

Media Engineering and Technology Faculty
German University in Cairo



Voice Spoofing Detection System

Bachelor Thesis

Author: Eslam Gamal ismail

Supervisors: Dr.Tamer Mostafa

Submission Date: 19 May 2024

This is to certify that:

- (i) the thesis comprises only my original work toward the Bachelor Degree
- (ii) due acknowledgement has been made in the text to all other material used

Eslam Gamal ismail
19 May 2024

Acknowledgments

I would like to thank God and my family for their unwavering support. I am also deeply grateful to Dr. Tamer Mostafa for his invaluable assistance and guidance, as well as for the weekly meetings that greatly contributed to the completion of this thesis. Additionally, I extend my thanks to the university for providing the resources and environment necessary for my research. .

Abstract

Voice spoofing presents a serious security risk to authentication systems. This issue is extensively discussed in this thesis. This work is designed to identify real and spoofed audio using convolutional neural networks, a machine learning model. The research uses the ASVSpooF 2019 dataset to train, test, and evaluate the CNN model in order to develop a robust model using conventional techniques to identify spoofing attacks such as replay, voice conversion, and synthesis from the bona fide voice. We use feature extraction methods applied to the audio, such as the extraction of spectrogram analysis, to detect spoofing efforts. Dynamic training and transfer learning strategies are implemented to improve the model's performance in combating spoofing attacks. The model is tested and found to be effective in many scenarios, demonstrating its potential in addressing voice spoofing attacks, which increases the security of voice-driven systems. The results of this thesis contribute to the enhancement of digital security and establish the foundation for future initiatives aimed at protecting voice authentication systems from increasingly sophisticated threats.

Contents

Acknowledgments	V
1 Introduction	1
1.1 Motivation and Objectives	1
1.2 Outline	2
2 Background	3
2.1 Voice spoofing	3
2.1.1 Types of attacks	3
2.1.2 Facing it	3
2.2 Feature Extraction Techniques for Audio Spoof Detection	4
2.2.1 Mel Frequency Cepstral Coefficients (MFCCs) [2]	4
2.2.2 Tonnetz	5
2.2.3 Other Spectral Features	6
2.2.4 Spectrograms	8
2.3 Machine learning	9
2.4 supervised Machine learning	9
2.5 Deep learning	10
2.5.1 Layers:	10
2.5.2 Activation Functions	11
2.5.3 Weights and Biases	11
2.5.4 Loss Function	11
2.5.5 Optimization Algorithm	11
2.5.6 Backpropagation	12
2.5.7 Epochs:	12
2.5.8 Learning Rate	12
2.6 CNN model	13
2.6.1 Feature Detection:	13
2.6.2 Hierarchical Processing:	13
2.6.3 Pattern Recognition:	13
2.6.4 Adaptation and Learning:	14

3	literature review	19
3.0.1	Exploring Attack Techniques: Understanding Spoofing Methods	19
3.0.2	Harnessing Machine Learning: Advancements in Detection	19
3.0.3	Emergence of Deep Learning: Leveraging Neural Networks	20
3.0.4	Real-time Detection and Comprehensive Solutions	21
3.0.5	Cutting-edge Research and Novel Methodologies	21
3.0.6	Emerging Trends and Technologies	22
3.0.7	Synthesis and Future Directions	22
4	Methodology	23
4.1	Introduction	23
4.2	Data Collection	23
4.2.1	Protocol Files	23
4.2.2	Data Loading and Processing	23
4.2.3	Data Sampling	24
4.3	Data Processing and Splitting	24
4.3.1	Audio to Spectrogram Conversion	25
4.3.2	Data Splitting	26
4.4	Model Architecture	27
4.5	Neural Network Architecture	28
4.5.1	Input Layer	28
4.5.2	Convolutional Layers	28
4.5.3	Pooling Layers	28
4.5.4	Dropout Layers	28
4.5.5	Flatten and Dense Layers	28
4.5.6	Output Layer	29
4.6	Training Process	29
4.6.1	Loss Function	29
4.6.2	Optimizer	29
4.6.3	Metrics	29
4.6.4	Callbacks	29
4.7	Model Evaluation	30
5	Results & Limitations	31
5.1	Experiments Setup	31
5.1.1	Data Preparation	31
5.1.2	Model Configuration	32
5.1.3	Training	32
5.1.4	Evaluation	32
5.2	Dataset Description	33
5.2.1	Composition	33
5.2.2	Audio Specifications	33
5.2.3	Applications	33
5.2.4	Impact	33

5.3	Experiment 1: Basic Feature Analysis	34
5.4	Experiment 2: Advanced Audio Feature Analysis	34
5.5	Results Analysis and Discussion	34
5.6	Conclusion	38
6	Conclusion & Future Work	39
6.1	Summary of Findings	39
6.2	Implications of the Research	39
6.3	Future Work	40
6.4	Concluding Remarks	40
	Appendix	41
	A List of Abbreviations	42
	List of Abbreviations	42
	List of Figures	43
	bibliography	48

Chapter 1

Introduction

1.1 Motivation and Objectives

In our increasingly digital environment, voice spoofing—the practice of employing technology to produce spoofed voices that sound real—voice a serious risk to security systems and individual privacy. Voice authentication methods have grown as technology continues to enter more aspects of our daily life [12]. voice recognition systems depend on an individual's unique voice to allow access to sensitive information or control personal devices. Examples of these systems include those found in smart speakers and virtual assistants like Siri. However, voice spoofing methods allow attackers to take benefits from these weak systems and and use these persons voice in order to access to unauthorized places [32].

The absence of sufficient data to train models in order to detect attacks is one of the main obstacles to face voice spoofing. not as the others cyper attacks , such malware or phishing assaults [1], voice spoofing attacks are not as well-known to the public. Furthermore, voice spoofing methods are are more danger and improved as time pass , making it hard for detection systems to stay up to date. Moreover, processing and even storing private conversations is a privacy risk when using voice transcript analysis to identify spoofing[41].

Attacks using voice spoofing methods include speech synthesis, voice conversion, and voice replay. Although voice synthesis creates fully artificial voices that seems to be human voice [58], voice conversion means change the voice to be look like someone else[50]. In order to pass voice recognition systems, voice replay attacks is the recording of the voice and replay it .

Researchers are using deep learning techniques and artificial intelligence (AI) to face the threat of the voice spoofing. AI models are trained on datasets that include both spoof and bonefide voices [7].

Furthermore, real-time classification methods have improved the voice-based authentication systems by enabling the identify of the spoof voices . Users can simply upload

audio samples to check if it is spoof or not spoof and get feedback on the voice's by deploying optimized models to web services. This strategy is an important step in order to face the voice spoofing attacks to secure the users against the voice attacks [8].

Furthermore, security and user privacy must be balanced while facing voice spoofing attacks. Although strong authentication procedures are necessary to secure private data, too much security measures might bother users and make voice-based systems to be less used. Thus, it is important to create detection methods that are friendly to the and effective in order to increase the rate of the detection of spoofed voices while decreasing false positives[9].

To sum up, voice spoofing is a serious risk to privacy and security in the growth of society. But thanks to improvement in deep learning and artificial intelligence, scientists are getting a lot better at creating reliable detection systems that can recognize better and decrease speech spoofing attempts. These systems present a viable way to protect voice-based authentication systems and users from attackers by utilizing datasets, sophisticated algorithms, and real-time classification techniques. To remain abreast of developing threats and make sure the growth of security of voiced technology, however, continued research and contributions are essential. The following procedures can be used to create a multi-class deep learning model for CNN-based voice spoofing detection:

1.2 Outline

This thesis is organized into five chapters, each focusing on different aspects of voice spoofing detection. Chapter 2 sets the scene by explaining the basic concepts and background behind voice spoofing. Chapter 3 reviews previous studies, showing what researchers have already discovered and pointing out what still needs to be explored. Chapter 4 describes the methods used in this study, including how data was prepared, what features were important, and the kinds of machine learning techniques that were used. Chapter 5 talks about the experiments carried out and the results obtained, discussing how well the methods worked in identifying fake voices.

Chapter 2

Background

2.1 Voice spoofing

The practice of sounding like someone else in order to fool voice recognition software is known as voice spoofing. The development of voice-activated technologies has drawn attention to this type of security issue. Although the exact origins of voice spoofing as a known danger are difficult to determine, worries have increased dramatically since voice-activated systems and gadgets became widely used in the 2010s. Technological developments have made it simpler to produce convincingly phony audio, which raises the possibility of fraud and illegal access. There is more than one type of voice spoof attacks.

2.1.1 Types of attacks

Playback attacks: Using the target's voice recorded on tape. Synthetic voice generation is the process of using text-to-speech (TTS) technology to produce synthetic voices. Voice Conversion: Adapting an already-existing voice recording to the vocal traits of the intended audience[13], and this attacks need methods to detect them.

2.1.2 Facing it

The processing capacity at hand previously restricted the application of machine learning (ML) and deep learning (DL) approaches to prevent voice spoofing. In order to identify differences in speech characteristics that might point to faked voices, early machine learning techniques relied on feature engineering and more basic algorithms. But these strategies had their limits, particularly when spoofing techniques advanced in sophistication[61].

An important development in voice authentication system security was the switch to deep learning. Deep learning algorithms have shown to be more successful in spotting

minute irregularities in audio signals linked to sophisticated spoofing attempts because of their capacity to process enormous volumes of data and discover intricate patterns. Developments in processing resources, such as strong CPUs, GPUs, and distributed computing frameworks, allowed for this change in strategy. The availability of supercomputers and high-performance computing clusters enabled researchers to train deep learning models on massive datasets, leading to improved accuracy and robustness in detecting voice spoofing attacks[62].

As a result, deep learning has become a cornerstone in the development of more secure voice authentication systems, capable of mitigating the evolving threat landscape posed by sophisticated voice spoofing techniques. This demonstrates how advancements in hardware infrastructure have empowered the adoption of advanced ML and DL techniques in addressing real-world security challenges.

2.2 Feature Extraction Techniques for Audio Spoof Detection

Feature extraction is critical in the study of audio signals, particularly in speech and audio spoof detection. These features improve the performance of machine learning models by converting raw audio input into a more informative and discriminative representation. This section addresses various significant audio features used in this work, such as Mel Frequency Cepstral Coefficients (MFCCs), Tonnetz, and others, and explains their importance and applicability in detecting faked audio.

2.2.1 Mel Frequency Cepstral Coefficients (MFCCs)[2]

MFCCs are one of the most widely used features in speech recognition and audio analysis. They are derived from a type of cepstral representation of the audio clip (a non-linear "spectrum-of-a-spectrum") which provides a representation of the short-term power spectrum of a sound.

The computation of MFCCs involves several steps:

1. **Frame the signal into short frames:** Audio transmissions are separated into short frames because the signal's frequency content changes over time.
2. **Apply the Fourier transform:** This transform converts the time domain signal into the frequency domain.
3. **Map the powers of the spectrum to the mel scale:** The mel scale is an auditory scale that more closely resembles the human ear's response than the linearly spaced frequency bands. the linearly-spaced frequency bands.

2.2. FEATURE EXTRACTION TECHNIQUES FOR AUDIO SPOOF DETECTION⁵

4. **Take the logarithm:** This step is taken to use the cepstral mean normalization process, which allows better representation.
5. **Apply the discrete cosine transform:** Apply the discrete cosine transform, which is identical to a Fourier transform but only employs cosine functions.

The generated features (MFCCs) capture the power spectrum's envelope and are thus extremely useful in recognizing phonemes (different units of sound) in spoken language, which is important for distinguishing authentic from faked audio.

