

NATIONAL CENTER FOR SCIENTIFIC RESEARCH
"DEMOKRITOS"



UNIVERSITY OF THE PELOPONNESE



Deep Learning

Prof: Gianakopoulos Theodoros

Asimoglou Menelaos (ID: 2022202200001)

Report on Automatic Speaker Verification (ASV) Anti-Spoofing using CNNs and GANs



1. Introduction

- 1.1 Background and Importance of ASV
- 1.2 The Challenge of Spoofing and the Need for Anti-Spoofing Measures

2. Overview of ASV and Spoofing Attacks

- 2.1 Detailed Description of ASV Systems
- 2.2 Understanding the Nature of Spoofing Attacks
- 2.3 The Significance of Anti-Spoofing Mechanisms

3. Deep Learning in ASV Anti-Spoofing

- 3.1 Introduction to Convolutional Neural Networks (CNNs)
- 3.2 Introduction to Generative Adversarial Networks (GANs)
- 3.3 Role and Advantages of CNNs and GANs in ASV Anti-Spoofing

4. Methodology

- 4.1 Description of the Dataset Used
- 4.2 Explanation of CNN-based ASV System
- 4.3 Description of GAN Techniques Applied for Spoofing
- 4.4 Explanation of CNN-based Anti-Spoofing Techniques

5. Experimental Setup

- 5.1 Hardware and Software Specifications used
- 5.2 Setup of CNN-based ASV System
- 5.3. Execution of GAN-based Spoofing Attacks
- 5.4 Implementation of CNN-based Anti-Spoofing Measures



6. Results and Discussion

- 6.1 Performance of CNN-based ASV System
- 6.2 Efficacy of GAN-based Spoofing Attacks
- 6.3 Effectiveness of CNN-based Anti-Spoofing Measures
- 6.4 Comparative Analysis of Results

7. Conclusion and Future Work

- 7.1 Summary of Findings
- 7.2 Implications and Recommendations



Abstract

Focusing on the usage of CNNs and GANs in ASV Anti-Spoofing, with each section detailing different aspects of the project. It begins with an introduction to ASV and the issues surrounding spoofing attacks, then moves on to the methodologies and experimental setup, and finally, presents the results, concluding remarks, and future recommendations.



1. Introduction

The evolving field of technology has led to significant advancements in many sectors, with communication being a pivotal one. With the prevalence of digital communication platforms, there is an increasing necessity to ensure that the identity of individuals taking in these platforms is authentic. This introduces the domain of Automatic Speaker Verification (ASV).

1.1 Background and Importance of ASV

Automatic Speaker Verification (ASV) is a system that uses biometric verification to confirm the identity of a speaker based on unique patterns and features of their voice. ASV takes place in various sectors including customer service, security, forensics, and personal virtual assistants. Its main strength lies in its ability to provide non-intrusive and natural biometric verification, which is exceptionally important in a world where digital communication is becoming more and more prevalent. It adds a layer of security by confirming the identity of the speaker ensuring that the access and exchange of information are being conducted by authorized individuals. As such ASV forms an integral part of identity verification in our current and future digital world.

1.2 The Challenge of Spoofing and the Need for Anti-Spoofing Measures

However, as ASV systems have become more sophisticated, so have the techniques to deceive them. Spoofing is one such deceptive practice where an impostor attempts to mimic the voice of a genuine speaker with the intention to fool the ASV system. Spoofing attacks pose a significant threat to the integrity and reliability of ASV systems, leading to potential security breaches and unauthorized access to sensitive information.

Given the serious implications of successful spoofing attacks, there is a crucial need for effective anti-spoofing measures. These measures are designed to equip ASV systems with the ability to detect and prevent spoofing attempts, thereby ensuring that the systems remain robust against such fraudulent activities. The development and implementation of anti-spoofing measures thus stand as a critical area of research and development in the field



of ASV. This report is into this topic, exploring the use of advanced techniques like Convolutional Neural Networks (CNNs) and Generative Adversarial Networks (GANs) in ASV anti-spoofing.

2. Overview of ASV and Spoofing Attacks

2.1 Detailed Description of ASV Systems

Automatic Speaker Verification (ASV) systems operate based on the unique vocal characteristics of an individual. They function on two primary processes: enrollment and verification. During the enrollment process, the system learns the unique voice patterns and acoustic characteristics of a registered speaker and creates a mathematical representation known as a voice model. The verification process occurs when a speaker attempts to access a system protected by ASV. Here, the ASV system extracts features from the speaker's voice and compares it to the pre-recorded voice model. If the features match closely, the system verifies the speaker's identity; otherwise, access is denied. Various algorithms and machine learning techniques, like Convolutional Neural Networks (CNNs), are used to build and improve the accuracy of these ASV systems.

