ds'
5 output_folder = "/home/ricardo.velasco/TT_kichwa/
6                 transfer_learning_model/split_ds"
7
8 #Stratified sampling for imbalanced data
9 splitfolders.ratio(input_folder, output=output_folder, seed=42, ratio
10                   =(.34, .33, .33), group_prefix=None)
11
12 #Split syntetic dataset
13 input_folder = '/home/ricardo.velasco/TT_kichwa/resources/syn_ds'
14 output_folder = "/home/ricardo.velasco/TT_kichwa/
15                 transfer_learning_model/split_syn_ds"
16
17 splitfolders.ratio(input_folder, output=output_folder, seed=42, ratio
18                   =(.7, .2, .1), group_prefix=None)
```



## C.2 Audio Dataset Augmentation, Resampling, and Conversion

```
1 from glob import glob
2 from audiomentations import Compose, AddGaussianNoise, TimeStretch,
   PitchShift, Shift, AdjustDuration, TimeMask
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import torchaudio
6 import os
7 import shutil
8 import torch
9 import librosa
10
11 default_sample_rate = 44100
12 augment_size = 100
13
14 #Define augmentation parameters
15 audio_augmentation = Compose([
16     AdjustDuration(duration_seconds=3, p=1.0),
17     Shift(min_shift=-2, max_shift=2, p=1.0, shift_unit="seconds",
18         rollover=True),
19     PitchShift(min_semitones=-4, max_semitones=4, p=0.5),
20     TimeStretch(min_rate=0.5, max_rate=1.5, p=0.5),
21     AddGaussianNoise(min_amplitude=0.001, max_amplitude=0.015, p=0.5),
22     TimeMask(min_band_part=0.05, max_band_part=0.15, fade=True, p=0.5)
23 ],
24 ])
25
26 #Given a training folder of audios, augment them and convert them into
27 spectrograms
28 def spectrogramAugmentation(path, save):
29     #Define spectrogram function
30     toSpectrogram = torchaudio.transforms.Spectrogram(n_fft=512)
31
32     #Create a list of .ogg and .wav audio files to convert
33     audio_files = glob(os.path.join(path, "*.ogg")) + glob(os.path.
34         join(path, "*.wav"))
35
36     for audio in audio_files:
37         #Extract audio name for spectrogram image name
38         spectrogram_name = (audio.split("/")[-1]).split(".")[0]
```

```

36     #Read audio file (normalization takes place by default)
37     waveform, sample_rate = torchaudio.load(audio)
38
39     #Define resampling audio function and resample
40     resamplingAudio = torchaudio.transforms.Resample(sample_rate,
41     default_sample_rate, dtype=waveform.dtype)
42     waveform = resamplingAudio(waveform)
43
44     #Save augmentations in .png format (100 augmentations for each
45     audio file)
46     for i in range(1, augment_size + 1):
47         #Augment audio file
48         aug = audio_augmentation(samples = np.asarray(waveform),
49         sample_rate=default_sample_rate)
50
51         #Transform raw audios into spectrograms
52         audio_spec = toSpectrogram(torch.from_numpy(np.asarray(aug
53         )))
54
55         #Save spectrograms in save folder
56         plt.imsave(os.path.join(save, spectrogram_name + "_aug_" +
57         str(i)) + ".png", librosa.power_to_db(audio_spec[0]))
58
59 #Given a validation and testing folder of audios, convert them into
60 spectrograms
61 def spectrogram(path, save):
62     #Define spectrogram function
63     toSpectrogram = torchaudio.transforms.Spectrogram(n_fft=512)
64
65     #Create a list of .ogg and .wav audio files to convert
66     audio_files = glob(os.path.join(path, "*.ogg")) + glob(os.path.
67     join(path, "*.wav"))
68
69     for audio in audio_files:
70         #Extract audio name for spectrogram image name
71         spectrogram_name = (audio.split("/")[-1]).split(".")[0]
72
73         #Read audio file (normalization takes place by default)
74         waveform, sample_rate = torchaudio.load(audio)
75
76         #Define resampling audio function and resample
77         resamplingAudio = torchaudio.transforms.Resample(sample_rate,
78         default_sample_rate, dtype=waveform.dtype)

```

```

72     waveform = resamplingAudio(waveform)
73
74     #Make all audios have a fixed duration of 3 seconds
75     adjust_duration = AdjustDuration(duration_seconds=3, p=1.0)
76     waveform = adjust_duration(np.asarray(waveform), sample_rate=
default_sample_rate)
77
78     #Transform raw audios into spectrograms
79     audio_spec = toSpectrogram(torch.from_numpy(np.asarray(
waveform)))
80
81     #Save spectrogram in .png format
82     plt.imsave(os.path.join(save, spectrogram_name) + ".png",
librosa.power_to_db(audio_spec[0]))
83
84 #-----
85 # DATASET CONVERSION
86 #-----
87
88 #Convert training data with 0 and 1 label
89 spectrogramAugmentation("/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/split_ds/train/0",
"/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/spec/train/0")
91
92
93 spectrogramAugmentation("/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/split_ds/train/1",
"/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/spec/train/1")
95
96 #Convert validation data with 0 and 1 label
97 spectrogram("/home/ricardo.velasco/TT_kichwa/transfer_learning_model/
split_ds/val/0",
"/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/spec/val/0")
99
100 spectrogram("/home/ricardo.velasco/TT_kichwa/transfer_learning_model/
split_ds/val/1",
"/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/spec/val/1")
102
103 #Convert testing data with 0 and 1 label
104 spectrogram("/home/ricardo.velasco/TT_kichwa/transfer_learning_model/

```

```

split_ds/test/0",
105         "/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/spec/test/0")
106
107 spectrogram("/home/ricardo.velasco/TT_kichwa/transfer_learning_model/
split_ds/test/1",
108         "/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/spec/test/1")
109
110
111 #-----
112 # SYNTHETIC DATASET CONVERSION
113 #-----
114
115 #Convert training data with 0 and 1 label
116 spectrogramAugmentation("/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/split_syn_ds/train/0",
117         "/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/syn_spec/train/0")
118
119
120 spectrogramAugmentation("/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/split_syn_ds/train/1",
121         "/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/syn_spec/train/1")
122
123 #Convert validation data with 0 and 1 label
124 spectrogram("/home/ricardo.velasco/TT_kichwa/transfer_learning_model/
split_syn_ds/val/0",
125         "/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/syn_spec/val/0")
126
127 spectrogram("/home/ricardo.velasco/TT_kichwa/transfer_learning_model/
split_syn_ds/val/1",
128         "/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/syn_spec/val/1")
129
130 #Convert testing data with 0 and 1 label
131 spectrogram("/home/ricardo.velasco/TT_kichwa/transfer_learning_model/
split_syn_ds/test/0",
132         "/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/syn_spec/test/0")
133
134 spectrogram("/home/ricardo.velasco/TT_kichwa/transfer_learning_model/

```