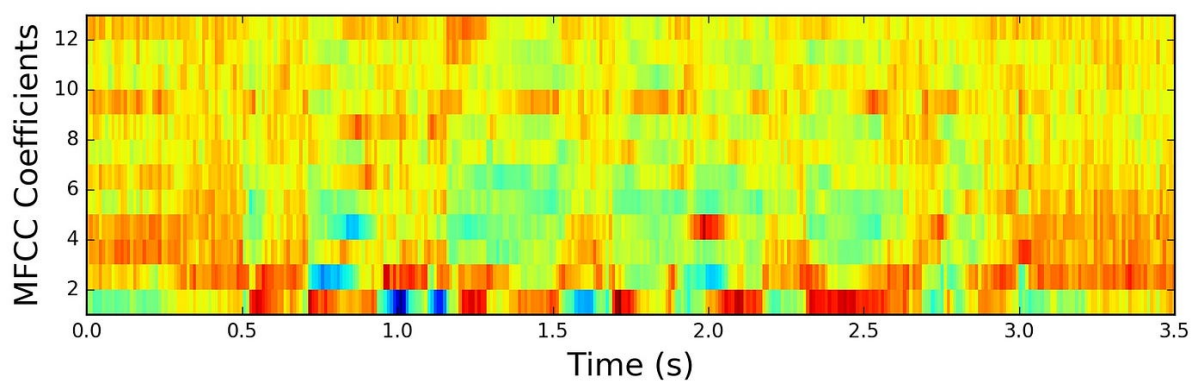


Figure 2.1: MFCCs

2.2.2 Tonnetz

The tonal centroid features, or tonnetz, are employed to represent the harmonic content of music. This representation is derived from the tonal pitch class profile, which captures harmonic relations in Western music theory[15].

Tonnetz features are calculated as follows:

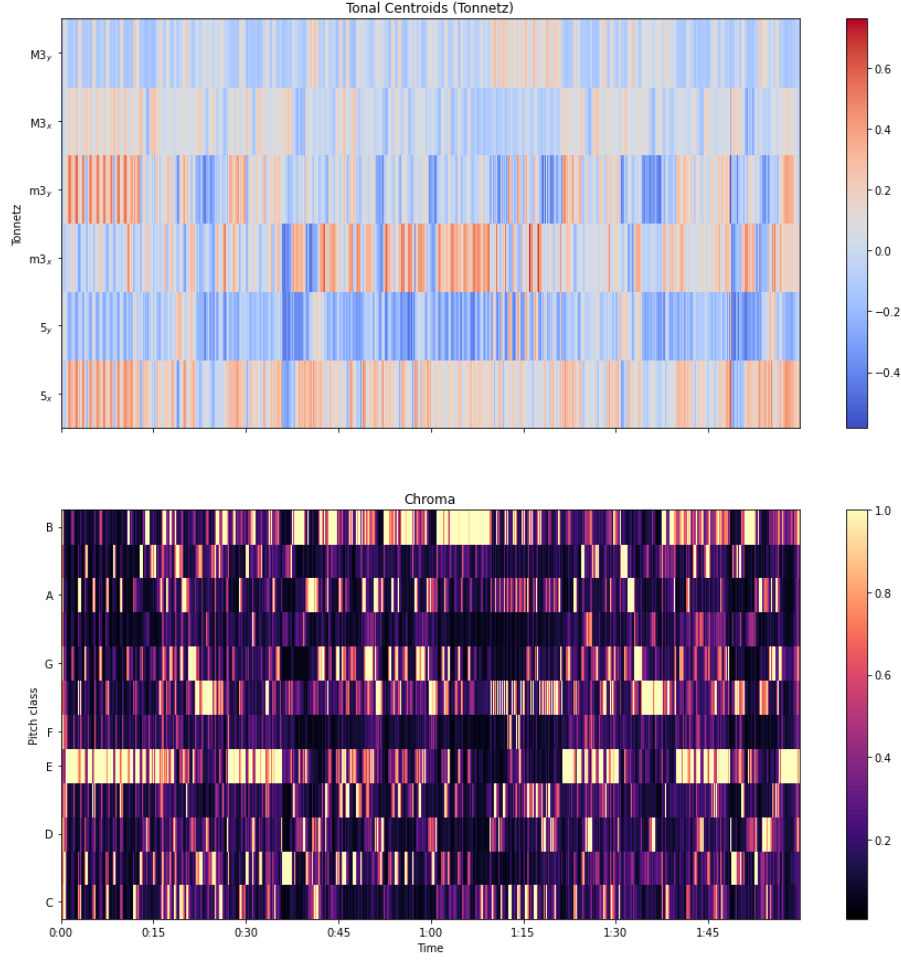


Figure 2.2: Tonnez

1. **Harmonic Pitch Class Profiles (HPCP):** These profiles are generated to depict the strengths of the twelve pitch classes (chroma) of the musical octave.
2. **Mapping to Tonal Centroid Space:** The six-dimensional tonnetz space represents the perfect fifth, minor third, and major third intervals.

These traits are very beneficial for music information retrieval tasks like as key detection, chord recognition, and emotion perception in music. Tonnez aids in the detection of audio spoofing by identifying abnormal shifts and anomalies in the harmonic structure that may suggest tampering or synthesis.

2.2.3 Other Spectral Features

This study also considered spectral contrast, spectral roll-off, and zero-crossing rate. Each of these characteristics captures distinct aspects of the audio signal.:

2.2. FEATURE EXTRACTION TECHNIQUES FOR AUDIO SPOOF DETECTION7

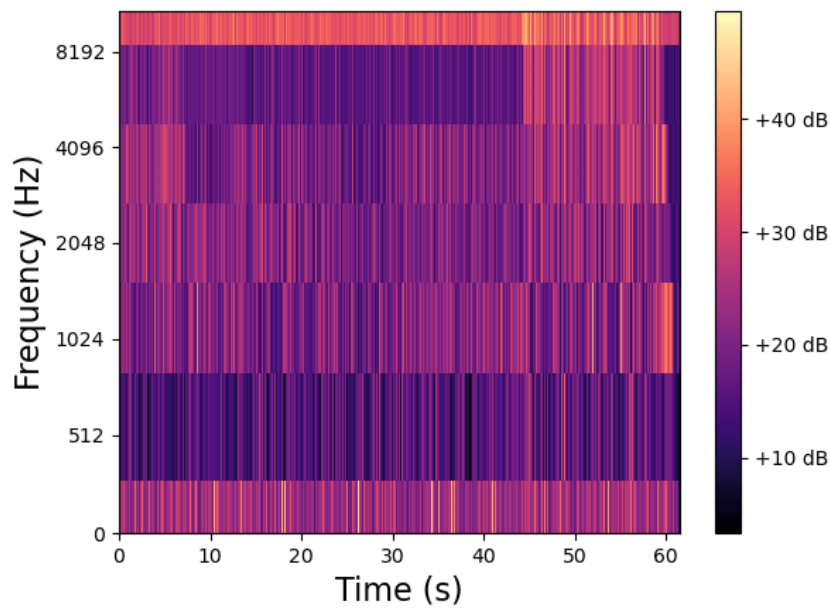


Figure 2.3: Spectral Contrast

- **Spectral Contrast:** Measures the difference between peaks and valleys in the sound spectrum, which is useful for discriminating between speech and noise[59].
- **Spectral Roll-off:** The spectral envelope's form is shown by the frequency below which a specific percentage of the total spectral energy, such as 85
- **Zero-Crossing Rate:** The rate at which the signal changes signs, providing insight into the frequency content of the signal[19].

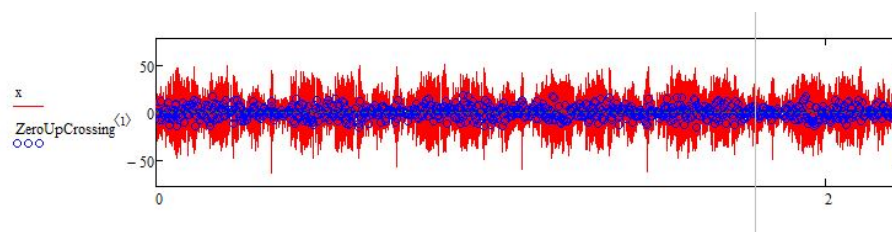


Figure 2.5: Zero-Crossing Rate

These attributes work together to provide a more robust analysis of audio signals, improving the model's capacity to recognize and distinguish between regular and faked audio samples.

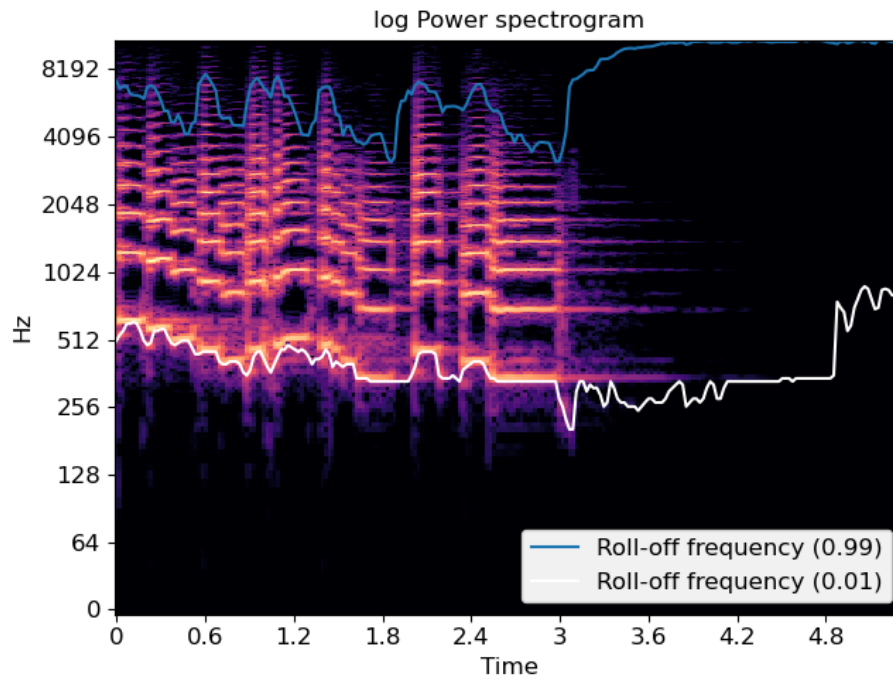


Figure 2.4: Spectral Roll-off[45]

2.2.4 Spectrograms

A spectrogram is a graphical representation of the frequencies that comprise a sound source as they change over time. It is commonly used to examine the frequency content of signals, particularly in acoustics, voice processing, and music technology[18].

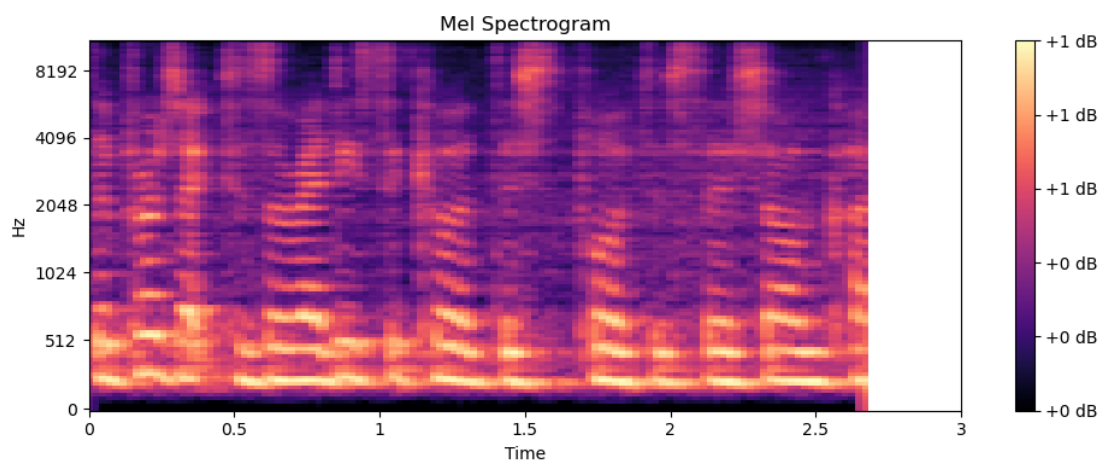


Figure 2.6: Spectrogram

Definition and Importance

A spectrogram plots the strength of frequencies in a signal across time, offering a visual representation of how these frequencies fluctuate. This makes spectrograms useful for a wide range of applications, including voice analysis, music analysis, and bioacoustics. In the case of speech spoofing detection, spectrograms are crucial for recognizing unique patterns that distinguish authentic from counterfeit sounds.

