

A Survey of Audio DeepFake Detection

Deep Sengupta¹, Rahul Saha¹, Anupam Mondal (Mentor)¹

¹Department of Computer Science and Engineering(Artificial Intelligence), Institute Of Engineering and Management , University Of Engineering and Management, Saltlake sector V, Kolkata, 700091, West Bengal, India.

Contributing authors: senguptadeep03@gmail.com;
rahul.rs99.saha@gmail.com; anupam.mondal@iem.edu.in;

Abstract

Deepfake detection using machine learning and deep learning is a rapidly growing field where artificial intelligence and machine learning algorithms generate fake content. Applications of Audio Deepfakes (AD) range from audiobook enhancement to public safety threats. This article will provide a study of ways to overcome AD using a combination of machine learning (ML), deep learning (DL), and other methods. This research covers many areas of depth perception, focusing on Mel Frequency Cepstral Coefficients (MFCC) techniques and deep learning. Preliminary experiments on fake or real data demonstrate the effectiveness of support vector machine (SVM) for short words, the possibility of gradient boosting on similar data, and the performance of the VGG-16 model.

In this study, Fake or Real (FoR) dataset is used to explore features and image-based methods in addition to deep audio. Deep learning, specifically Temporal Convolutional Networks (TCN), outperforms machine learning with 92 percent accuracy. Compared to traditional CNN models such as VGG16 and XceptionNet, the proposed model shows greater accuracy in classifying sounds as falsetto or real.

They can be used to spread false information, deceive the public, or harm individuals or organizations. We conduct a comprehensive review of the existing literature, including numerical analysis, simulated and synthetic AD attacks, and quantitative comparisons of detection methods.

Keywords: Audio Deepfakes (ADs), Machine Learning (ML), Deep Learning (DL), Imitated Audio

1 Introduction

In recent years, the application of artificial intelligence in many fields, including cloning, has grown and created a lot of noise. The growth of different industries has also led to the growth of audio-fake. Nowadays, the term deepfake has become a curse and has led to the destruction of information that affects personal security. News leads to acts of violence such as Deepfakes, slander and even violence. Deepfake is a combination of using deepfake techniques to create fake content such as faces in photos, videos, or recordings. It is a type of digital content exchange in which the original face in a photo, video, or recording is replaced with a fake face. DeepFake is similar to changing the head area (i.e. the upper part) of the synthesized target so that it behaves the same as the source. Deepfake threats include the creation of revenge videos featuring the faces of victims, real videos showing national leaders compromising on false statements, stock market executives and online butchers coming face to face in a video chat. The seriousness of these risks has attracted worldwide media attention and led to two public hearings in the last two years.

1. Using machine learning for AD detection has the following advantages: • Advantages: SVM model performs well on short sounds, provides gradient enhancement on original data, while VGG-16 performs well on raw data. • Disadvantages: Search is limited to deep audio based on specific models for specific situations. 2. Siamese Architecture for Deepfake Multimedia Recognition: • Advantages: State-of-the-art techniques achieving high AUC scores on DFDC and DF-TIMIT data. • Disadvantages: Limited content on specific challenges addressed. 3. Integration of visual and auditory models: • Advantages: Global search improves performance and emphasizes the integration of visual and auditory decision. • Disadvantages: Specific description of the search function and characteristics of the data are not specified. 4. Deep Voice Search (FoR) using fake or real data: • Advantages: The proposed model is better than traditional CNN and solves the voice communication threat. • Disadvantages: Limited comparison with advanced models, lack of in-depth analysis of FoR. 5. Deep Audio Synthesis Detection Challenge (ADD) 2022: • Pros: Solve a real-life situation and show us how to compete. • Disadvantages: Lack of understanding of challenges and successes. 6. Audio Deepfakes (AD) Review: • Strengths: Provides an overview of available techniques, highlighting the need for robust AD detection. • Disadvantages: There are no specific guidelines for developing AD diagnostic criteria. 7. Evaluation of CNN Architectures for Noise Analysis: • Pros: The custom model shows realism and allows experimentation with different sounds. • Disadvantages: context dependency, suggesting that there must be many architectures

This work shows that deep learning and triple decay occur from Siamese architectures. This new approach analyzes the similarities between audiovisual and film theory for in-depth research. The proposed model surpassed the state-of-the-art method and achieved a single-video AUC score of 84.4 percent on DFDC and a best-of-video AUC score of 96.6 percent on DF -TIMIT, full audio, video, and video combined. perceptual performance dataset. View. . Unity of thought.