```

split_syn_ds/test/1",
135         "/home/ricardo.velasco/TT_kichwa/
transfer_learning_model/syn_spec/test/1")

```

## C.3 CNN Features-based Model

```

1 from collections import OrderedDict
2 from glob import glob
3 from PIL import Image
4 import os
5 import torch
6 import torchvision
7 import numpy as np
8 import torchextractor
9 import itertools
10 import Orange
11 import matplotlib
12 import matplotlib.pyplot as plt
13
14 from torchvision.datasets import ImageFolder
15 from torch.utils.data import DataLoader
16 from torchvision.transforms import transforms
17 from sklearn.metrics import precision_score, make_scorer, log_loss
18 from sklearn.neighbors import KNeighborsClassifier
19
20 #Feature selection methods
21 from skfeature.function.statistical_based.chi_square import chi_square
22     , feature_ranking
23
24 matplotlib.rcParams["figure.dpi"] = 500
25
26 num_cores = 64
27 device = "cuda" if torch.cuda.is_available() else "cpu"
28
29 #-----
30 #Alexnet Feature Extraction
31 #-----
32
33 def alexnetFeatureExtraction(path):
34     #Define resampling function
35     resample = transforms.Compose([

```

```

36         transforms.Resize(size=(224, 224)),
37         transforms.ToTensor(),
38         transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229,
0.224, 0.225]),
39     ])
40
41     # Array containing datasets using corresponding layers
42     d = np.array([{"train":np.empty([0, 4096]), "val":np.empty([0,
4096])}, {"train":np.empty([0, 4096])}, {"train":np.empty([0, 4096]), "val":np.empty([0,
4096])}, {"train":np.empty([0, 4096])}])
43
44
45     l = {"train":np.array([]), "val":np.array([]), "test":np.array([])
}
46
47     #Target layers in order to extract features, select inference mode
and send it to GPU
48     layers = ["classifier.1", "classifier.4"]
49     model = torchvision.models.alexnet(weights='AlexNet_Weights.
IMAGENET1K_V1')
50     model = torchextractor.Extractor(model, layers)
51     model.eval()
52     model.cuda()
53
54     #Perform feature extraction by batches
55     for i in ["train", "val", "test"]:
56         data = ImageFolder(os.path.join(path, i), transform=resample)
57         dataloader = DataLoader(dataset=data, batch_size=1024,
num_workers=num_cores)
58         for j, (inputs, labels) in enumerate(dataloader):
59             with torch.no_grad():
60                 inputs = inputs.to(device)
61                 outputs, features = model(inputs)
62
63                 f0 = features[layers[0]].detach().cpu().numpy()
64                 f1 = features[layers[1]].detach().cpu().numpy()
65
66                 d[0][i] = np.concatenate((d[0][i], f0))
67                 d[1][i] = np.concatenate((d[1][i], f1))
68                 l[i] = np.concatenate((l[i], labels.numpy()))
69
70     return d, l
71
72 #-----

```

```

73 #VGG19 Feature Extraction
74 #-----
75
76 def vggFeatureExtraction(path):
77     #Define resampling function
78     resample = transforms.Compose([
79         transforms.Resize(size=(224, 224)),
80         transforms.ToTensor(),
81         transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229,
0.224, 0.225]),
82     ])
83
84     # Array containing datasets using corresponding layers
85     d = np.array([{"train":np.empty([0, 4096]), "val":np.empty([0,
4096])}, {"train":np.empty([0, 4096]), "val":np.empty([0,
4096])}, {"train":np.empty([0, 4096]), "val":np.empty([0,
4096])}], "test":np.empty([0, 4096])}]
86
87
88     l = {"train":np.array([]), "val":np.array([]), "test":np.array([])
}
89
90     #Target layers in order to extract features, select inference mode
and send it to GPU
91     layers = ["classifier.0", "classifier.3"]
92     model = torchvision.models.vgg19(weights='VGG19_Weights.
IMAGENET1K_V1')
93     model = torchextractor.Extractor(model, layers)
94     model.eval()
95     model.cuda()
96
97     #Perform feature extraction by batches
98     for i in ["train", "val", "test"]:
99         data = ImageFolder(os.path.join(path, i), transform=resample)
100         dataloader = DataLoader(dataset=data, batch_size=1024,
num_workers=num_cores)
101         for j, (inputs, labels) in enumerate(dataloader):
102             with torch.no_grad():
103                 inputs = inputs.to(device)
104                 outputs, features = model(inputs)
105
106                 f0 = features[layers[0]].detach().cpu().numpy()
107                 f1 = features[layers[1]].detach().cpu().numpy()
108
109                 d[0][i] = np.concatenate((d[0][i], f0))

```

```

110         d[1][i] = np.concatenate((d[1][i], f1))
111         l[i] = np.concatenate((l[i], labels.numpy()))
112
113     return d, l
114
115
116 #-----
117 #VGG19 Feature Extraction
118 #-----
119
120 def resnetFeatureExtraction(path):
121     #Define resampling function
122     resample = transforms.Compose([
123         transforms.Resize(size=(224, 224)),
124         transforms.ToTensor(),
125         transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229,
126     0.224, 0.225]),
127     ])
128
129     # Array containing datasets using corresponding layers
130     d = {"train":np.empty([0, 2048]), "val":np.empty([0, 2048]), "test
131     ":np.empty([0, 2048])}
132
133     l = {"train":np.array([]), "val":np.array([]), "test":np.array([])
134     }
135
136     #Target layers in order to extract features, select inference mode
137     #and send it to GPU
138     layers = ["avgpool"]
139     model = torchvision.models.resnet50(weights="ResNet50_Weights.
140     IMAGENET1K_V2")
141     model = torchextractor.Extractor(model, layers)
142     model.eval()
143     model.cuda()
144
145     #Perform feature extraction by batches
146     for i in ["train", "val", "test"]:
147         data = ImageFolder(os.path.join(path, i), transform=resample)
148         dataloader = DataLoader(dataset=data, batch_size=1024,
149     num_workers=num_cores)
150         for j, (inputs, labels) in enumerate(dataloader):
151             with torch.no_grad():
152                 inputs = inputs.to(device)
153                 outputs, features = model(inputs)

```