2.2 Understanding the Nature of Spoofing Attacks

Spoofing attacks pose a significant threat to ASV systems. These attacks involve an impostor attempting to mimic the voice of a genuine speaker or using sophisticated voice conversion or synthesis techniques to deceive the system. There are various types of spoofing attacks, such as replay attacks (where a previously recorded voice sample of the genuine speaker is used), voice conversion (where an impostor's voice is artificially altered to sound like the target speaker), and text-to-speech synthesis (where a synthetic voice that sounds like the target speaker is generated). The sophistication of these attacks



means they can often deceive even advanced ASV systems, leading to unauthorized access and potential security breaches.

2.3 The Significance of Incorporating Anti-Spoofing Mechanisms

The increasing complexity of spoofing attacks has underscored the importance of incorporating robust anti-spoofing mechanisms into ASV systems. Anti-spoofing mechanisms are designed to detect unusual patterns or discrepancies that could indicate a spoofing attempt. For example, they might look for signs of artificially altered or synthetic voices, or inconsistencies in ambient noise that might suggest a replay attack. By identifying these indicators anti-spoofing mechanisms can help to prevent unauthorized access, thereby enhancing the security and reliability of ASV systems. Furthermore, as spoofing techniques continue to evolve, so too must anti-spoofing measures. This necessitates ongoing to develop countermeasures that can effectively neutralize them.

3. Deep Learning in ASV Anti-Spoofing

3.1 Introduction to Convolutional Neural Networks (CNNs)

Convolutional Neural Networks (CNNs) are a class of deep learning models particularly effective in analyzing visual and audio input. CNNs are designed to automatically and adaptively learn spatial hierarchies of features directly from data. In the context of audio analysis, like ASV systems, a CNN can learn discriminative spectral patterns of a voice signal, which are robust to variations and therefore useful for speaker recognition tasks.

The CNN architecture is composed of one or more convolutional layers, followed by pooling layers, fully connected layers, and finally a classification



layer. Convolutional layers apply a series of filters to the input data to create feature maps, pooling layers reduce the spatial size of these feature maps, and fully connected layers interpret these features and classify them into various categories.

3.2. Introduction to Generative Adversarial Networks (GANs)

Generative Adversarial Networks (GANs) are a type of machine learning system introduced by Ian Goodfellow and his colleagues in 2014. GANs consist of two parts: a generator, which creates new data instances, and a discriminator, which tries to distinguish between real and synthetic data. The two components are trained together, with the generator trying to produce data that the discriminator cannot distinguish from real data, and the discriminator striving to improve its ability to differentiate real data from the fake ones. This adversarial process leads to the generator producing high-quality synthetic data.

3.3. Role and Advantages of CNNs and GANs in ASV Anti-Spoofing

In ASV anti-spoofing, CNNs and GANs play a crucial role in enhancing the system's robustness against spoofing attacks. CNNs can learn discriminative features from voice signals that are robust to various types of spoofing attacks. They can also adapt to new or unseen types of attacks by learning new representations during training. GANs, on the other hand, can be used to augment training data by generating synthetic voice samples. This can help in



improving the robustness of the system, especially in situations where the training data is limited or imbalanced.

CNNs and GANs bring several advantages to ASV anti-spoofing. They can model complex and non-linear relationships, which can be beneficial for tasks like speaker verification. They are capable of automatic feature learning, reducing the need for manual feature engineering. Lastly, due to their capability to learn directly from raw data, they can potentially outperform traditional methods that rely on handcrafted features. The effectiveness of CNNs and GANs in ASV anti-spoofing, however, largely depends on the availability of sufficient and high-quality training data.

4. Methodology

4.1. Description of the Dataset Used

The ASVspoof 2021 dataset is widely recognized as a benchmark for ASV anti-spoofing research. This dataset contains a diverse range of genuine and spoofed speech samples to mimic a variety of realistic scenarios. The genuine speech samples are recorded from several speakers in different languages while the spoofed samples are generated using a variety of voice conversion and text-to-speech synthesis techniques. The diversity of the dataset allows the ASV system to learn and generalize from a broad range of voice features making it an excellent resource for training and testing ASV and anti-spoofing systems.