2 Literature Review

[1] This work uses machine learning and deep learning, specifically Mel Frequency Cepstral Coefficient (MFCC), to identify deep sounds in Fake or Real dataset. Experimental results show that Support Vector Machines (SVM) perform well on exposure to short sounds, gradient boosting performs well for old data, and the VGG-16 model performs better in other cases, especially on raw data. . [2] This work introduces a deep learning method to recognize deepfake multimedia content by analyzing audiovisual similarities and emotions in videos. The design inspired by Siamese architecture and triple loss outperforms the state-of-the-art method by achieving a good single video AUC score of 84.4 percent on DFDC and 96.6 percent on the DF-TIMIT dataset. and pioneered the sound. Video and visualization for deeper knowledge discovery. [3] This study addresses two vision and hearing threats in depth, proposing a common search operation that leads to a common combination of these models. These tests show the best performance and detail compared to the training model itself, highlighting the importance of visual and auditory judgment in deep exploration. [4] This research focuses on the use of false or true information (FoR) generated by advanced text-to-speech models to deal with the threat of voice communication. Two methods based on visuals and images were investigated for deep voice search. .. The proposed model shows greater accuracy in classifying sounds as fake or real compared to traditional CNN models such as VGG16 and XceptionNet. [5] This article introduces the Audio Deep Synthesis Detection Challenge (ADD) 2022, which addresses different real-life and complex situations for deep sound detection. This challenge includes three methods: Perfect Data Discovery (LF), Perfect Data Discovery (PF), and Perfect Data Discovery (FG). This article provides an overview of data, metrics, and methods and highlights recent advances and findings in the field of deep language search.[6]This article provides an overview of audio deepfakes (AD) and the possibilities for continued improvement of detection methods when there are concerns about their impact on public safety. It examines existing machine learning and deep learning methods, compares audio data errors, and identifies important trade-offs between accuracy and measurement methods. This review highlights the need for further research to resolve these inconsistencies and suggests potential guidelines for more robust AD detection models, particularly in addressing noise and the sound of the world.[7] This study evaluates various CNN architectures for deep sound detection, including concepts such as size, technique, and accuracy. The customized architecture of Malik et al. The most accurate is 100 percent. But it seems to depend on the context, indicating the need for different architectures. Experiments with different sound representations demonstrate the consistency of the customized model. Although these standards lag behind legal standards, they have paved the way for the creation of effective standards that will meet legal restrictions while solving deep-rooted problems.[8] This study addresses the threat of misuse of synthetic speech by arguing that the real voice is different from the synthetic voice in group discussion. The system uses deep neural networks combining negative speech, speech binarization, and CNN-based methods to achieve high accuracy and effective speech analysis. [9] This study demonstrates the application of securing transmission points using manual and automatic extraction. It uses CNN for histogram analysis and shows the performance evaluation of the

model-based application and Deep Voice speech recognition. [10] This paper solves the problem of detecting deep voice spoofs, using the ASVspoof dataset, and combining data enhancement and hybrid feature extraction. The proposed model adopts LSTM in the backend, uses MFCC + GTCC and SMOTE, and achieves 99 percent testing accuracy and 1.6 percent EER performance on ASVspoof 2021 deepfake sections. Noise evaluations and experiments on the DECRO dataset further demonstrate the effectiveness of the model. [11] Motivation: This paper demonstrates the challenges of XAI image classification by focusing on the quality score of similar objects. Recognize the difference between human and machine understanding and the impact on interpreting XAI output. [12] Feature extraction: Audio analysis by Fourier transform becomes a spectrogram, converting the audio signal into a visual form. Analysis of mel spectrograms for human speech analysis, providing visual and audio interpretation using scores based on spectrogram scores. [13] Model: Compared to CNN and LSTM models, spectrogram-based features are used for deep speech detection. Emphasis on simplicity rather than realism, using general communication techniques for speech recognition. [14] Explainable Artificial Intelligence (XAI): Introduce the XAI method based on Taylor decomposition, using the gradient integral to evaluate the accuracy of the correlation method by the integral method and the correlation redistribution to achieve the explanation. [15] Speech Generation: Describe the Griffin-Lim algorithm for generating voice from scores of spectrogram scores, for simplicity and ability to influence the factor properties of the voice even without perfect segmentation. [16] Understanding with Humans: Inspired by XAI's human assistance in visual processing, this paper explores the concept of classification of scores in complex language, absorption detection and comparison with spectrogram-based audio.