```

148
149         f0 = features[layers[0]].detach().cpu().numpy()
150
151         d[i] = np.concatenate((d[i], np.squeeze(f0)))
152         l[i] = np.concatenate((l[i], labels.numpy()))
153
154     return d, l
155
156 #-----
157 # Discretization function
158 #-----
159
160 #Discretize continuous data into categorical data by method of Fayyad
    and Irani [FI93]. Returns the discretized dataset.
161 def discretize(X, y):
162     domain = Orange.data.Domain.from_numpy(X, y)
163     data = Orange.data.Table.from_numpy(domain, X, y)
164
165     disc = Orange.preprocess.Discretize() #Define discretize function
166     disc.method = Orange.preprocess.discretize.EntropyMDL() #Select
    desired method
167
168     discrete_data = disc(data) #Discretize dataset
169
170     return discrete_data.X
171
172 #Discretize the entire dataset and split the training, validation and
    testing sets again
173 def discretizeDataset(X_train, X_val, X_test, y_train, y_val, y_test):
174     #Concatenate all sets to form a unique dataset for discretization
175     X = np.concatenate((X_train, X_val, X_test))
176     y = np.concatenate((y_train, y_val, y_test))
177
178     #Discretize the dataset
179     X_disc = discretize(X, y)
180
181     #Retrieve discretized datasets from training, validation and
    testing sets
182     X_train_disc = X_disc[:X_train.shape[0], :]
183     X_val_disc = X_disc[X_train.shape[0]:X_train.shape[0] + X_val.
    shape[0], :]
184     X_test_disc = X_disc[X_train.shape[0] + X_val.shape[0]:, :]
185
186     return X_train_disc, X_val_disc, X_test_disc

```

```

187
188
189 def metric(y_true, y_pred):
190     ppv = precision_score(y_true, y_pred)
191     npv = precision_score(y_true, y_pred, pos_label=0)
192     bpv = (ppv + npv) / 2
193
194     return ppv, npv, bpv
195
196
197 def knnWithoutFeatureSelection(X_train, X_val, X_test, y_train, y_val,
198     y_test, title):
199     val_results = np.array([])
200     val_best = 0
201     best_k = 0
202
203     for k in range(1, 11):
204         #Declare KNN Neighbours Classifier
205         knn_classifier = KNeighborsClassifier(n_neighbors=k, n_jobs =
206     num_cores)
207
208         #Train the model with training set.
209         knn_classifier.fit(X_train, y_train)
210
211         #Obtain prediction using validation set and testing set
212         y_pred_val = knn_classifier.predict(X_val)
213         y_pred_test = knn_classifier.predict(X_test)
214
215         #Obtain validation bpv
216         val_bpv = metric(y_val, y_pred_val)[2]
217
218         #Save best k with best bpv
219         if val_best <= val_bpv:
220             test_ppv, test_npv, test_bpv = metric(y_test, y_pred_test)
221             val_best = val_bpv
222             best_k = k
223
224         val_results = np.append(val_results, val_bpv)
225
226     print("PPV: ", test_ppv, "NPV:", test_npv, "BPV: ", test_bpv, "k: "
227     , best_k, "val_bpv", val_best)
228     plt.plot([i for i in range(1,11)], val_results, label="Validation
229     BPV")

```

```

227 plt.xlabel("Number of Neighbours (k)")
228 plt.ylabel("Balanced Predictive Value")
229 plt.title(title)
230 plt.legend()
231 plt.show()
232
233
234 def knnChi(X_train, X_val, X_test, y_train, y_val, y_test, k ,title,
last):
235     #Create array to save results
236     val_results = np.array([])
237     val_best = 0
238     best_number_features = 0
239
240     #Discretize datasets
241     X_train, X_val, X_test = discretizeDataset(X_train, X_val, X_test,
y_train, y_val, y_test)
242
243     #Get selected features by ranking
244     ranking = feature_ranking(chi_square(X_train, y_train))
245
246     for i in range(10, 401):
247         #Declare KNN Neighbours Classifier
248         knn_classifier = KNeighborsClassifier(n_neighbors=k, n_jobs =
num_cores)
249
250         #Best ith features
251         best_features = ranking[:i]
252
253         #Train the model with training set.
254         knn_classifier.fit(X_train[:,best_features], y_train)
255
256         #Obtain prediction using validation set and testing set
257         y_pred_val = knn_classifier.predict(X_val[:,best_features])
258         y_pred_test = knn_classifier.predict(X_test[:,best_features])
259
260         #Obtain metrics for validation and testing set
261         val_bpv = metric(y_val, y_pred_val)[2]
262
263         #Save best k with best bpv
264         if val_best <= val_bpv:
265             test_ppv, test_npv, test_bpv = metric(y_test, y_pred_test)
266             val_best = val_bpv
267             best_number_features = i

```

```

268
269         val_results = np.append(val_results, val_bpv)
270
271
272         print("PPV: ", test_ppv, "NPV:", test_npv, "BPV: ", test_bpv, "
num_features: ", best_number_features, "val_bpv", val_best)
273         plt.plot([i for i in range(10,401)], val_results, label="
Validation BPV with Feature Selection")
274         plt.axhline(y = last, color = 'r', linestyle = '--', label = "
Validation BPV without Feature Selection")
275         plt.xlabel("Number of Features")
276         plt.ylabel("Balanced Predictive Value")
277         plt.title(title)
278         plt.legend()
279         plt.show()
280
281 #Obtain dataset matrix and labels for dataset
282 alexData, alexLabels = alexnetFeatureExtraction("/home/ricardo.velasco
/TT_kichwa/transfer_learning_model/spec")
283 vggData, vggLabels = vggFeatureExtraction("/home/ricardo.velasco/
TT_kichwa/transfer_learning_model/spec")
284 resData, resLabels = resnetFeatureExtraction("/home/ricardo.velasco/
TT_kichwa/transfer_learning_model/spec")
285
286 #Obtain transformed dataset matrix and labels for syntetic dataset
287 synalexData, synalexLabels = alexnetFeatureExtraction("/home/ricardo.
velasco/TT_kichwa/transfer_learning_model/syn_spec")
288 synvggData, synvggLabels = vggFeatureExtraction("/home/ricardo.velasco
/TT_kichwa/transfer_learning_model/syn_spec")
289 synresData, synresLabels = resnetFeatureExtraction("/home/ricardo.
velasco/TT_kichwa/transfer_learning_model/syn_spec")
290
291
292 #KNN HyperParameter on Original Dataset, plot graphics and test
parameter with testing set
293
294 knnWithoutFeatureSelection(alexData[0]["train"], alexData[0]["val"],
alexData[0]["test"], alexLabels["train"], alexLabels["val"],
alexLabels["test"], "AlexNet Hyperparameter Tuning without Feature
Selection in Original Dataset: FC1")
295 knnWithoutFeatureSelection(alexData[1]["train"], alexData[1]["val"],
alexData[1]["test"], alexLabels["train"], alexLabels["val"],
alexLabels["test"], "AlexNet Hyperparameter Tuning without Feature
Selection in Original Dataset: FC2")

```