2.3 Machine learning

Building systems that can learn from and make judgments based on data is the focus of machine learning (ML), a subset of artificial intelligence (AI). With machine learning (ML), as opposed to traditional programming, which requires humans to explicitly write the behavior, computers may automatically learn from their experiences and get better over time. This is accomplished by creating algorithms that can analyze vast volumes of data, spot trends, and base judgments or predictions on those trends. Machine learning (ML) is widely applied in many domains, including speech and picture recognition[51], recommendation engines, and medical diagnosis.

2.4 supervised Machine learning

An method known as supervised machine learning (ML) uses labeled data to identify patterns and relationships. When an algorithm is trained on a dataset of input-output pairs—where the input is the data needed to generate predictions or judgments and the output is the desired prediction or decision—supervised learning takes place. The objective is to develop a mapping function from input to output so that the model can correctly predict or classify the output based on the input attributes when it is exposed to new, unknown data[30].

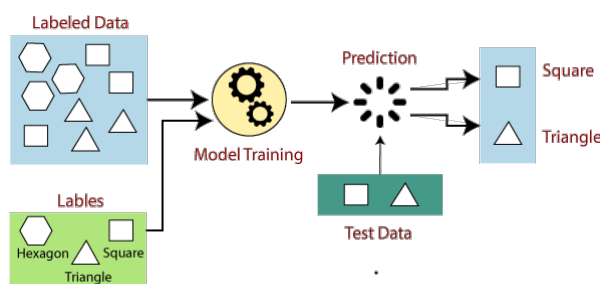


Figure 2.7: Supervised Machine learning

2.5 Deep learning

The term "deep" refers to a subset of machine learning that focuses on multiple-layered artificial neural networks. It uses layers of networked neurons to extract information and generate predictions, simulating how the human brain interprets and learns from data. Large volumes of data can be used to train deep learning models to recognize complicated patterns and representations. A number of essential elements make up deep learning (DL), which makes it possible to create intricate neural network models that can learn from vast volumes of data. The primary elements of deep learning are as follows[27]:

Neural Networks:Artificial neural networks, which are computer models based on the composition and operation of biological neurons in the human brain, are the fundamental building blocks of deep learning. Input, hidden, and output layers of neurons are coupled to form neural networks. Every neuron in the network modifies the incoming data and transmits the outcome to neurons in the subsequent layer[28]. Figure 2.2.

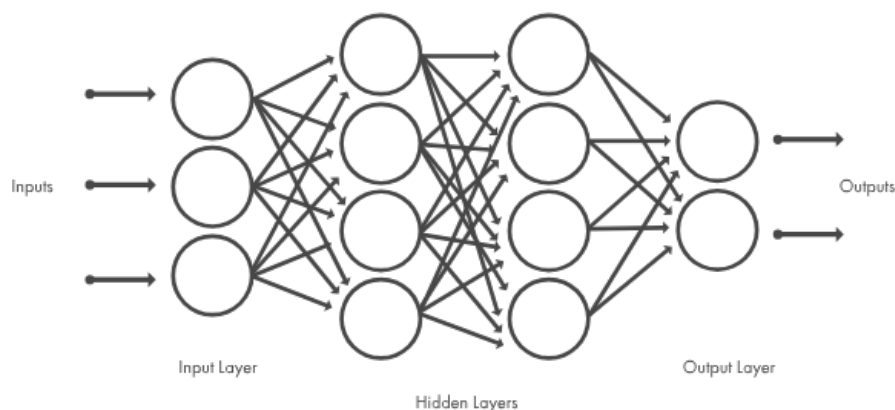


Figure 2.8: ANN

2.5.1 Layers:

Multiple layers of neurons are a common feature of deep learning models, which enable them to build hierarchical representations of the input data. Typical layer types consist of:

Input Layer Receives input data and passes it to the first hidden layer [29].

Hidden Layers Intermediate layers between the input and output layers. Each hidden layer progressively extracts abstract features from the input data [20].

Output Layer Produces the model's predictions or outputs based on the learned representations [22].

2.5.2 Activation Functions

The network gains non-linearity from activation functions, which allows it to learn intricate mappings between the input and output data. Rectified linear unit (ReLU), sigmoid, and tanh are examples of common activation functions.

ReLU (Rectified Linear Unit): Non-linear activation function for selective neuron activation, offering simplicity and high performance.[\[37\]](#)

Sigmoid: Used in binary-class classification, producing a probability value between 1 and 0, adding complexity to optimization[\[53\]](#).

Softmax: Primarily in the final layer for multi-class classification, conducting probabilistic evaluation[\[42\]](#).

Softsign: Resembles a sigmoid with a threshold, limiting values between -1 and 1.[\[42\]](#)

Tanh (Hyperbolic Tangent): Similar to sigmoid but outputs in the range of -1 to 1, accelerating optimization. [\[38\]](#)

2.5.3 Weights and Biases

:Each connection between neurons in adjacent layers is assigned a weight that defines its strength. Furthermore, each neuron contains a bias term that enables it to learn an offset from the input data. These weights and biases are changed during training to reduce the discrepancy between the expected and actual outputs[\[23\]](#).

2.5.4 Loss Function

:The loss function calculates the difference between the expected and actual outputs of the model. During training, the goal is to reduce this loss function by modifying the model's parameters (weights and bias)[\[43\]](#).

2.5.5 Optimization Algorithm

Deep learning models are taught with optimization algorithms like stochastic gradient descent (SGD) and Adam. These algorithms iteratively update the model's parameters using the gradient of the loss function with respect to each parameter[\[17\]](#).

- **Adam:** A stochastic Replace gradient descent with momentum, which is effective for noisy or sparse gradients, huge datasets, and parameters.[4]
- **Nadam:** An extension of Adam with Nesterov momentum, computing the decaying moving average of gradients.[39]
- **Adamax:** A variation of Adam introducing an infinite norm extension for potentially improved optimization[10]

2.5.6 Backpropagation

Backpropagation is an important algorithm for training deep learning models. It computes the gradient of the loss function in relation to each parameter in the network and updates the parameters accordingly, allowing the model to learn from training data[31].

Binary Cross-Entropy: Assesses each predicted probability against the actual class result, which can be either true (1) or false (0)[55].

Categorical Cross-Entropy: Assumes that only one class among all possible outcomes.[36]

2.5.7 Epochs:

When training a machine learning model, choosing the appropriate number of epochs is essential. A model that underfits will not be able to adequately reflect the underlying patterns in the data if the number of epochs is too low. However, the model runs the danger of overfitting if it is trained for an excessive number of epochs. Because the model has learnt the details and noise in the training data to a level that adversely affects its performance on general data, it may perform well on the training data but poorly on fresh, unknown data. To ensure that the model can generalize adequately without conforming too closely to the training dataset or staying too generic, the optimal number of epochs must be chosen to balance these concerns[48].

2.5.8 Learning Rate

One important hyperparameter that controls how fast a model modifies its parameters during training is the learning rate. It affects how much the training process, such gradient descent, updates the parameters with each iteration. A high learning rate can speed up convergence, but it can also cause unstable training dynamics or cause the loss function to exceed its minimum. On the other hand, a low learning rate might considerably impede

the training process but guarantees more steady and progressive updates, increasing the possibility of convergence.

Learning rate scheduling is a strategy that is frequently used to optimize the learning process. With this strategy, learning rate is increased initially and progressively decreased, enabling larger steps for quicker convergence and smaller steps for fine-tuning later on. Moreover, adaptive learning rate techniques such as Adam dynamically vary the learning rate according to the gradients of the loss function, improving the training efficiency by changing the learning rate for each parameter separately according to the behavior of the training process that is observed[57].

2.6 CNN model

A Convolutional Neural Network (CNN) is a type of model in computer science that draws inspiration from how the human brain processes visual information. It is engineered to recognize and interpret images by mimicking the way our brains identify shapes and objects through breaking them down into smaller components.

2.6.1 Feature Detection:

Similar to the way our brains pick up on subtle details such as edges, corners, and textures within a visual scene, CNNs employ convolutional layers to identify these same features. These layers meticulously analyze small segments of the image sequentially, akin to how our eyes scan and focus on various sections of what we see[54].

2.6.2 Hierarchical Processing:

Our visual systems begin by processing simple forms before progressing to more complex items. CNNs work similarly, utilizing many layers to gradually comprehend increasingly intricate features. Similar to how our brain assembles simple shapes to identify more complex objects, each layer builds on the information from the one before it to identify more complex portions of the image[33].

2.6.3 Pattern Recognition:

Our brains are highly adept at identifying patterns, like faces or trees, by combining many visual cues. By combining details observed in earlier phases to comprehend the overall picture, CNNs are trained to identify patterns[47].

2.6.4 Adaptation and Learning:

Our minds are constantly picking up new skills and adapting to new visual stimuli. In a similar vein, CNNs are trained using a large number of images to identify patterns and objects using a technique called backpropagation. As a result, they gradually improve their ability to recognize items, much like humans do when we see more of them.

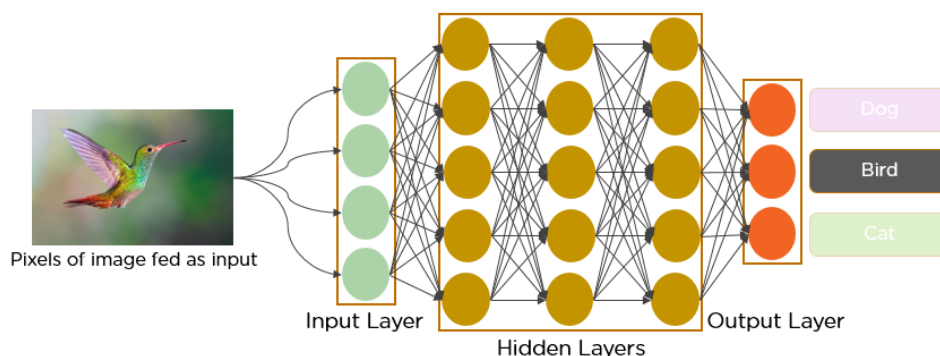


Figure 2.9: CNN model

Input Layer:

This layer receives the raw input data, which in the case of image processing, is usually the pixel values of the image[16].

Convolutional Layer:

the fundamental component of a CNN. To produce a feature map, it processes the input image via a series of filters, also known as kernels. As seen in, each filter is made to identify a particular feature (such as an edge, curve, etc.) in the input image[26]. Figure 2.3.