3 Dataset based Survey

[1] This work uses machine learning and deep learning for deep voice search and focuses on fake or real documents. In this study, the performance of SVM on short sounds, gradient boosting on original data, and the performance of the VGG-16 model on raw data were analyzed. The research delves into the content of fake or real data, examining in depth its composition, size, and impact on the performance of different systems and learning models. [2] This research presents a deep learning method for deep content analysis using DFDC and DF-TIMIT datasets. This research examines these materials, exploring their properties, diversity and implications for in-depth research. It evaluates Siamese architecture-inspired models in this literature that have demonstrated their effectiveness in a variety of situations. [3] To eliminate the threat of sight and sound, this research presents a collaborative research project. This study explores the dataset involved in training the joint model, examining the characteristics of the data and how the combination improves performance compared to a single model. [4] Focusing on false-or-true (FoR) data generated from text-to-speech models, this research investigates two deep speech search methods. The investigation analyzed the FoR dataset, determining the nature of false or fabricated data. It evaluates how these changes affect the performance of the proposed model and compares it with traditional CNN models. [5] Introducing the 2022 Audio Deep Synthesis Detection

Challenge (ADD), this article discusses a variety of real-life and complex detection scenarios: deep audio. This study evaluates competing materials (LF, PF, FG), investigating their properties, diversity, and implications for assessing the robustness of deep acoustic measurement models. [6] This review highlights public safety concerns by examining deep data (AD) and research techniques. This research examines current machine learning and deep learning by analyzing data used for training and testing. It explores inconsistencies in audio data and trade-offs between accuracy and measurement methods. [7] Evaluating various CNN architectures for deep sound detection, this research introduces the design of Malik et al. The research examines the data used to train and test these models, assessing how different sounds affect the consistency and performance of the models. [8] To solve the threat of speech misuse, this research uses deep neural networks. This research evaluates whether speech distortion and speech binarization lead to model accuracy by analyzing the data used to train these networks. [9] To demonstrate the use of anti-virus content, this study uses CNN for histogram analysis and Deep Voice recognition. This research explores the use of data to prevent content transfer, evaluates the impact of guidance and enables operational models to be derived. [10] To solve the problem of deep speech spoofing detection, this paper uses ASV spoofing data. The study analyzed the ASVspoof dataset to evaluate its relevance and features in training models for deep speech recognition. [11] Introducing the XAI image classification problem, this paper focuses on similar objects with good scores. The survey explores the datasets used for XAI image classification, analyzing the differences in human and machine understanding. [12] Incorporating audio analysis by Fourier transform, this study uses spectrogram-based features for deep speech detection. The survey examines the datasets employed for deep speech detection, evaluating the effectiveness of spectrogram-based features compared to CNN and LSTM models. [13] Introducing XAI based on Taylor decomposition, this paper evaluates the accuracy of the correlation method. The survey explores the datasets used for XAI based on Taylor decomposition, analyzing how the gradient integral and correlation redistribution contribute to explanation accuracy. [14] Describing the Griffin-Lim algorithm for speech generation, this paper explores voice generation from spectrogram scores. The survey scrutinizes the datasets used for speech generation, assessing the impact of the Griffin-Lim algorithm on generating voice. [15] Inspired by XAI's human assistance, this paper explores the classification of scores in complex language. The survey investigates the datasets used for classification of scores, analyzing how XAI assists in visual processing and its implications in complex language understanding. [16] This study evaluates various CNN architectures for deep sound detection, emphasizing the customized architecture of Malik et al. The research delves into the data used to train and test these models, assessing the effects of different sounds and the consistency of the design. It shows how this model provides a good solution for deep acoustic sensing.

4 Technology Based Survey

[1] Using Machine Learning and Deep Learning Using the MFCC learning process, this research analyzed the deep. SVM is good at processing short words, gradient boosting