```

296 knnWithoutFeatureSelection(vggData[0]["train"], vggData[0]["val"],
    vggData[0]["test"], vggLabels["train"], vggLabels["val"], vggLabels
    ["test"], "VGG19 Hyperparameter Tuning without Feature Selection in
    Original Dataset: FC1")
297 knnWithoutFeatureSelection(vggData[1]["train"], vggData[1]["val"],
    vggData[1]["test"], vggLabels["train"], vggLabels["val"], vggLabels
    ["test"], "VGG19 Hyperparameter Tuning without Feature Selection in
    Original Dataset: FC2")
298 knnWithoutFeatureSelection(resData["train"], resData["val"], resData["
    test"], resLabels["train"], resLabels["val"], resLabels["test"], "
    ResNet50 Hyperparameter Tuning without Feature Selection in
    Original Dataset: AveragePool")
299
300 #KNN HyperParameter on Synthetic Dataset, plot graphics and test
    parameter with testing set
301 knnWithoutFeatureSelection(synalexData[0]["train"], synalexData[0]["
    val"], synalexData[0]["test"], synalexLabels["train"],
    synalexLabels["val"], synalexLabels["test"], "AlexNet
    Hyperparameter Tuning without Feature Selection in Synthetic
    Dataset: FC1")
302 knnWithoutFeatureSelection(synalexData[1]["train"], synalexData[1]["
    val"], synalexData[1]["test"], synalexLabels["train"],
    synalexLabels["val"], synalexLabels["test"], "AlexNet
    Hyperparameter Tuning without Feature Selection in Synthetic
    Dataset: FC2")
303 knnWithoutFeatureSelection(synvggData[0]["train"], synvggData[0]["val"
    ], synvggData[0]["test"], synvggLabels["train"], synvggLabels["val"
    ], synvggLabels["test"], "VGG19 Hyperparameter Tuning without
    Feature Selection in Synthetic Dataset: FC1")
304 knnWithoutFeatureSelection(synvggData[1]["train"], synvggData[1]["val"
    ], synvggData[1]["test"], synvggLabels["train"], synvggLabels["val"
    ], synvggLabels["test"], "VGG19 Hyperparameter Tuning without
    Feature Selection in Synthetic Dataset: FC2")
305 knnWithoutFeatureSelection(synresData["train"], synresData["val"],
    synresData["test"], synresLabels["train"], synresLabels["val"],
    synresLabels["test"], "ResNet50 Hyperparameter Tuning without
    Feature Selection in Synthetic Dataset: AveragePool")
306
307
308 #KNN Feature Selection in Original Dataset with Chi Score
309 knnChi(alexData[0]["train"], alexData[0]["val"], alexData[0]["test"],
    alexLabels["train"], alexLabels["val"], alexLabels["test"], 9 , "
    AlexNet Hyperparameter Tuning with Chi-Square Feature Selection in
    Original Dataset: FC1", 0.6239935587761675)

```

```

310 knnChi(alexData[1]["train"], alexData[1]["val"], alexData[1]["test"],
      alexLabels["train"], alexLabels["val"], alexLabels["test"], 4 ,
      AlexNet Hyperparameter Tuning with Chi-Square Feature Selection in
      Original Dataset: FC2", 0.5754750175932442)
311
312 knnChi(vggData[0]["train"], vggData[0]["val"], vggData[0]["test"],
      vggLabels["train"], vggLabels["val"], vggLabels["test"], 9 ,"VGG19
      Hyperparameter Tuning with Chi-Square Feature Selection in Original
      Dataset: FC1", 0.6018518518518519)
313 knnChi(vggData[1]["train"], vggData[1]["val"], vggData[1]["test"],
      vggLabels["train"], vggLabels["val"], vggLabels["test"], 5 ,"VGG19
      Hyperparameter Tuning with Chi-Square Feature Selection in Original
      Dataset: FC2", 0.6323902027027027)
314
315 knnChi(resData["train"], resData["val"], resData["test"], resLabels["
      train"], resLabels["val"], resLabels["test"], 5 ,"ResNet50
      Hyperparameter Tuning with Chi-Square Feature Selection in Original
      Dataset: AveragePool", 0.6414285714285715)
316
317 #KNN Feature Selection in Synthetic Dataset with Chi Score
318 knnChi(synalexData[0]["train"], synalexData[0]["val"], synalexData[0][
      "test"], synalexLabels["train"], synalexLabels["val"],
      synalexLabels["test"], 9 ,"AlexNet Hyperparameter Tuning with Chi-
      Square Feature Selection in Synthetic Dataset: FC1", 0.9)
319 knnChi(synalexData[1]["train"], synalexData[1]["val"], synalexData[1][
      "test"], synalexLabels["train"], synalexLabels["val"],
      synalexLabels["test"], 9 ,"AlexNet Hyperparameter Tuning with Chi-
      Square Feature Selection in Synthetic Dataset: FC2",
      0.8827450980392156)
320
321 knnChi(synvggData[0]["train"], synvggData[0]["val"], synvggData[0]["
      test"], synvggLabels["train"], synvggLabels["val"], synvggLabels["
      test"], 10 ,"VGG19 Hyperparameter Tuning with Chi-Square Feature
      Selection in Synthetic Dataset: FC1", 0.8896920175989944)
322 knnChi(synvggData[1]["train"], synvggData[1]["val"], synvggData[1]["
      test"], synvggLabels["train"], synvggLabels["val"], synvggLabels["
      test"], 1 ,"VGG19 Hyperparameter Tuning with Chi-Square Feature
      Selection in Synthetic Dataset: FC2", 0.8952941176470588)

```

## C.4 CNN Transfer Learning-based AlexNet Hyperparameter Tuning

```

1 from functools import partial
2 import os
3 import torch
4 import torch.nn as nn
5 import torch.nn.functional as F
6 import torch.optim as optim
7 import torchvision
8 import torchvision.transforms as transforms
9 from torchvision.datasets import ImageFolder
10 from torch.utils.data import DataLoader
11 import ray
12
13 from ray import tune
14 from ray.air import session, RunConfig, CheckpointConfig, Checkpoint
15 from ray.tune.schedulers import ASHAScheduler
16
17 device = "cuda" if torch.cuda.is_available() else "cpu"
18
19 num_cores = 64
20 gpus_per_trial = 1
21 max_num_epochs = 64
22 num_samples = 32 #Number of items in hyperparameter search
23
24 trainingPath = "/home/ricardo.velasco/TT_kichwa/
    transfer_learning_model/spec/train"
25 validationPath = "/home/ricardo.velasco/TT_kichwa/
    transfer_learning_model/spec/val"
26 testingPath = "/home/ricardo.velasco/TT_kichwa/transfer_learning_model
    /spec/test"
27
28 synTrainingPath = "/home/ricardo.velasco/TT_kichwa/
    transfer_learning_model/syn_spec/train"
29 synValidationPath = "/home/ricardo.velasco/TT_kichwa/
    transfer_learning_model/syn_spec/val"
30 synTestingPath = "/home/ricardo.velasco/TT_kichwa/
    transfer_learning_model/syn_spec/test"
31
32 #Load training, testing and validation function for hyperparameter
    searching
33 def load_data(train, val):
34
35     #Input resampling
36     transform = transforms.Compose([
37         transforms.Resize(size=(224, 224)),

```