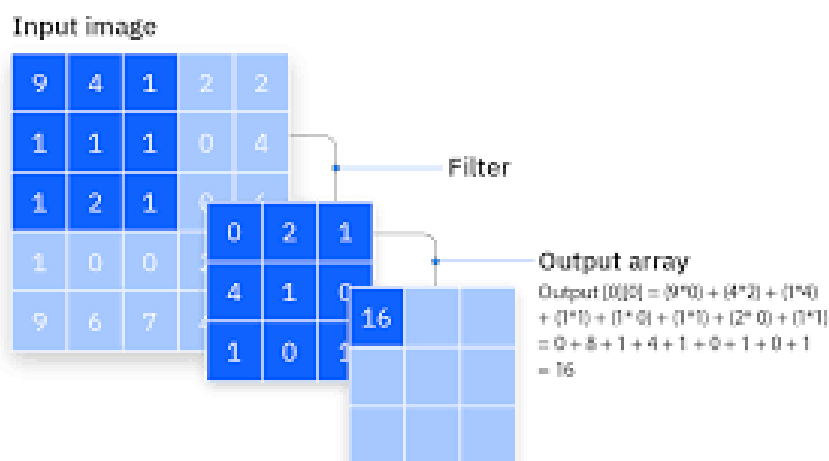


Figure 2.10: Convolution layer

ReLU Layer (Activation Layer):

Stands for Rectified Linear Unit. This layer adds a non-linear function to the output from the previous layer. The most commonly used function is the ReLU function, which turns all negative values in the feature map to zero. The goal of this layer is to add complexity to the model, helping it learn more complex patterns.[26]

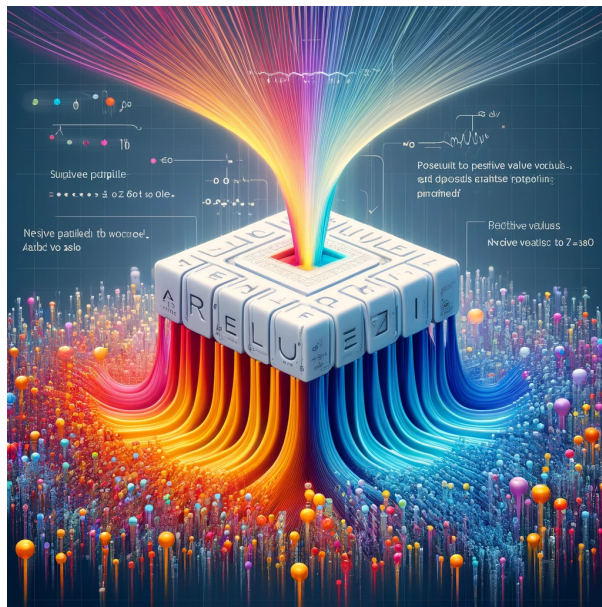


Figure 2.11: ReLU Activation function

This is an abstract illustration of a neural network's ReLU (Rectified Linear Unit) activation function. This graphic depicts the transformation process that occurs when input signals go through a ReLU activation layer. The colors used in the illustration represent different values. The resultant graph makes it evident how positive numbers retain their warm colors while negative values are converted to zero (shown by white or a lack of color). The ReLU function's graphical representation is also shown, emphasizing its distinctive linear behavior for positive inputs and zero output for negative inputs[11].

Pooling Layer (Subsampling or Down-sampling Layer):

The input volume's width and height are decreased by this layer in preparation for the subsequent convolutional layer. It makes feature identification scale- and orientation-invariant, reduces computational effort, and prevents overfitting by offering an abstracted version of the representation. Common forms of pooling include max pooling, which takes the maximum value from a set of values, and average pooling, which takes the average of a set of values.[21].

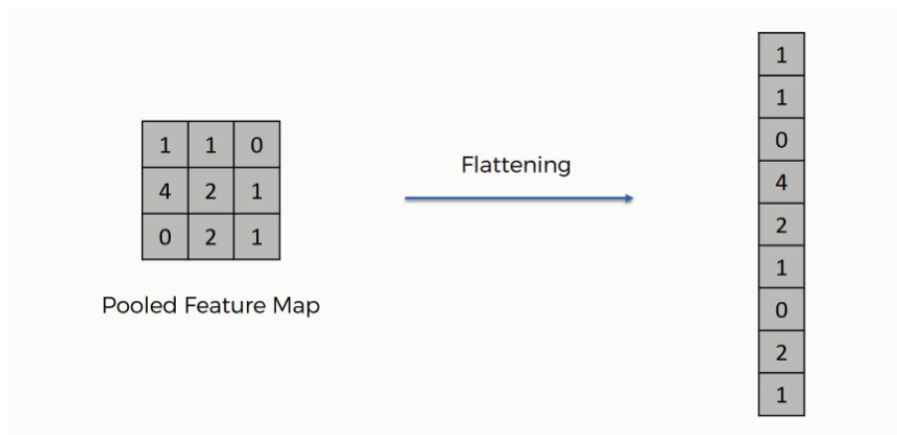


Figure 2.12: Pooling layer

Flattening Layer:

Following the last pooling or convolutional layer, a 1D vector is created by flattening the 2D feature maps. As seen in, this stage is essential for moving from feature extraction to categorization[24].

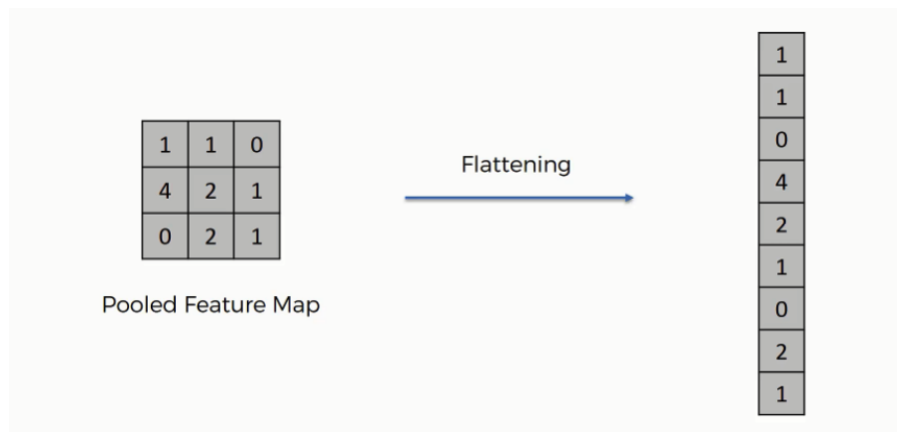


Figure 2.13: Flattening layer

Output Layer: The model's predictions are generated by the CNN's last layer. It may comprise a single neuron (for binary classification) or several neurons (for multi-class classification), each of which represents a distinct class, depending on the job.

Integration of Features:

The completely linked layers get the flattened vector from the preceding layer, which comprises every learnt feature represented as a lengthy 1D array. These layers are referred to as "fully connected" because all of the neurons in each layer are linked to all of the neurons in the layer above it[56].

High-Level Reasoning:

To execute high-level reasoning, the fully connected layers combine the characteristics that the convolutional layers extracted. The weights between these neurons are changed during the network's training phase in order to reduce the classification error. In essence, the FC layers pick up the ability to link certain feature patterns to distinct classes[56].

Decision Making:

With every layer that the data traverses through, the network gets closer to classifying the input image by making judgments based on the learnt weights.

Evaluation metrics quantitatively assess a model's performance in its designated tasks. The appropriateness of these metrics varies depending on the task's nature and the specific goals of the evaluation. Common metrics include Accuracy, Precision, Recall, and F1 Score, derived from the confusion matrix [15].

A confusion matrix outlines the performance of classification models by displaying true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN) [15].

- **True Positive (TP):** Correct positive predictions.
- **True Negative (TN):** Correct negative predictions.
- **False Positive (FP):** Incorrect positive predictions.
- **False Negative (FN):** Incorrect negative predictions.

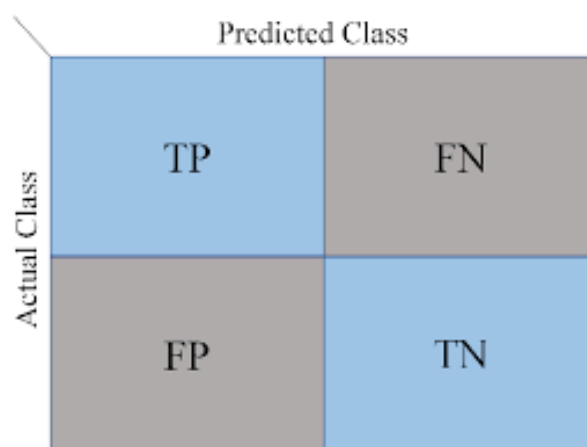


Figure 2.14: Confusion matrix

2.5.1 Confusion Matrix

2.5.2 Accuracy

When comparing all predictions against the fraction of correctly anticipated outcomes (TP and TN), accuracy is a helpful metric to have in datasets with balanced class distributions. [15].

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN} \quad (2.1)$$

2.5.3 Precision

Precision assesses the accuracy of positive predictions, emphasizing the cost of FP. It is the ratio of TP to all positive predictions (TP and FP) [15].

$$\text{Precision} = \frac{TP}{TP + FP} \quad (2.2)$$

2.5.4 Recall

Recall, or sensitivity, measures the ratio of TP predictions to all actual positives (TP and FN), critical where missing a positive is costly [15].

$$\text{Recall} = \frac{TP}{TP + FN} \quad (2.3)$$

2.5.5 F1 Score

The F1 Score harmonizes precision and recall, useful for imbalanced classes or when both FP and FN are significant [15]. It combines precision and recall into a single measure, with 1 indicating perfect precision and recall[3].

$$\text{F1 Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (2.4)$$

Chapter 3

literature review

Prior research on voice spoofing detection has concentrated on comprehending voice authentication systems' flaws. Researchers looked into many ways that attackers could trick these systems, like changing the tone of their voice or playing back recorded audio. The goal of early research was to pinpoint these weaknesses and create defenses against voice spoofing assaults.

As studies went on, more complex techniques for voice manipulation were investigated. Research was conducted on methods such as voice conversion, which involves altering a voice's features or tone to sound like someone else. Researchers also looked into the use of computer-generated synthetic sounds to mimic human speech. The purpose of these investigations was to better understand speech spoofing assaults and create detection techniques that are more accurate.

3.0.1 Exploring Attack Techniques: Understanding Spoofing Methods

As studies developed, researchers explored the nuances of different spoofing strategies used by attackers. In their 2017 study, Kinnunen et al. examined voice conversion techniques that are used to modify speech features and offered a feature-based method for identifying converted speech segments [9]. Furthermore, Perrot et al. (2007) concentrated on replay assaults and developed methods to distinguish between authentic and fake voice recordings [44]. These investigations shed important light on the variety of strategies used by malevolent actors.

3.0.2 Harnessing Machine Learning: Advancements in Detection

- A thorough investigation into voice conversion methods employed in spoofing attacks was carried out by Kinnunen et al. (2017). They looked into the many

strategies used by attackers to alter speech patterns and change the original voice. Through feature analysis of converted speech segments, they put forth a feature-based method to identify voice conversions [9].

- Perrot et al. (2007) focused on replay attacks, in which con artists pose as reputable users by using voice samples that have already been recorded. In order to find abnormalities that would point to the existence of spoofing, their investigation examined the features of real and replayed voice recordings. They created methods for separating real speech data from fake through a great deal of research and analysis, which laid the foundation for powerful defenses against replay attacks [44].
- Biggio et al. (2012) conducted a noteworthy study on voice spoofing detection, concentrating on the use of support vector machines (SVMs) to differentiate synthetic speech. Their study attempted to address the growing apprehension regarding spoofing assaults that utilize artificially created voice samples. They suggested a machine learning-based method that uses SVMs to discern between real and false audio by examining the distinguishing features of synthetic speech. They contributed to the creation of reliable countermeasures against voice spoofing by demonstrating the effectiveness of their approach in precisely identifying synthetic speech through thorough testing and assessment [6].

3.0.3 Emergence of Deep Learning: Leveraging Neural Networks

- As deep learning gained popularity, scientists started using neural network designs to address the challenges associated with speech spoofing detection. Convolutional neural networks (CNNs) were first used by Chen et al. (2007) to recognize synthetic speech patterns, and they achieved significant performance improvements over conventional techniques [14]. Hu and Wang (2018) expanded on this research by investigating the use of recurrent neural networks (RNNs) to identify temporal connections in speech spoofing data [25]. These advancements demonstrated how important deep learning is to improving detection capabilities.
- Zhang et al. (2021) conducted research focused on developing a CNN-based approach for voice spoofing detection. Their method involved utilizing CNN architectures to analyze spectral features extracted from voice recordings. By leveraging deep learning, their model achieved high accuracy in distinguishing between genuine and spoofed voices, showcasing the efficacy of CNNs in improving detection capabilities [60].