4.2. Explanation of CNN-based ASV System

In our CNN-based ASV system, the primary task is to classify whether a given speech sample belongs to a claimed speaker or

not. The system uses a CNN architecture that learns a robust representation of speech features, which are used for speaker verification. The raw speech signal is first converted into a spectrogram which is used as the input to the CNN. The CNN then can learn to extract essential features from the spectrogram for the speaker verification task. The output layer of the CNN gives a score indicating the likelihood that the speech sample belongs to the claimed speaker. A threshold is then applied to this score to make the final verification decision.

4.3. Description of GAN Techniques Applied for Spoofing

In our experiment, we utilized GANs for generating spoofed speech samples. The GAN is trained on the genuine speech samples from the ASVspoof 2021 dataset, and the generator component of the GAN learns to produce synthetic speech samples that mimic the features of genuine speech. These synthetic samples are then used to test the robustness of the ASV system against spoofing attacks.

4.4. Explanation of CNN-based Anti-Spoofing Techniques

In the anti-spoofing setup, the objective is to detect whether a given speech sample is genuine or spoofed. The CNN-based anti-spoofing system similar to the ASV system uses a CNN to learn features from speech spectrograms. However, instead of learning features for speaker verification the anti-spoofing system learns features that can discriminate between genuine and spoofed speech. The output layer of the CNN gives a score indicating the likelihood that the speech sample is genuine. A threshold is then applied to this score to make the final decision about whether the sample is genuine or spoofed. The system is trained on a combination of genuine and spoofed



speech samples from the ASVspoof 2021 dataset, allowing it to learn to distinguish between the two classes effectively.

5.1.Exersise Setup

Software Specifications:

Operating System: Windows 11

Python: Python 3.8

Python libraries: Keras, TensorFlow, numpy, pandas.

IDE: Jupyter notebook, Visual Studio Code

Google Colab Pro

Colab Pro 1 V100 GPU

"Standard" RAM 13GB RAM and 2 CPUs

Software Specifications: Python.

Environment: Jupyter notebook.

5.2. Setup of CNN-based ASV System

The first step in setting up the CNN-based ASV system is to preprocess the audio data. The raw audio signals are transformed into spectrograms, which serve as the input for the CNN. Each spectrogram represents the temporal evolution of the frequency content in the audio signal, which is crucial information for speaker verification tasks.

The CNN architecture is then defined, typically including several convolutional layers, pooling layers, and fully connected layers, finally ending with a classification layer.



The CNN is trained on the training data of the ASVspoof 2021 dataset with the spectrograms of audio samples serving as inputs and the speaker identities as targets. The network learns to extract essential features from the spectrograms and map them to the corresponding speaker identities.

After the CNN is trained it can be tested on the testing portion of the ASVspoof 2021 dataset. The performance of the system can be evaluated based on metrics such as accuracy, false rejection rate, and false acceptance rate.

5.2 Setup of a Convolutional Neural Network (CNN)

Prepare the dataset: Genuine and spoofed audio samples must be preprocessed and labeled correctly. The genuine samples are labeled as '0' (not spoofed) and the spoofed samples are labeled as '1' (spoofed). The samples are then split into training and testing sets.

The CNN is designed to classify an audio sample as either genuine or spoofed.

The CNN is trained on the labeled training dataset.

The trained CNN is then tested on the testing dataset to evaluate its performance.

6. Results and Discussion

6.1 Performance of CNN-based ASV System

After trained CNN-based ASV System, we will assess the models performance not only based on its accuracy but also its ability to correctly classify genuine and spoofed instances.

Metrics:



Confusion Matrix:

Visualizing Training History:

6.2 Efficacy of GAN-based Spoofing Attacks

calculate the error rate of the ASV system when subject to GAN-generated spoofed samples. feeding the GAN-generated samples to the ASV system and computing how many of these samples are incorrectly classified as genuine.

accuracy, precision, recall, F1-score on a validation set that includes both genuine and spoofed samples. If the anti-spoofing measure is effective it should correctly classify a high proportion of both the genuine and spoofed samples.

Analyzing and interpreting results comparing the performance of the CNN-based ASV system when subjected to genuine speech GAN-based spoofed speech, and the effectiveness of CNN-based anti-spoofing measures. This comparison will provide a comprehensive understanding of the system's capabilities and weaknesses.

6.3 Effectiveness of CNN-based Anti-Spoofing Measures

The effectiveness of CNN-based anti-spoofing by measuring the accuracy, precision, recall, F1-score on a validation set that includes both genuine and spoofed samples.

achieving a balance between them is key.

6.4 Comparative Analysis and Interpretation of Results

1. Comparison of genuine and spoofed speech recognition by the ASV system
2. Evaluation of anti-spoofing measures



3. Interpretation of results

7. Conclusion

7.1 Summary of Findings