```

38     transforms.ToTensor(),
39     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229,
40     0.224, 0.225]),
41 ]
42
43     trainset = ImageFolder(train, transform=transform)
44
45     valset = ImageFolder(val, transform=transform)
46
47     return trainset, valset
48
49 def metrics(outputs, labels):
50     TP = ((outputs == 1) & (labels == 1)).sum().item()
51     FP = ((outputs == 1) & (labels == 0)).sum().item()
52     TN = ((outputs == 0) & (labels == 0)).sum().item()
53     FN = ((outputs == 0) & (labels == 1)).sum().item()
54
55     ppv = (TP / (TP + FP)) if (TP + FP) > 0 else 0
56     npv = (TN / (TN + FN)) if (TN + FN) > 0 else 0
57
58     return ppv, npv
59
60 def training(config, train=None, val=None):
61     #Import Alexnet pretrained model
62     net = torchvision.models.alexnet(weights='AlexNet_Weights.
63     IMAGENET1K_V1')
64
65     #Freeze features' layer for transfer learning
66     for param in net.features.parameters():
67         param.requires_grad = False
68
69     #Reset and modify classification layer suitable for binary
70     classification
71     net.classifier = torch.nn.Sequential(
72         torch.nn.Dropout(p=0.5, inplace=False),
73         torch.nn.Linear(in_features=9216, out_features=config["hn1"],
74         bias=True),
75         torch.nn.ReLU(inplace=True),
76         torch.nn.Dropout(p=0.5, inplace=False),
77         torch.nn.Linear(in_features=config["hn1"], out_features=config
78         ["hn2"], bias=True),
79         torch.nn.ReLU(inplace=True),
80         torch.nn.Linear(in_features=config["hn2"], out_features=2,

```



```

bias=True))
77
78 #Send the model to GPU and put it in training mode
79 net.to(device)
80 net.train()
81
82 #Run always on GPU, if there are more GPUs, allow parallel data
parallel training
83 if torch.cuda.device_count() > 1:
84     net = nn.DataParallel(net)
85
86 #Define Loss Function
87 criterion = nn.CrossEntropyLoss()
88
89 #Define optimizer and establish a trainable parameter for
learning rate
90 optimizer = optim.SGD(net.classifier.parameters(), lr=config["lr"
], momentum=0.9)
91
92 # To restore a checkpoint, use 'session.get_checkpoint()'.
93 checkpoint = session.get_checkpoint()
94
95 if checkpoint:
96     checkpoint_state = checkpoint.to_dict()
97     start_epoch = checkpoint_state["epoch"]
98     net.classifier.load_state_dict(checkpoint_state["
classifier_state_dict"])
99     optimizer.load_state_dict(checkpoint_state["
optimizer_state_dict"])
100 else:
101     start_epoch = 0
102
103
104 #Import training, validation and testing datasets
105 trainset, valset= load_data(train, val)
106
107 #Load training dataset
108 trainloader = DataLoader(
109     dataset=trainset,
110     batch_size=int(config["batch_size"]),
111     num_workers=num_cores,
112     shuffle = True
113 )
114

```

```

115     #Load validation dataset
116     valloader = DataLoader(
117         dataset=valset,
118         batch_size=len(valset), #Due to hardware, pick largest batch
size possible to compute metrics
119         num_workers=num_cores,
120         shuffle = True
121     )
122
123
124     for epoch in range(start_epoch, max_num_epochs):
125         # -----
126         # TRAINING STEP
127         # -----
128         train_loss, train_ppv, train_npv = 0.0, 0.0, 0.0
129         for i, (inputs, labels) in enumerate(trainloader):
130             # get the inputs; data is a list of [inputs, labels]
131             inputs, labels = inputs.to(device), labels.to(device)
132
133             # zero the parameter gradients
134             optimizer.zero_grad()
135
136             # forward + backward + optimize
137             outputs = net(inputs)
138             loss = criterion(outputs, labels)
139             loss.backward()
140             optimizer.step()
141
142             #Accumulate statistics across a single batch
143             _, y = torch.max(outputs.data, 1)
144             train_loss += loss.item()
145             train_ppv += metrics(y, labels)[0]
146             train_npv += metrics(y, labels)[1]
147
148         #Compute average of metrics across all batches
149         train_loss /= len(trainloader)
150         train_ppv /= len(trainloader)
151         train_npv /= len(trainloader)
152
153         #-----
154         # VALIDATION STEP
155         #-----
156         val_loss, val_ppv, val_npv = 0.0, 0.0, 0.0
157

```

```

158     for i, (inputs, labels) in enumerate(valloader):
159         with torch.no_grad():
160             inputs, labels = inputs.to(device), labels.to(device)
161
162             #Compute outputs
163             outputs = net(inputs)
164             _, y = torch.max(outputs.data, 1)
165
166             #Compute loss
167             loss = criterion(outputs, labels)
168
169             #Accumulate statistics across a single batch
170             val_loss+= loss.item()
171             val_ppv += metrics(y, labels)[0]
172             val_npv += metrics(y, labels)[1]
173
174             #Compute average of metrics across all batches
175             val_loss /= len(valloader)
176             val_ppv /= len(valloader)
177             val_npv /=len(valloader)
178
179             # Here we save a checkpoint. It is automatically registered
with
180             # Ray Tune and can be accessed through 'session.get_checkpoint
()‘
181
182             # API in future iterations.
183             #Checkpoint data
184             checkpoint_data = {
185                 "epoch": epoch,
186                 "classifier_state_dict": net.classifier.state_dict(),
187                 "optimizer_state_dict": optimizer.state_dict()
188             }
189
190             checkpoint = Checkpoint.from_dict(checkpoint_data)
191             session.report({"tr_loss":train_loss, "v_loss": val_loss,
192                 "tr_bpv": (train_ppv + train_npv) / 2,
193                 "v_bpv": (val_ppv + val_npv)/2,
194                 "tr_ppv":train_ppv, "tr_npv":train_npv,
195                 "v_ppv": val_ppv, "v_npv": val_npv},
checkpoint=checkpoint)
196
197 #-----
198 #Hyperparameter Tuning and Testing for Dataset

```