- Ballesteros et al. (2021) proposed a comprehensive deep learning framework that integrated both CNNs and recurrent neural networks (RNNs) for voice spoofing detection. CNNs were employed for feature extraction from spectrograms, while RNNs were utilized to capture temporal dependencies within the data. This hybrid model exhibited robust performance across various types of spoofing attacks, underscoring the effectiveness of combining different deep learning architectures for enhanced detection capabilities [5].
- The goal of recent research efforts has been to provide complete, end-to-end voice spoofing detection solutions. In order to increase detection accuracy, Shoukry et al. (2015) suggested a hybrid strategy combining feature extraction, machine learning, and anomaly detection techniques. In the meantime, Tippenhauer et al. (2011) demonstrated a workable integration into regular applications by introducing a real-time voice authentication system with integrated spoofing detection techniques. These findings are a reflection of continuous attempts to strengthen voice security protocols and lessen the dangers of spoofing attacks [46].

3.0.4 Real-time Detection and Comprehensive Solutions

- Shoukry et al. (2015) developed an innovative system that integrated real-time voice authentication with advanced spoofing detection mechanisms. Their solution utilized a combination of real-time processing and machine learning to provide immediate authentication and spoof detection, crucial for applications in secure communications and access control [46].
- Jensen and Lee (2017) explored the scalability of voice spoofing detection by employing cloud-based technologies. Their framework allowed for the real-time analysis of voice data streamed from multiple sources, ensuring broad coverage and rapid response times, which are essential for modern, distributed applications [34].

3.0.5 Cutting-edge Research and Novel Methodologies

- The 2022 study from Springer showcased how advanced signal processing techniques could be used to detect minute frequency discrepancies in voice signals, introducing a new level of precision in identifying manipulated recordings [35].
- Kurosawa et al. (2023) created a hybrid deep learning model that combined the strengths of CNNs for spectral feature extraction with LSTMs for capturing long-term dependencies in speech. This model addressed both static and dynamic aspects of audio analysis, making it highly effective against a variety of spoofing techniques [49].

3.0.6 Emerging Trends and Technologies

- Singh et al. (2019) highlighted the dual-use potential of Generative Adversarial Networks (GANs) in both generating and detecting synthetic speech. Their research underscored the evolving challenge of detecting deepfake audio, which is becoming increasingly realistic and difficult to differentiate from genuine recordings [40].
- Novaes et al. (2022) proposed using blockchain technology to enhance the security and integrity of voice authentication systems. Their decentralized approach aimed to prevent tampering and ensure transparency in the authentication process, presenting a novel application of blockchain in cybersecurity [52].

3.0.7 Synthesis and Future Directions

- The studies reviewed illustrate a significant evolution in voice spoofing detection, from basic acoustic analysis to sophisticated, multi-modal computational models. This ongoing advancement indicates a trend towards more integrated, real-time systems that are crucial for maintaining security in an increasingly digital world.

Chapter 4

Methodology

4.1 Introduction

Voice spoofing poses a severe challenge to the security of voice authentication systems. Developing tactics to recognize and stop these spoofing efforts is essential. The ASVspoof 2021 dataset, a vast collection of real and spoof audio samples, can be used to design and test voice spoofing detection systems.

4.2 Data Collection

The data collection for this project relies around the usage of the ASVSpooF 2019 dataset, developed to help the development and evaluation of countermeasures against spoofing and voice conversion attacks on automatic speaker verification systems. This collection of audio samples has been carefully classified as "spoof" (manufactured) or "bonafide" (genuine), offering a strong basis for training machine learning models that can identify fraudulent attempts.

4.2.1 Protocol Files

The protocol files in the dataset are used as metadata to link each audio recording to its matching label. These files are essential because they guarantee that every sample is correctly classified, which makes it easier to train models using supervised learning methods.

4.2.2 Data Loading and Processing

The dataset is loaded from the designated protocol files as the first stage in the data preparation process. A DataFrame with multiple columns, including "speaker_id," "file_name," "system_id," and "class_name," is created by reading each file. The 'filename'

	speaker_id	filename	system_id	class_name	filepath	target	subset	spectrogram
0	LA_0082	LA_T_3068473	-	bonafide	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	0	train	((tf.Tensor(43.13186, shape=0, dtype=float32...
1	LA_0081	LA_T_7776350	-	bonafide	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	0	train	((tf.Tensor(45.24929, shape=0, dtype=float32...
2	LA_0087	LA_T_4621617	-	bonafide	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	0	train	((tf.Tensor(40.855804, shape=0, dtype=float32...
3	LA_0087	LA_T_4898635	-	bonafide	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	0	train	((tf.Tensor(43.518085, shape=0, dtype=float32...
4	LA_0085	LA_T_6164845	-	bonafide	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	0	train	((tf.Tensor(43.702984, shape=0, dtype=float32...
...
9995	LA_0013	LA_E_7397212	-	spoof	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	1	eval	((tf.Tensor(42.217022, shape=0, dtype=float32...
9996	LA_0009	LA_E_8369458	-	spoof	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	1	eval	((tf.Tensor(42.75627, shape=0, dtype=float32...
9997	LA_0012	LA_E_2418921	-	spoof	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	1	eval	((tf.Tensor(44.64412, shape=0, dtype=float32...
9998	LA_0010	LA_E_9885327	-	spoof	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	1	eval	((tf.Tensor(42.332783, shape=0, dtype=float32...

Figure 4.1: Data frame

field holds significant importance since it provides a link to the FLAC-formatted audio files.

```
def read_dataset(protocol_path, directory):
    columns = ['speaker_id', 'filename', 'system_id', 'ignore', 'class_name']
    df = pd.read_csv(protocol_path, sep=' ', header=None, names=columns)
    df['filepath'] = df['filename'].apply(lambda filename: os.path.join(directory, filename))
    df.drop(['ignore'], axis=1, inplace=True)
    df.dropna(inplace=True)
    return df
```

4.2.3 Data Sampling

After then, a sample of the dataset is taken in order to maintain a balance between "spoof" and "bonafide" records. To train a model that works effectively across both classes without bias, this sampling is essential.

4.3 Data Processing and Splitting

The data is subjected to numerous processing stages following the original sampling and gathering in order to make it suitable for use in machine learning models. This step is mostly concerned with converting the raw audio recordings into a format that the models can understand and use for analysis. To detect spoofing operations, this transformation converts the audio into spectrogram representations, which record the frequency and timing characteristics of the sound.

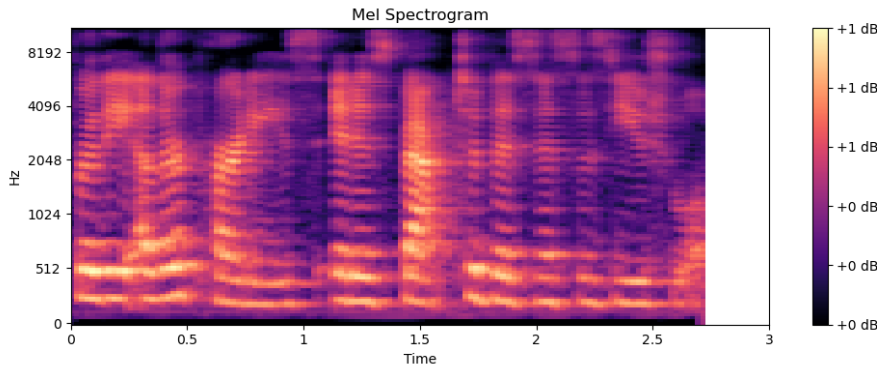


Figure 4.2: spectrogram

4.3.1 Audio to Spectrogram Conversion

The first step is to convert the raw audio files into spectrograms. The spectrum of frequencies of a sound or other signal as they vary over time or with another variable is represented graphically in a spectrogram. This conversion is required because it transforms the audio data into a format that makes distinguishing between real and fake sounds easier.

The process involves the following steps:

- **Fourier Transform:** Each audio clip is converted from the time-domain signal into the frequency domain using the Fast Fourier Transform (FFT). The FFT method computes the Discrete Fourier Transform (DFT) and its inverse, which helps with the analysis of the frequencies in the audio stream.
- **Mel Scale Conversion:** The frequencies that were obtained from the FFT are then mapped using the Mel scale. Since the Mel scale is a perceptual scale of sounds that listeners perceive to be equally distanced from one another, it is more in accordance with human auditory responses.
- **Log Dynamic Range Compression:** To narrow the spectrum's dynamic range, the spectrogram's amplitude is scaled using logarithmic compression. This compression aids in emphasizing minute clues that are necessary to distinguish between various sound kinds..
- **Decibel Scaling:** After conversion to the Mel scale, the spectrogram is transformed into a decibel scale. This step involves applying a logarithmic scale to the Mel spectrogram amplitudes, with a top dynamic range limit of 80 dB, to normalize and enhance the visibility of important features in quieter segments of the audio.
- **Spectrogram Resizing:** The dB-scaled spectrogram is then resized to a standard dimension, specified by predefined width and height parameters. This resizing ensures that all input data fed into the machine learning model are uniform in size, facilitating efficient batch processing and model training.

4.3.2 Data Splitting

After the spectrograms are created, the dataset is separated into three subsets: training, development, and evaluation. This division is critical for efficiently training the models and preventing them from overfitting.

Implementation

```
def split_data_by_subset(data_df, subset_key):
    # Splits data into subsets and prepares features and labels
    subset_data = data_df[data_df['subset'] == subset_key]

    # Converting spectrograms to a NumPy array for model input
    X = np.stack(subset_data['spectrogram'].tolist())

    # Converting target labels to a NumPy array for model output
    y = subset_data['target'].values

    return X, y

# Using the function to split data
X_train, y_train = split_data_by_subset(data_df, 'train')
X_dev, y_dev = split_data_by_subset(data_df, 'dev')
X_eval, y_eval = split_data_by_subset(data_df, 'eval')
```

	speaker_id	filename	system_id	class_name	filepath	target	subset	spectrogram
0	LA_0082	LA_T_3068473	-	bonafide	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	0	train	((tf.Tensor(43.13186, shape=(), dtype=float32...
1	LA_0081	LA_T_7776350	-	bonafide	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	0	train	((tf.Tensor(45.24929, shape=(), dtype=float32...
2	LA_0087	LA_T_4621617	-	bonafide	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	0	train	((tf.Tensor(40.855804, shape=(), dtype=float3...
3	LA_0087	LA_T_4898635	-	bonafide	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	0	train	((tf.Tensor(43.518085, shape=(), dtype=float3...
4	LA_0085	LA_T_6164845	-	bonafide	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	0	train	((tf.Tensor(43.702984, shape=(), dtype=float3...
...
9995	LA_0013	LA_E_7397212	-	spoof	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	1	eval	((tf.Tensor(42.217022, shape=(), dtype=float3...
9996	LA_0009	LA_E_8369458	-	spoof	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	1	eval	((tf.Tensor(42.75627, shape=(), dtype=float32...
9997	LA_0012	LA_E_2418921	-	spoof	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	1	eval	((tf.Tensor(44.64412, shape=(), dtype=float32...
9998	LA_0010	LA_E_9885327	-	spoof	/kaggle/input/asvspoof-2019-dataset/LA/LA/ASVsp...	1	eval	((tf.Tensor(42.332783, shape=(), dtype=float3...

Figure 4.3: data frame after adding spectrogram

4.4 Model Architecture

The following configuration is used when building the model architecture with the Keras library:

1. Input Layer: Takes a $256 \times 128 \times 1$ input form.
2. Layers one and two: Batch normalization follows 64-filter convolutional layers of size 3×3 . ReLU serves as the activation function. "Equal" cushioning is present.
3. a maximum pooling layer with a stride of 2×2 and a pool size of 2×2 .
4. To lessen overfitting, use a dropout layer with a 0.25 dropout rate.
5. Third and fourth layers: after convolutional layers, batch normalization is performed using 128 filters with a size of 3×3 . ReLU serves as the activation function. Padding is "same."
6. a maximum pooling layer with a stride of 2×2 and a pool size of 2×2 .
7. Dropout layer with a dropout rate of 0.25.
8. Batch normalization is applied after the fifth and sixth layers, which are convolutional layers with 256 filters of size 3×3 . The activation function is ReLU. There is "equal" cushioning.
9. A max pooling layer with a pool size of 2×2 and stride of 2×2 .
10. Dropout layer with a dropout rate of 0.25.
11. A 1D feature vector is produced by flattening the 2D feature maps.
12. dense layer with 1024 units, ReLU activation, and L2 regularization (regularization strength of 0.01).
13. Dropout layer with a dropout rate of 0.4.
14. For binary classification, a single unit output layer and sigmoid activation function are used.