```

199 #-----
200
201 config = {
202     "lr": tune.loguniform(1e-4, 1e-1),
203     "batch_size": tune.choice([64, 128, 256, 512]),
204     "hn1": tune.choice([32, 64, 128, 256, 512, 1024, 2048, 4096]),
205     "hn2": tune.choice([32, 64, 128, 256, 512, 1024, 2048, 4096])
206 }
207
208 scheduler = ASHAScheduler(
209     max_t=max_num_epochs,
210     grace_period=3,
211     reduction_factor=2,)
212
213
214 tuner = tune.Tuner(
215     tune.with_resources(
216         tune.with_parameters(training, train=trainingPath, val=
validationPath),
217         resources={"cpu": num_cores, "gpu": gpus_per_trial}
218     ),
219     tune_config=tune.TuneConfig(
220         metric="v_loss",
221         mode="min",
222         scheduler=scheduler,
223         num_samples=num_samples,
224     ),
225     run_config=RunConfig(
226         checkpoint_config=CheckpointConfig(
227             num_to_keep = None, #SAVE ALL CHECKPOINTS
228         )
229     ),
230     param_space=config,
231 )
232
233 results = tuner.fit()
234
235 #-----
236 #Hyperparameter Tuning and Testing for Syntetic Dataset
237 #-----
238 scheduler2 = ASHAScheduler(
239     max_t=max_num_epochs,
240     grace_period=3,
241     reduction_factor=2,)

```

```

242
243 tuner2 = tune.Tuner(
244     tune.with_resources(
245         tune.with_parameters(training, train=synTrainingPath, val=
synValidationPath),
246         resources={"cpu": num_cores, "gpu": gpus_per_trial}
247     ),
248     tune_config=tune.TuneConfig(
249         metric="v_loss",
250         mode="min",
251         scheduler=scheduler2,
252         num_samples=num_samples,
253     ),
254     run_config=RunConfig(
255         checkpoint_config=CheckpointConfig(
256             num_to_keep = None #SAVE ALL CHECKPOINTS
257         )
258     ),
259     param_space=config,
260 )
261
262 results2 = tuner2.fit()

```

## C.5 CNN Transfer Learning-based AlexNet Testing

```

1 from functools import partial
2 import os
3 import torch
4 import torch.nn as nn
5 import torch.nn.functional as F
6 import torch.optim as optim
7 import torchvision
8 import torchvision.transforms as transforms
9 from torchvision.datasets import ImageFolder
10 from torch.utils.data import DataLoader
11 import ray
12
13 from ray import tune
14 from ray.air import session, RunConfig, CheckpointConfig, Checkpoint
15 from ray.tune.schedulers import ASHAScheduler
16
17 device = "cuda" if torch.cuda.is_available() else "cpu"

```

```

18
19 num_cores = 64
20 gpus_per_trial = 2
21 max_num_epochs = 64
22 num_samples = 64 #Number of items in hyperparameter search
23
24 trainingPath = "/home/ricardo.velasco/TT_kichwa/
    transfer_learning_model/spec/train"
25 validationPath = "/home/ricardo.velasco/TT_kichwa/
    transfer_learning_model/spec/val"
26 testingPath = "/home/ricardo.velasco/TT_kichwa/transfer_learning_model
    /spec/test"
27
28 synTrainingPath = "/home/ricardo.velasco/TT_kichwa/
    transfer_learning_model/syn_spec/train"
29 synValidationPath = "/home/ricardo.velasco/TT_kichwa/
    transfer_learning_model/syn_spec/val"
30 synTestingPath = "/home/ricardo.velasco/TT_kichwa/
    transfer_learning_model/syn_spec/test"
31
32 #Load training, testing and validation function for hyperparameter
    searching
33 def load_data(train, val):
34
35     #Input resmampling
36     transform = transforms.Compose([
37         transforms.Resize(size=(224, 224)),
38         transforms.ToTensor(),
39         transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229,
0.224, 0.225]),
40     ])
41
42     trainset = ImageFolder(train, transform=transform)
43
44     valset = ImageFolder(val, transform=transform)
45
46     return trainset, valset
47
48
49 def metrics(outputs, labels):
50     TP = ((outputs == 1) & (labels == 1)).sum().item()
51     FP = ((outputs == 1) & (labels == 0)).sum().item()
52     TN = ((outputs == 0) & (labels == 0)).sum().item()
53     FN = ((outputs == 0) & (labels == 1)).sum().item()

```

```

54
55     ppv = (TP / (TP + FP)) if (TP + FP) > 0 else 0
56     npv = (TN / (TN + FN)) if (TN + FN) > 0 else 0
57
58     return ppv, npv
59
60 def training(config, train=None, val=None):
61     #Import Alexnet pretrained model
62     net = torchvision.models.alexnet(weights='AlexNet_Weights.
IMAGENET1K_V1')
63
64     #Freeze features' layer for transfer learning
65     for param in net.features.parameters():
66         param.requires_grad = False
67
68     #Reset and modify classification layer suitable for binary
classification
69     net.classifier = torch.nn.Sequential(
70         torch.nn.Dropout(p=0.5, inplace=False),
71         torch.nn.Linear(in_features=9216, out_features=config["hn1"],
bias=True),
72         torch.nn.ReLU(inplace=True),
73         torch.nn.Dropout(p=0.5, inplace=False),
74         torch.nn.Linear(in_features=config["hn1"], out_features=config
["hn2"], bias=True),
75         torch.nn.ReLU(inplace=True),
76         torch.nn.Linear(in_features=config["hn2"], out_features=2,
bias=True))
77
78     #Send the model to GPU and put it in training mode
79     net.to(device)
80     net.train()
81
82     #Run always on GPU, if there are more GPUs, allow parallel data
parallel training
83     if torch.cuda.device_count() > 1:
84         net = nn.DataParallel(net)
85
86     #Define Loss Function
87     criterion = nn.CrossEntropyLoss()
88
89     #Define optimizer and establish a trainable parameter for
learning rate
90     optimizer = optim.SGD(net.classifier.parameters(), lr=config["lr"]

```

```

], momentum=0.9)
91
92 # To restore a checkpoint, use 'session.get_checkpoint()'.
93 checkpoint = session.get_checkpoint()
94
95 if checkpoint:
96     checkpoint_state = checkpoint.to_dict()
97     start_epoch = checkpoint_state["epoch"]
98     net.classifier.load_state_dict(checkpoint_state["
classifier_state_dict"])
99     optimizer.load_state_dict(checkpoint_state["
optimizer_state_dict"])
100 else:
101     start_epoch = 0
102
103
104 #Import training, validation and testing datasets
105 trainset, valset= load_data(train, val)
106
107 #Load training dataset
108 trainloader = DataLoader(
109     dataset=trainset,
110     batch_size=int(config["batch_size"]),
111     num_workers=num_cores,
112     shuffle = True
113 )
114
115 #Load validation dataset
116 valloader = DataLoader(
117     dataset=valset,
118     batch_size=len(valset), #Due to hardware, pick largest batch
size possible to compute metrics
119     num_workers=num_cores,
120     shuffle = True
121 )
122
123
124 for epoch in range(start_epoch, max_num_epochs):
125     # -----
126     # TRAINING STEP
127     # -----
128     train_loss, train_ppv, train_npv = 0.0, 0.0, 0.0
129     for i, (inputs, labels) in enumerate(trainloader):
130         # get the inputs; data is a list of [inputs, labels]

```