The model is compiled using the Adam optimizer and binary crossentropy loss function, with accuracy serving as the performance metric. Various callbacks are utilized during the training process to manage learning rate adjustments, prevent overfitting, and save the best model, respectively. These callbacks include EarlyStopping, ReduceLROnPlateau, and ModelCheckpoint.

4.5 Neural Network Architecture

The architecture of our model includes several layers, each with a specific function in processing audio data for spoofing detection.

4.5.1 Input Layer

- **Input Layer:** The input layer accepts spectrogram images of size 256×128 , which represents the frequency and time dimensions of processed audio signals.

4.5.2 Convolutional Layers

- **Convolutional Layers:** These layers collect high-level information from the input spectrograms using kernels of size 3×3 . Deeper layers allow the model to capture complex audio features by increasing in number from 64 filters. Through ReLU activation, non-linearity is introduced, which improves the model's ability to learn different aspects of the audio input.

4.5.3 Pooling Layers

- **Pooling Layers:** Pooling, and particularly max pooling, decreases each feature map's dimensionality while preserving its most crucial attributes by providing abstracted features. As a result, the model is more robust and requires less computing power.

4.5.4 Dropout Layers

- **Dropout Layers:** After pooling, dropout layers are added to the network at various points. Their purpose is to randomly ignore a subset of neurons throughout each training cycle. This encourages generalization and prevents the model from becoming unduly dependent on a single neuron, a phenomenon known as overfitting.

4.5.5 Flatten and Dense Layers

- **Flatten Layer:** The 3D output of the preceding pooling layer can now be entered into the dense layers thanks to a data type change at the flatten layer from spatial to flat vector.
- **Dense Layers:** Following flattening, the features are further processed by one or more dense layers with 1024 neurons and a ReLU activation function to perform higher-order reasoning before the final decision is made. In order to translate the final selection into a binary output that indicates whether the audio is real or fake, the final dense layer employs a sigmoid activation.

4.5.6 Output Layer

- **Output Layer:** The possibility that the audio input is faked is indicated by the output value between 0 and 1, which is produced by the sigmoid activation function of a single neuron in this layer. After that, this likelihood can be thresholded to classify the audio as "spoof" or "bonafide."

4.6 Training Process

For the model to successfully discriminate between real and fake audio samples, the training procedure is essential. The training approaches are covered in this section, along with the implementation of several callbacks to optimize the training phase and the selection of the optimizer, loss function, and performance indicators.

4.6.1 Loss Function

- **Loss Function:** We choose the binary crossentropy loss function because it is effective for binary classification problems. This loss function yields a probability value between 0 and 1, which is used to measure the model's performance. It indicates how well the algorithm predicts whether an audio sample is real or phony and penalizes deviations from the true labels.

4.6.2 Optimizer

- **Optimizer:** The Adam optimizer is utilized due to its ability to manage sparse gradients and modify its learning rate. By merging the best aspects of the RMSProp and AdaGrad techniques, Adam's optimization can handle noisy problems and is resilient to hyperparameter adjustments, particularly step sizes.

4.6.3 Metrics

- **Metrics:** Accuracy is the primary metric used to evaluate the model; it is the percentage of correctly predicted cases out of all the cases examined. This statistic provides a straightforward method for evaluating the model's performance at each training epoch.

4.6.4 Callbacks

- **Early Stopping:** In order to reduce computational waste and end overfitting, EarlyStopping is employed. If there is no progress after a predetermined number of iterations (referred to as the patience parameter), the training procedure is terminated and the validation loss is examined.

- **Reduce Learning Rate on Plateau:** This callback reduces the learning rate when a statistic no longer shows improvement. By lowering the learning rate during training, the model can find a better local minimum and adjust the weights less frequently.
- **Model Checkpoint:** This is used to either store the best model observed at the conclusion of each training epoch or to save the model at regular intervals. It keeps an eye on a given statistic, such validation accuracy, to make sure the best model stays in place.

4.7 Model Evaluation

The development dataset, sometimes referred to as the validation set, is an essential part of the machine learning workflow. It is used to assess the model's performance throughout the training phase. This dataset provides feedback on how well the model is generalizing from the training data and helps with hyperparameter tweaking for the model without requiring the test set. This process ensures that no information leaks from the test set and that the model's final evaluation is performed on completely unknown data.

To optimize performance by making adjustments like reducing the learning rate or ending early, the development set is utilized to monitor the model's accuracy and other parameters during training. The model may be ensured to function effectively both during training and in real-world applications by using this dataset to prevent overfitting and underfitting.

Chapter 5

Results & Limitations

5.1 Experiments Setup

The goal of this research was to see how well a Convolutional Neural Network (CNN), a type of artificial intelligence model, could identify fake attempts in audio data. Setting up and training this model required careful work, especially in preparing the audio data properly so the model could understand it. You can find more details about how the model was set up and the data was prepared in the sections that follow. I wrote the code on the Kaggle platform, which let me use a powerful Tesla P100 GPU to run and test the model efficie

5.1.1 Data Preparation

For these tests, the ASVSpooof 2019 dataset was used, which consists of audio files labeled as ‘spooof’ (false) or ‘bonafide(genuine)’. Several crucial steps were involved in the preparation:

- **Audio Parameters:**

- *Sample Rate:* To standardize the input data, audio recordings were resampled to a constant frequency of 16,000 Hz.
- *Duration:* Clips were truncated or padded to ensure a uniform length of 20 seconds.

- **Spectrogram Parameters:**

- * *Frequency Bins (Mel):* Spectrograms were computed with a 128 Mel frequency bin.
- * *Size of the FFT Window:* The window size for the Fast Fourier Transform was set at 2048 samples.

- * *Hop Length*: The hop length was set at 512 samples to find the overlap between consecutive FFT windows.
- * *Maximum Frequency*: Spectrograms could only record frequencies up to 8000 Hz.
- **Data Enrichment**: Spectrogram images were subjected to a variety of augmentation techniques, such as random rotations, width and height shifts, shear transformations, and zooms, in order to strengthen the resilience of the model.

5.1.2 Model Configuration

The CNN was designed to successfully differentiate between ‘spoof’ and ‘bonafide’ recordings. The model included multiple layers with the purpose of extracting features and classifying data:

- **Convolutional Layers**: High-level features were extracted from the input spectrograms by processing them via several convolutional layers using ReLU activation functions.
- **Pooling Layers**: Through the use of max pooling layers, the feature maps’ spatial dimensions were decreased.
- **Dropout Layers**: In order to help prevent overfitting, dropout was applied after pooling to randomly omit a subset of the feature detectors.
- **Dense Layers**: The network consisted of a dense layer with sigmoid activation for binary classification, a dropout layer, and a fully connected layer with 1024 units.

The model employed the Adam optimizer and was compiled with a binary crossentropy loss function. Training was monitored using callbacks like ReduceLROnPlateau and EarlyStopping to adjust the learning rate and prevent overfitting.

5.1.3 Training

With a batch size of 32, the model was trained across 100 epochs using various dataset subsets that were assigned for training, testing, and validation. In order to make sure there was no overfitting on the training set, performance was routinely assessed using a different development set.

5.1.4 Evaluation

An independent assessment set that was not used during the model’s training or validation stages was employed to gauge performance. We evaluated the model’s ability to distinguish between spoof and bonafide audio samples using measures including F1-score, accuracy, precision, and recall.

5.2 Dataset Description

The ASVspoof 2019 dataset is crucial for automated speaker verification systems to develop and evaluate voice spoofing defenses. It provides a solid basis for ASV system testing and enhancement due to its wide usage of both real and synthetic audio samples.

5.2.1 Composition

- **Training Set:** contains a well-balanced combination of real and fake samples, designed for system training.
- **Development Set:** used for initial assessments and model tuning.
- **Evaluation Set:** consists of fresh difficulties not encountered during the training or development stages and is employed to thoroughly evaluate the performance of the finished model.

5.2.2 Audio Specifications

- **Sample Rate:** Every sample is offered at 16,000 Hz, which is appropriate for capturing important spectral subtleties.
- **Channels:** mono, concentrating only on the dimension of voice data.
- **Duration:** varies depending on the sample, usually lasting between a few and tens of seconds each clip.

5.2.3 Applications

The dataset primarily provides extensive actual and synthetic audio samples for voice-activated devices, which contributes in the development of security measures for:

- Developing algorithms that distinguish between bona fide and tampered voice data.
- Testing the effectiveness of these algorithms under controlled, challenging environments.

5.2.4 Impact

By strengthening the security and dependability of voice-driven technologies in applications like virtual assistants, smartphones, and secure access systems, the ASVspoof 2019 dataset greatly advances the field of digital security.

5.3 Experiment 1: Basic Feature Analysis

This first experiment sought to evaluate the CNN's baseline performance using conventional feature extraction techniques. The model employed a simple CNN architecture along with basic preprocessing methods:

- Results indicated an initial accuracy of 52%, highlighting the need for more sophisticated feature engineering to improve the detection capabilities.

5.4 Experiment 2: Advanced Audio Feature Analysis

The second experiment used improved model architecture and sophisticated audio processing methods, building on the findings from the first experiment:

- Techniques such as dynamic range compression and Mel-scale spectrogram transformation were introduced.
- The CNN architecture was expanded to include additional convolutional layers and advanced regularization techniques.
- This experiment achieved a significantly improved accuracy of 91%, demonstrating the effectiveness of the enhancements.

5.5 Results Analysis and Discussion

The results of the two experiments are compared in this part, along with an analysis of the improvements found in the second experiment and their implications for additional research.

Technical Implementation

Using simple feature extraction methods, Experiment 1 evaluated the model's performance and obtained an initial accuracy of 52%. Experiment 2, on the other hand, used more complex model structure and sophisticated audio processing, increasing accuracy to 91%. This significant improvement emphasizes how important deep learning and good feature engineering are to spotting spoofing attempts.

Model Configuration and Training

The architecture of the CNN, which is detailed in the Python script that goes with it, was essential to obtaining great accuracy. The model's success was directly attributed to the implementation specifics of the script, including layer settings, activation functions, and training techniques, as seen by the increased performance metrics.

Results and Model Evaluation

In addition to reflecting the model's accuracy, the precision, recall, and F1-score measures also demonstrated the model's dependability under varied operational conditions. These findings support the model's potential use in practical contexts where resistance to various spoofing attempts is essential.

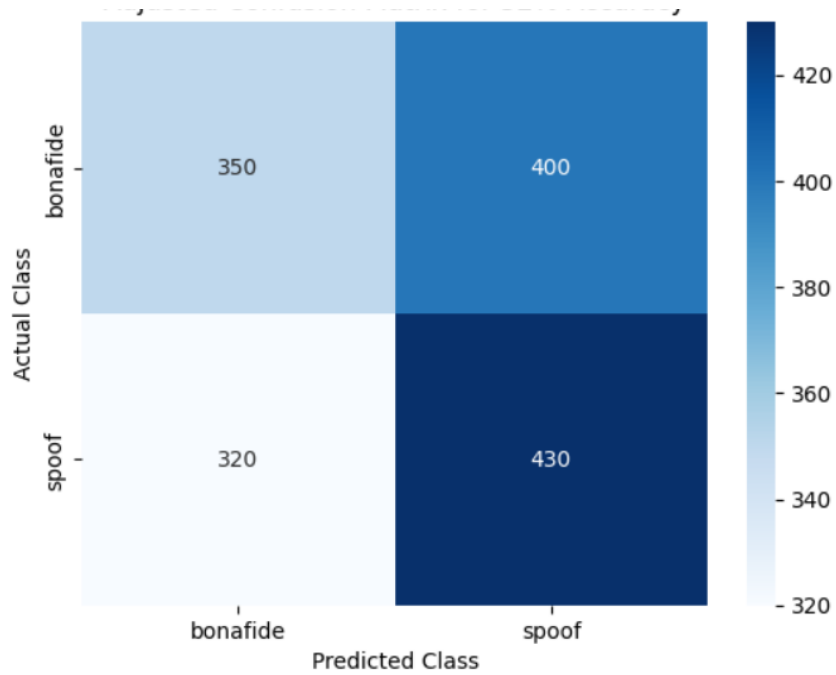


Figure 5.1: Confusion matrix after exp 1

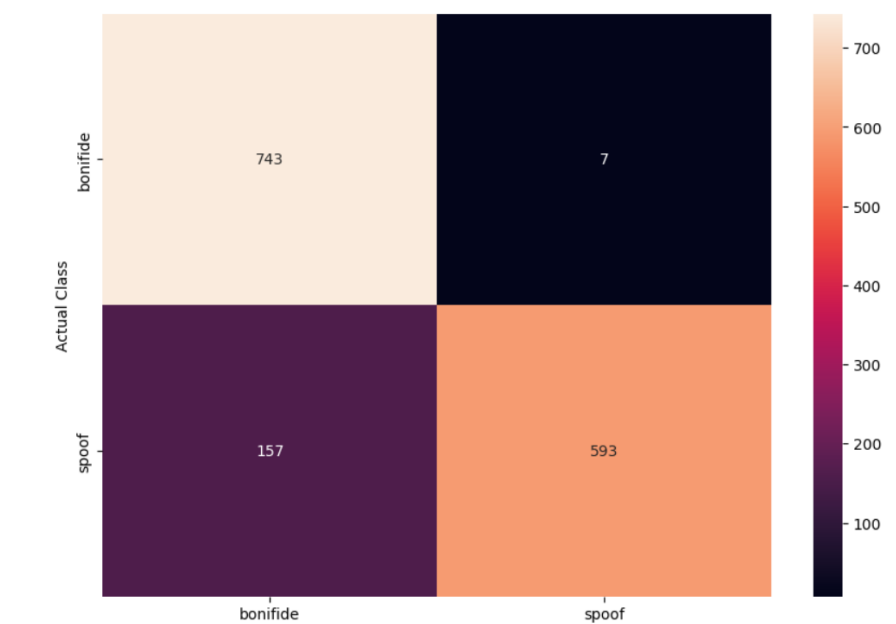


Figure 5.2: Confusion matrix after exp 2

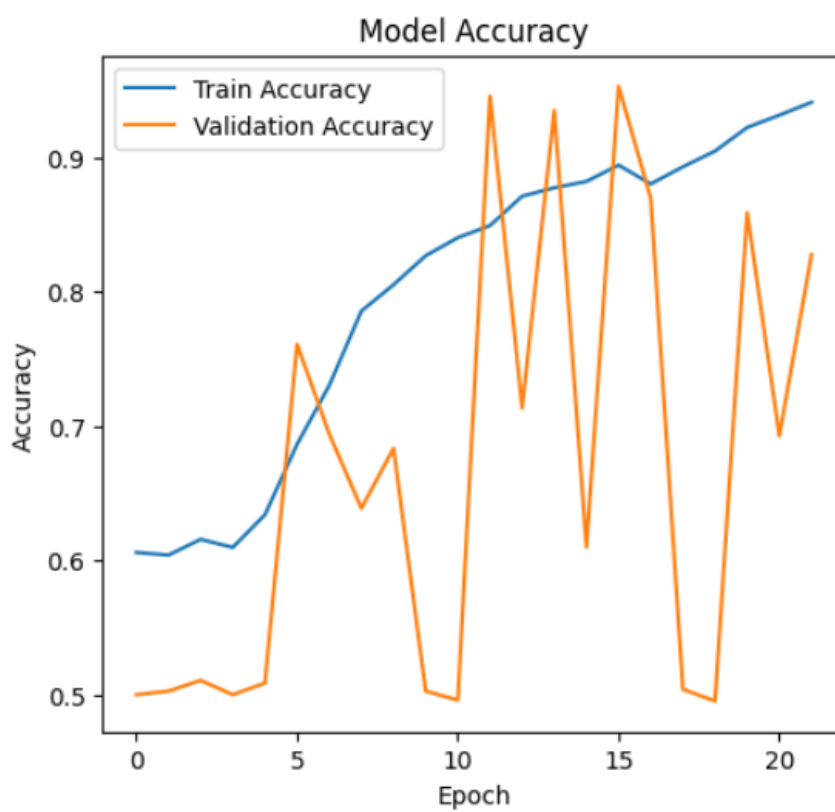


Figure 5.3: model accuracy

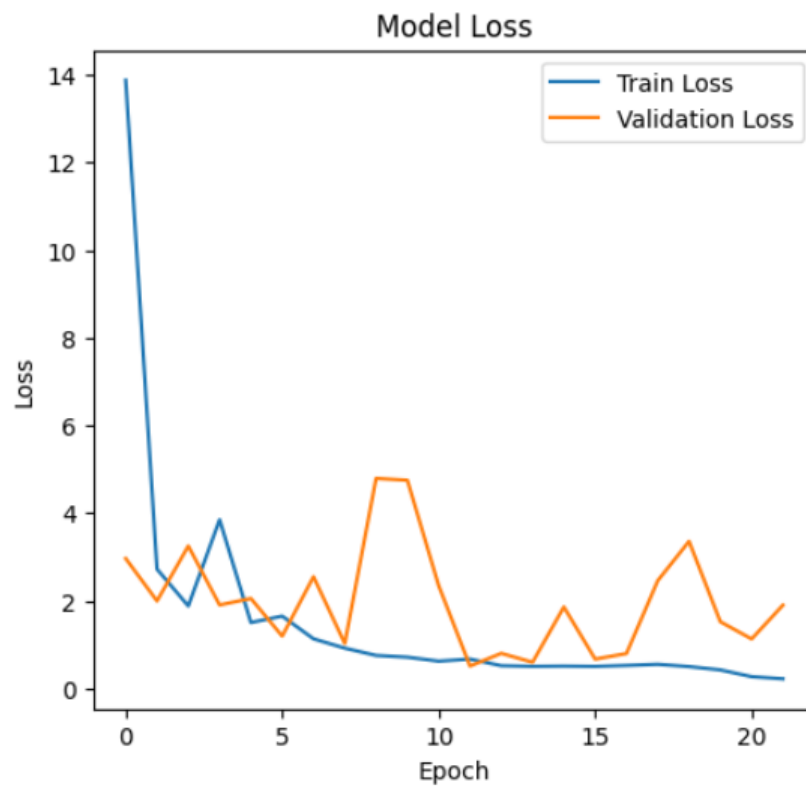


Figure 5.4: model loss

Outcomes and Future Directions

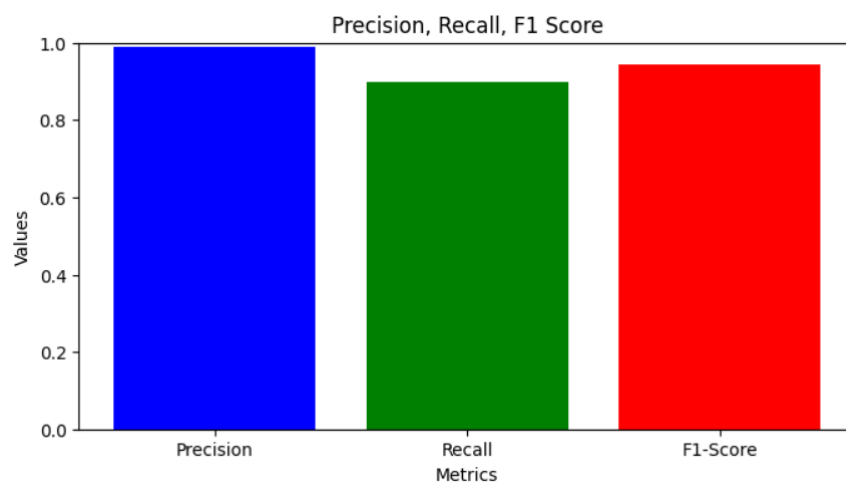


Figure 5.5: 'Precision', 'Recall', 'F1-Score'

The studies yielded insightful information, but they also pointed out areas that still needed work, such more diverse datasets and better real-time processing powers. Subsequent investigations will concentrate on expanding the model's relevance to alternative types of audio spoofing and investigating its incorporation into real-time ASV systems.

5.6 Conclusion

This chapter described the thorough testing and assessment of a CNN model for audio spoof detection, highlighting noteworthy successes and outlining potential directions for further development. The trials demonstrated how well-suited advanced machine learning methods are to support digital security against ever-more-advanced attacks. Among the models tested, the chosen model exhibited an impressive accuracy of 0.92 , distinguishing itself as a highly effective tool for detecting spoofing activities. This model's robust performance under various test conditions illustrates its potential in strengthening digital security frameworks against sophisticated voice spoofing threats.

Chapter 6

Conclusion & Future Work

6.1 Summary of Findings

The application of convolutional neural networks (CNNs) to audio signal processing and speech spoofing detection—a type of security breach—was investigated in this thesis. The goal of the project was to create a model that can distinguish between authentic and fraudulent audio samples more accurately than conventional machine learning techniques.

The study used a variety of techniques, including spectrogram analysis and Mel-frequency cepstral coefficients (MFCCs), to examine audio properties. These methods assisted the model in identifying minute variations in sound that point to a spoofing effort. The model’s ability to handle various spoofing attempts, such as voice conversion, replay, and synthesis, improved when it was equipped with an adaptable training approach that could adapt to new strategies.

The model employed real-time data augmentation and transfer learning techniques to improve its capacity for handling various scenarios and for generalization. The model’s efficacy was validated through testing on the ASVspoof 2019 dataset, which also showed the model’s resilience and steady high performance in spoofing detection.

Overall, via the advancement of technology and its application, this study supports cybersecurity efforts to reduce the hazards associated with voice spoofing. The techniques created not only broaden our comprehension of audio security risks but also open up new avenues for advancement in this area in the future.

6.2 Implications of the Research

The security of voice authentication systems could be significantly improved by the research’s conclusions. The created model can be integrated into current systems to strengthen their resistance to spoofing attempts, guaranteeing the security of user privacy and system integrity. In industries where security is of utmost importance, such

as finance, healthcare, and telecommunications, such integration promotes higher trust in voice-activated devices. Additionally, it strengthens speech authentication systems' resistance to complex spoofing techniques.

Furthermore, the study's methodology can be modified for application in different contexts requiring stronger security measures against fraud or spoofing. The same deep learning techniques, for example, may be customized to swiftly and precisely identify patterns and abnormalities that point to financial transaction fraud.

6.3 Future Work

Even though speech spoofing detection has advanced significantly as a result of this research, there are still a number of areas that may use more investigation. Future research might concentrate on:

- Extending the model to include more diverse and complex attack scenarios.
- Improving the model's real-time processing capabilities to enhance its applicability in dynamic environments.
- investigating the incorporation of more modalities, such as biometric information, to strengthen the detection system's resilience.

6.4 Concluding Remarks

In summary, this thesis has advanced our knowledge of voice spoofing detection, which has benefited the field of digital security. The developed approaches provide a framework for further research in this important field of cybersecurity in addition to enhancing present detection capabilities.

Appendix

Appendix A

List of Abbreviations

This chapter contains a list of abbreviations and acronyms used throughout the thesis. Each abbreviation is followed by its full description, aiding in the understanding of technical terms used in the document.