```

131         inputs, labels = inputs.to(device), labels.to(device)
132
133         # zero the parameter gradients
134         optimizer.zero_grad()
135
136         # forward + backward + optimize
137         outputs = net(inputs)
138         loss = criterion(outputs, labels)
139         loss.backward()
140         optimizer.step()
141
142         #Accumulate statistics across a single batch
143         _, y = torch.max(outputs.data, 1)
144         train_loss += loss.item()
145         train_ppv += metrics(y, labels)[0]
146         train_npv += metrics(y, labels)[1]
147
148         #Compute average of metrics across all batches
149         train_loss /= len(trainloader)
150         train_ppv /= len(trainloader)
151         train_npv /= len(trainloader)
152
153         #-----
154         # VALIDATION STEP
155         #-----
156         val_loss, val_ppv, val_npv = 0.0, 0.0, 0.0
157
158         for i, (inputs, labels) in enumerate(valloader):
159             with torch.no_grad():
160                 inputs, labels = inputs.to(device), labels.to(device)
161
162                 #Compute outputs
163                 outputs = net(inputs)
164                 _, y = torch.max(outputs.data, 1)
165
166                 #Compute loss
167                 loss = criterion(outputs, labels)
168
169                 #Accumulate statistics across a single batch
170                 val_loss+= loss.item()
171                 val_ppv += metrics(y, labels)[0]
172                 val_npv += metrics(y, labels)[1]
173
174         #Compute average of metrics across all batches

```

```

175     val_loss /= len(valloader)
176     val_ppv /= len(valloader)
177     val_npv /=len(valloader)
178
179     # Here we save a checkpoint. It is automatically registered
with
180     # Ray Tune and can be accessed through 'session.get_checkpoint
()'
181
182     # API in future iterations.
183     #Checkpoint data
184     checkpoint_data = {
185         "epoch": epoch,
186         "classifier_state_dict": net.classifier.state_dict(),
187         "optimizer_state_dict": optimizer.state_dict()
188     }
189
190     checkpoint = Checkpoint.from_dict(checkpoint_data)
191     session.report({"v_loss": val_loss, "v_ppv": val_ppv,
192                  "v_npv": val_npv, "tr_loss":train_loss,
193                  "tr_ppv":train_ppv, "tr_npv":train_npv},
checkpoint=checkpoint)
194
195
196 def testing(path, hn1, hn2, c_path):
197     #Resample function
198     transform = transforms.Compose([
199         transforms.Resize(size=(224, 224)),
200         transforms.ToTensor(),
201         transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229,
0.224, 0.225]),
202     ])
203
204     #Define testing dataset from folder
205     test = ImageFolder(path, transform=transform)
206
207     #Load testing dataset
208     testloader = torch.utils.data.DataLoader(
209         test, batch_size=len(test), shuffle=True, num_workers=
num_cores
210     )
211
212     #Load pretrained AlexNet model
213     net = torchvision.models.alexnet(weights='AlexNet_Weights.

```

```

IMAGENET1K_V1')
214
215 #Modify classification layer with best parameters
216 net.classifier = torch.nn.Sequential(
217     torch.nn.Dropout(p=0.5, inplace=False),
218     torch.nn.Linear(in_features=9216, out_features=hn1, bias=True)
219 ,
220     torch.nn.ReLU(inplace=True),
221     torch.nn.Dropout(p=0.5, inplace=False),
222     torch.nn.Linear(in_features=hn1, out_features=hn2, bias=True),
223     torch.nn.ReLU(inplace=True),
224     torch.nn.Linear(in_features=hn2, out_features=2, bias=True))
225
226 net.to(device)
227 net.eval()
228
229 #Load checkpoint and weights
230 checkpoint = Checkpoint.from_directory(c_path)
231 checkpoint_state = checkpoint.to_dict()
232 net.classifier.load_state_dict(checkpoint_state["
classifier_state_dict"])
233
234 test_ppv, test_npv = 0, 0
235
236 with torch.no_grad():
237     for data in testloader:
238         images, labels = data
239         images, labels = images.to(device), labels.to(device)
240
241         outputs = net(images)
242         _, y = torch.max(outputs.data, 1)
243
244         test_ppv += metrics(y, labels)[0]
245         test_npv += metrics(y, labels)[1]
246
247     test_ppv /= len(testloader)
248     test_npv /= len(testloader)
249
250 return test_ppv, test_npv, (test_ppv + test_npv)/2
251
252 #Save table with best results from each trial given a metric and
253 filter_mode
254 def saveFile(results, metric, filter_mode, save):
255     df = results.get_dataframe(filter_metric = metric, filter_mode =

```

```

    filter_mode)
254     df.to_csv(save)
255
256
257 #Restore tuners given training folders corresponding to original
    dataset and synthetic dataset
258 tuner = tune.Tuner.restore("/home/ricardo.velasco/ray_results/
    training_2023-09-18_16-14-45", trainable=training)
259 syn_tuner = tune.Tuner.restore("/home/ricardo.velasco/ray_results/
    training_2023-09-18_19-57-26", trainable=training)
260
261 #Get results from each tuner
262 results = tuner.get_results()
263 syn_results = syn_tuner.get_results()
264
265 #Save results given metrics, filter_mode and save path for dataset
266 saveFile(results, "v_loss", "min", "/home/ricardo.velasco/TT_kichwa/
    transfer_learning_model/tuning/tables/alexnet_ds_loss.csv")
267 saveFile(results, "v_bpv", "max", "/home/ricardo.velasco/TT_kichwa/
    transfer_learning_model/tuning/tables/alexnet_ds_bpv.csv")
268
269 #Save results given metrics, filter_mode and save path for synthetic
    dataset
270 saveFile(syn_results, "v_loss", "min", "/home/ricardo.velasco/
    TT_kichwa/transfer_learning_model/tuning/tables/alexnet_syn_ds_loss
    .csv")
271 saveFile(syn_results, "v_bpv", "max", "/home/ricardo.velasco/TT_kichwa
    /transfer_learning_model/tuning/tables/alexnet_syn_ds_bpv.csv")
272
273 #Test data given checkpoints
274 #Original dataset best TOP 5 MODELS LOSS
275 a, b, c = testing(testingPath, 4096, 512, "/home/ricardo.velasco/
    ray_results/training_2023-09-18_16-14-45/
    training_80559_00021_21_batch_size=64,hn1=4096,hn2=512,lr=0.0021
    _2023-09-18_16-14-53/checkpoint_000032/")
276 print(a, b, c)
277
278 a, b, c = testing(testingPath, 2048, 256, "/home/ricardo.velasco/
    ray_results/training_2023-09-18_16-14-45/
    training_80559_00011_11_batch_size=64,hn1=2048,hn2=256,lr=0.0004
    _2023-09-18_16-14-53/checkpoint_000021")
279 print(a, b, c)
280
281 a, b, c = testing(testingPath, 512, 2048, "/home/ricardo.velasco/

```

```

ray_results/training_2023-09-18_16-14-45/
training_80559_00009_9_batch_size=64,hn1=512,hn2=2048,lr=0.0006
_2023-09-18_16-14-53/checkpoint_000010")
282 print(a, b, c)
283
284 a, b, c = testing(testingPath, 512, 1024, "/home/ricardo.velasco/
ray_results/training_2023-09-18_16-14-45/
training_80559_00000_0_batch_size=64,hn1=512,hn2=1024,lr=0.0052
_2023-09-18_16-14-53/checkpoint_000042")
285 print(a, b, c)
286
287 a, b, c = testing(testingPath, 1024, 2048, "/home/ricardo.velasco/
ray_results/training_2023-09-18_16-14-45/
training_80559_00006_6_batch_size=512,hn1=1024,hn2=2048,lr=0.0635
_2023-09-18_16-14-53/checkpoint_000046")
288 print(a, b, c)
289
290
291 #Original dataset best TOP 5 MODELS BPV
292 a, b, c = testing(testingPath, 4096, 2048, "/home/ricardo.velasco/
ray_results/training_2023-09-18_16-14-45/
training_80559_00016_16_batch_size=128,hn1=4096,hn2=2048,lr=0.0438
_2023-09-18_16-14-53/checkpoint_000000")
293 print(a, b, c)
294
295 a, b, c = testing(testingPath, 64, 64, "/home/ricardo.velasco/
ray_results/training_2023-09-18_16-14-45/
training_80559_00028_28_batch_size=128,hn1=64,hn2=64,lr=0.0015_2023
-09-18_16-14-53/checkpoint_000007")
296 print(a, b, c)
297
298 a, b, c = testing(testingPath, 1024, 2048, "/home/ricardo.velasco/
ray_results/training_2023-09-18_16-14-45/
training_80559_00006_6_batch_size=512,hn1=1024,hn2=2048,lr=0.0635
_2023-09-18_16-14-53/checkpoint_000009")
299 print(a, b, c)
300
301 a, b, c = testing(testingPath, 512, 1024, "/home/ricardo.velasco/
ray_results/training_2023-09-18_16-14-45/
training_80559_00000_0_batch_size=64,hn1=512,hn2=1024,lr=0.0052
_2023-09-18_16-14-53/checkpoint_000008")
302 print(a, b, c)
303
304 a, b, c = testing(testingPath, 2048, 256, "/home/ricardo.velasco/

```

```

ray_results/training_2023-09-18_16-14-45/
training_80559_00011_11_batch_size=64,hn1=2048,hn2=256,lr=0.0004
_2023-09-18_16-14-53/checkpoint_000005")
305 print(a, b, c)
306
307
308 #Synthetic dataset best TOP 5 LOSS TESTING
309 a, b, c = testing(synTestingPath, 256, 2048, "/home/ricardo.velasco/
ray_results/training_2023-09-18_19-57-26/
training_9792e_00012_12_batch_size=128,hn1=256,hn2=2048,lr=0.0034
_2023-09-18_19-57-27/checkpoint_000012")
310 print(a, b, c)
311
312 a, b, c = testing(synTestingPath, 128, 128, "/home/ricardo.velasco/
ray_results/training_2023-09-18_19-57-26/
training_9792e_00000_0_batch_size=64,hn1=128,hn2=128,lr=0.0016_2023
-09-18_19-57-27/checkpoint_000044")
313 print(a, b, c)
314
315 a, b, c = testing(synTestingPath, 2048, 128, "/home/ricardo.velasco/
ray_results/training_2023-09-18_19-57-26/
training_9792e_00013_13_batch_size=64,hn1=2048,hn2=128,lr=0.0102
_2023-09-18_19-57-27/checkpoint_000012")
316 print(a, b, c)
317
318 a, b, c = testing(synTestingPath, 4096, 64, "/home/ricardo.velasco/
ray_results/training_2023-09-18_19-57-26/
training_9792e_00008_8_batch_size=512,hn1=4096,hn2=64,lr=0.0114
_2023-09-18_19-57-27/checkpoint_000002")
319 print(a, b, c)
320
321 a, b, c = testing(synTestingPath, 1024, 64, "/home/ricardo.velasco/
ray_results/training_2023-09-18_19-57-26/
training_9792e_00007_7_batch_size=128,hn1=1024,hn2=64,lr=0.0354
_2023-09-18_19-57-27/checkpoint_000003")
322 print(a, b, c)
323
324 #Synthetic dataset TOP 5 BEST BPV TESTING
325 a, b, c = testing(synTestingPath, 2048, 128, "/home/ricardo.velasco/
ray_results/training_2023-09-18_19-57-26/
training_9792e_00013_13_batch_size=64,hn1=2048,hn2=128,lr=0.0102
_2023-09-18_19-57-27/checkpoint_000011")
326 print(a, b, c)
327

```

```
328 a, b, c = testing(synTestingPath, 1024, 64, "/home/ricardo.velasco/  
ray_results/training_2023-09-18_19-57-26/  
training_9792e_00007_7_batch_size=128,hn1=1024,hn2=64,lr=0.0354  
_2023-09-18_19-57-27/checkpoint_000003")  
329 print(a, b, c)  
330  
331 a, b, c = testing(synTestingPath, 1024, 64, "/home/ricardo.velasco/  
ray_results/training_2023-09-18_19-57-26/  
training_9792e_00014_14_batch_size=256,hn1=1024,hn2=64,lr=0.0265  
_2023-09-18_19-57-27/checkpoint_000003")  
332 print(a, b, c)  
333  
334 a, b, c = testing(synTestingPath, 4096, 64, "/home/ricardo.velasco/  
ray_results/training_2023-09-18_19-57-26/  
training_9792e_00008_8_batch_size=512,hn1=4096,hn2=64,lr=0.0114  
_2023-09-18_19-57-27/checkpoint_000002")  
335 print(a, b, c)  
336  
337 a, b, c = testing(synTestingPath, 4096, 128, "/home/ricardo.velasco/  
ray_results/training_2023-09-18_19-57-26/  
training_9792e_00024_24_batch_size=64,hn1=4096,hn2=128,lr=0.0008  
_2023-09-18_19-57-27/checkpoint_000005")  
338 print(a, b, c)
```