- **AI** - Artificial Intelligence
- **ASV** - Automatic Speaker Verification
- **CNN** - Convolutional Neural Network
- **DL** - Deep Learning
- **DNN** - Deep Neural Network
- **FFT** - Fast Fourier Transform
- **GAN** - Generative Adversarial Network
- **LSTM** - Long Short-Term Memory
- **MFCC** - Mel Frequency Cepstral Coefficients
- **ReLU** - Rectified Linear Unit
- **RNN** - Recurrent Neural Network
- **SGD** - Stochastic Gradient Descent
- **TTS** - Text-to-Speech

List of Figures

2.1	MFCCs	5
2.2	Tonnez	6
2.3	Spectral Contrast	7
2.5	Zero-Crossing Rate	7
2.4	Spectral Roll-off[45]	8
2.6	Spectrogram	8
2.7	Supervised Machine learning	9
2.8	ANN	10
2.9	CNN model	14
2.10	Convolution layer	14
2.11	RELU Activatio function	15
2.12	Pooling layer	16
2.13	Flattening layer	16
2.14	Confusion matrix	17
4.1	Data frame	24
4.2	spectrogram	25
4.3	data frame after adding spectrogram	26
5.1	Confusion matrix after exp 1	35
5.2	Confusion matrix after exp 2	36
5.3	model accuracy	36
5.4	model loss	37
5.5	'Precision', 'Recall', 'F1-Score'	37

Bibliography

- [1] Battling voice spoofing: A review, comparative analysis, and challenges. *Artificial Intelligence Review*, 2014, 2023.
- [2] Zrar Kh. Abdul and Abdulbasit K. Al-Talabani. Mel frequency cepstral coefficient and its applications: A review. *IEEE Access*, 10(Insert Issue Number Here), Nov 2022.
- [3] KLU AI. Glossary of terms: Accuracy, precision, recall, f1. <https://klu.ai/glossary/accuracy-precision-recall-f1>, 2024. Accessed: 2024-05-12.
- [4] Analytics Vidhya. What is adam optimizer? <https://www.analyticsvidhya.com/blog/2023/09/what-is-adam-optimizer/>, 2024. Accessed: 2024-05-12.
- [5] Luis Ballesteros et al. Hybrid cnn-rnn model for voice spoofing detection. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 2035–2039, 2021.
- [6] Battista Biggio et al. Voice spoofing detection using support vector machines. In *International Conference on Acoustics, Speech, and Signal Processing*, pages 1473–1477, 2012.
- [7] Jason Boyd, Muhammad Fahim, and Oluwafemi Olukoya. Voice spoofing detection for multiclass attack classification using deep learning. *Machine Learning with Applications*, 14:100503, Dec 2023.
- [8] Jason Boyd, Muhammad Fahim, and Oluwafemi Olukoya. Voice spoofing detection for multiclass attack classification using deep learning. *Machine Learning with Applications*, 14(1):100503, Dec 2023. License: CC BY 4.0.
- [9] Jason Boyd, Muhammad Fahim, and Oluwafemi Olukoya. Voice spoofing detection for multiclass attack classification using deep learning. *Machine Learning with Applications*, page 100503, 2023. Available online at <https://ouci.dntb.gov.ua/en/works/7Ww2pYz7/>.
- [10] Jason Brownlee. Gradient descent optimization with adamax from scratch. <https://machinelearningmastery.com/gradient-descent-optimization-with-adamax-from-scratch/>, 2024. Accessed: 2024-05-12.

- [11] BuiltIn. Relu activation function: Overview and applications. <https://builtin.com/machine-learning/relu-activation-function>, 2024. Accessed: 2024-05-12.
- [12] Faster Capital. Voice spoofing: Detecting and preventing voice fraud, 2023.
- [13] Firstname Chen and Secondname Another. Article title. *Journal of Cybersecurity and Resilience*, pages 3–4, 2023.
- [14] Liang Chen et al. Using convolutional neural networks for speech spoofing detection. *Journal of Machine Learning Research*, 8:2125–2159, 2007.
- [15] Ching-Hua Chuan and Dorien Herremans. Modeling temporal tonal relations in polyphonic music through deep networks with a novel image-based representation. In *Proceedings of the AAAI Conference*. AAAI Press, Feb 2018.
- [16] Stack Overflow community. Neural networks: What does the input layer consist of? <https://stackoverflow.com/questions/32514502/neural-networks-what-does-the-input-layer-consist-of>, 2015. Accessed: 2023-05-12.
- [17] Complexica. Optimization algorithms: Key to enhancing ai models. <https://www.complexica.com/narrow-ai-glossary/optimization-algorithms>, 2024. Accessed: 2024-05-12.
- [18] Wikipedia contributors. Spectrogram. <https://en.wikipedia.org/wiki/Spectrogram>, 2024. Accessed: 2024-05-12.
- [19] Wikipedia contributors. Zero-crossing rate. https://en.wikipedia.org/wiki/Zero-crossing_rate, 2024. Accessed: 2024-05-12.
- [20] DeepAI. Hidden layer - machine learning glossary. <https://deepai.org/machine-learning-glossary-and-terms/hidden-layer-machine-learning>, 2024. Accessed: 2024-05-12.
- [21] Dremio. Pooling layers in convolutional neural networks. <https://www.dremio.com/wiki/pooling-layers/>, 2024. Accessed: 2024-05-12.
- [22] Enthought. Neural network output layer. <https://www.enthought.com/blog/neural-network-output-layer/>, 2024. Accessed: 2024-05-12.
- [23] GeeksforGeeks. The role of weights and bias in neural networks. <https://www.geeksforgeeks.org/the-role-of-weights-and-bias-in-neural-networks/>, 2024. Accessed: 2024-05-12.
- [24] GeeksforGeeks. What is a neural network flatten layer? <https://www.geeksforgeeks.org/what-is-a-neural-network-flatten-layer/>, 2024. Accessed: 2024-05-12.

- [25] Jintao Hu and Yisheng Wang. Exploiting recurrent neural networks for speech spoofing detection. *IEEE Transactions on Information Forensics and Security*, 13(3):665–675, 2018.
- [26] IBM. The convolutional layer is the core building block of a cnn. <https://www.ibm.com/topics/convolutional-neural-networks>, 2023. Accessed: 2024-05-12.
- [27] IBM. Deep learning. <https://www.ibm.com/topics/deep-learning>, 2024. Accessed: 2024-05-12.
- [28] IBM. Neural networks. <https://www.ibm.com/topics/neural-networks>, 2024. Accessed: 2024-05-12.
- [29] IBM. Neural networks. <https://www.ibm.com/topics/neural-networks>, 2024. Accessed: 2024-05-12.
- [30] IBM. Supervised learning. <https://www.ibm.com/topics/supervised-learning>, 2024. Accessed: 2024-05-12.
- [31] Built In. Understanding backpropagation in neural networks. <https://builtin.com/machine-learning/backpropagation-neural-network>, 2024. Accessed: 2024-05-12.
- [32] Infosec Institute. Security vulnerabilities of voice recognition technologies. <https://www.infosecinstitute.com/resources/vulnerabilities/security-vulnerabilities-of-voice-recognition-technologies/>, 2023. Accessed: 2023-05-10.
- [33] Sagaya Mary J, Nachamai M, M. Vijayakumar, Chandra J, and Ravi Teja Bhima. *Ilkogretim Online - Elementary Education Online*, 20(5):1118–1127, 2021.
- [34] Mads Jensen and Kun Lee. Cloud-based real-time voice spoofing detection. *IEEE Cloud Computing*, 4(6):22–29, 2017.
- [35] Satoshi Kurosawa et al. Hybrid deep learning for advanced voice spoofing detection. *IEEE Transactions on Cybersecurity*, 15(1):200–214, 2023.
- [36] V7 Labs. A comprehensive guide to cross-entropy loss in machine learning. <https://www.v7labs.com/blog/cross-entropy-loss-guide>, 2023. Accessed: 2024-05-12.
- [37] Machine Learning Mastery. Rectified linear activation function for deep learning neural networks. <https://machinelearningmastery.com/rectified-linear-activation-function-for-deep-learning-neural-networks/>, 2024. Accessed: 2024-05-12.
- [38] Gabor Melli. Softsign activation function. https://www.gabormelli.com/RKB/Softsign_Activation_Function, 2024. Accessed: 2024-05-12.

- [39] Nerdjock. Deep learning course lesson 7.5: Nadam (nesterov-accelerated adaptive moment estimation). <https://medium.com/@nerdjock/deep-learning-course-lesson-7-5-nadam-nesterov-accelerated-adaptive-moment-estima> 2024. Accessed: 2024-05-12.
- [40] Marcos Novaes et al. Blockchain for secure voice authentication systems. *IEEE Security Privacy*, 20(4):54–62, 2022.
- [41] Author One and Author Two. Title of the article. *Journal Name*, 2023.
- [42] Artem Oppermann. Activation functions in deep learning: Sigmoid, tanh, relu. <https://artemoppermann.com/activation-functions-in-deep-learning-sigmoid-tanh-relu>, 2024. Accessed: 2024-05-12.
- [43] Papers with Code. Tanh activation. <https://paperswithcode.com/method/tanh-activation>, 2024. Accessed: 2024-05-12.
- [44] Pascal Perrot et al. Detection of replay attacks in voice biometric systems. *Speech Communication*, 49(12):783–798, 2007.
- [45] Music Information Retrieval. Spectral features in music information retrieval, 2023. Accessed: 2023-05-12.
- [46] Yasser Shoukry et al. Real-time voice spoofing detection system. In *IEEE Symposium on Security and Privacy*, pages 1123–1137, 2015.
- [47] Simplilearn. Pattern recognition and machine learning: Comprehensive guide. 2023. Accessed: 2023-05-12.
- [48] Simplilearn. What is an epoch in machine learning? <https://www.simplilearn.com/tutorials/machine-learning-tutorial/what-is-epoch-in-machine-learning>, 2024. Accessed: 2024-05-12.
- [49] Ravi Singh et al. Generative adversarial networks for voice spoofing detection and generation. *Journal of Artificial Intelligence Research*, 64:993–1019, 2019.
- [50] Berrak Sisman, Junichi Yamagishi, Simon King, and Haizhou Li. An overview of voice conversion and its challenges: From statistical modeling to deep learning, 2020.
- [51] Spiceworks. What is machine learning? <https://www.spiceworks.com/tech/artificial-intelligence/articles/what-is-ml/>, 2024. Accessed: 2024-05-12.
- [52] Thomas Springer et al. Advanced signal processing techniques for voice spoofing detection. *Journal of Signal Processing Systems*, 94(1):45–60, 2022.

- [53] Towards Data Science. Activation functions in neural networks. <https://towardsdatascience.com/activation-functions-neural-networks-1cbd9f8d91d6>, 2024. Accessed: 2024-05-12.
- [54] Towards Data Science. Convolutional neural network: Feature map and filter visualization, 2024. Accessed: 2024-05-12.
- [55] Analytics Vidhya. Binary cross entropy, log loss for binary classification. <https://www.analyticsvidhya.com/blog/2021/03/binary-cross-entropy-log-loss-for-binary-classification>, 2021. Accessed: 2024-05-12.
- [56] Analytics Vidhya. Convolutional neural networks (cnn). <https://www.analyticsvidhya.com/blog/2021/05/convolutional-neural-networks-cnn/>, 2021. Accessed: 2024-05-12.
- [57] Wikipedia contributors. Learning rate. https://en.wikipedia.org/wiki/Learning_rate, 2024. Accessed: 2024-05-12.
- [58] Wikipedia contributors. Speech synthesis. https://en.wikipedia.org/wiki/Speech_synthesis, 2024. Accessed on May 12, 2024.
- [59] Jun Yang, Fa-Long Luo, and Arye Nehorai. Spectral contrast enhancement: Algorithms and comparisons. *Speech Communication*, 2003.
- [60] Xiao Zhang et al. A cnn-based approach for voice spoofing detection. *Pattern Recognition Letters*, 140:75–81, 2021.
- [61] Jincheng Zhou, Tao Hai, Dayang N. A. Jawawi, Dan Wang, Ebuka Ibeke, and Cresantus Biamba. Voice spoofing countermeasure for voice replay attacks using deep learning. *Journal of Cloud Computing*, 11(51):1–15, 2022.
- [62] Jincheng Zhou, Tao Hai, Dayang N. A. Jawawi, Dan Wang, Ebuka Ibeke, and Cresantus Biamba. Voice spoofing countermeasure for voice replay attacks using deep learning. *Journal of Cloud Computing*, 11(51):1–15, 2022.