on model datasets, and VGG-16 on raw datasets. This research uses depth of field processing by comparing SVM, gradient boosting, and VGG-16 in sound depth detection. [2] This research introduces a deep learning method that uses Siamese architecture and triplet loss to detect deep content. It outperforms the state-of-the-art method. This research explores the techniques used, examining design and performance loss and how they can help improve depth perception. [3] Regarding both face and deep face threat, this research presents a joint study together. This research uses synchronization between models to make it more efficient. This research examines techniques used for joint detection and explores how to combine visual and auditory cues to improve detection capabilities. [4] Focusing on the deep voice threat, this research uses fake or real (FoR) data created to exploit advanced speech patterns. It explores images and processes based on images, and TCN works better than machine learning. In this review, these methods are evaluated, images and image-based methods are compared, and the superiority of TCN is evaluated. [5] Introducing the 2022 Deep Sound Synthesis Detection Challenge (ADD), this research paper looks at different aspects of deep sound perception in real life. There are three ways to race: LF, PF and FG. The survey evaluates the technology used in the competition by comparing the performance and technological advances of different models. [6] This paper studies deep audio (AD) and updates the front-end detection technique. It compares ML and DL methods, compares fake data, and analyzes the trade-off between accuracy and probability. This study examines the methods used in current guidelines and explores how different methods contribute to the stability of AD diagnostic criteria. [7] This research evaluates various CNN architectures for deep sound detection, focusing on terms such as size, technology, and accuracy. The customized architecture of Malik et al. Focus on their specific facts. This research investigates the working process, examines different CNN architectures, including Malik et al. special model helps you find deep sound. Examines the role of size, technology, and accuracy in the development of deep acoustic sensing technology. [8] This study employs deep neural networks for speech analysis, combining negative speech, speech binarization, and CNN-based methods. This survey explores the technologies used, analyzing the architecture and techniques employed by deep neural networks for accurate and effective synthetic speech analysis. [9] Utilizing CNN for histogram analysis and model-based application, this study secures transmission points and performs Deep Voice speech recognition. This survey delves into the technologies employed, analyzing the role of CNN in histogram analysis and the techniques used in the model-based application for securing transmission points and speech recognition. [10] Addressing deep voice spoofs, this paper combines data enhancement and hybrid feature extraction, adopting LSTM in the backend. Measurement accuracy is high when MFCC+GTCC and SMOTE are used. This research investigates the techniques used in deep speech recognition by examining the role of modeling and extraction techniques such as LSTM, MFCC, GTCC and SMOTE. [11] Due to the difficulty faced in supporting XAI image distribution, this article will focus on the quality score of similar products. It acknowledges the difference between human and machine understanding by exploring the implications of interpreting XAI output. This research examines the technologies used and evaluates how different technologies can help solve problems in XAI image distribution. [12] In audio analysis, this article uses the Fourier

transform to create spectrograms and convert audio signals into visual form. It provides visual and audio descriptions using spectrogram scores as scores. This research explores the techniques used by examining how Fourier transforms and spectrogram-based techniques contribute to sound extraction. [13]Proposed model for deep speech detection using spectrogram-based features for simplicity rather than realism. This research compares CNN and LSTM models, examines their processing methods, and investigates the role of CNN and LSTM in spectrogram-based deep speech recognition. [14]Introducing Explainable Artificial Intelligence (XAI) using Taylor decomposition, this article evaluates the accuracy of the relevant methods. It uses gradient integration for evaluation and correlation analysis to achieve translation. This research investigates the techniques used in the XAI model for image classification by examining the role of Taylor decomposition, gradient integration and correlation redistribution. [15]Griffin-Lim algorithm identifies speech generated by spectrogram score; This article focuses on simplicity and the ability to influence expressive speech without perfect segmentation. This research delves into the process used, examining the role of the Griffin-Lim algorithm, spectrogram scores, and their impact on speech. [16]Inspired by XAI's aid in human visualization, this article investigates the classification of fractions in complex languages, detects and compares it to spectrogram-based audio. This research examines the strategies included in this study, examining how different methods contribute to the scoring, detection, and scoring of difficult words compared to spectrogram-based audio.

5 Evolutionary Mechanisms Survey

[1]This research uses machine learning and deep learning using MFCC to determine the depth of misinformation or true information. This research traces the evolution of detection mechanisms and examines how SVM, gradient boosting, and VGG-16 have evolved in the application of audio deepfake detection.[2]Introducing the Deepfake learning method, this research uses Siamese architecture and triplet loss to detect Deepfake multimedia content. This research explores the evolution of detection tools by examining how the use of Siamese architecture and triplet loss represent advances in audiovisual deepfake detection. [3]Regarding both facial and deep facial threat, this research presents a joint study together. This research examines the evolution of search mechanisms by examining how synchronization of vision and hearing represents an advance in deep search. [4]Focusing on the threat of deep voice, this research uses fake or real (FoR) data generated by high speech standards. This study explores the evolution of search mechanisms by comparing the evolution of feature-based and image-based methods and the emergence of TCN as a superior technology. [5]Introduction to Sound Depth Comprehensive Detection Competition (ADD) 2022, this article really addresses the difference of sound depth perception in real life. This research examines the evolution of detection mechanisms and examines how challenges drive innovation in noise detection techniques. [6]This article examines deep voice (AD) and the possibility of changes in perception. This research explores the evolution of search mechanisms, examines how current machine learning and deep learning methods have evolved, and identifies ongoing issues and patches in AD research.[7]

This study evaluates various CNN architectures for deep sound detection, taking into account aspects such as size, technology and accuracy. The customized architecture of Malik et al. Achieving 100 percent accuracy shows the importance of this. This research traces the evolution of detection mechanisms and evaluates how experiments with different sounds and CNN architectures (especially the proposed ones) have led to advances in sound perception. It shows that different design models are needed depending on the context, and these models can lead to effective solutions that meet legal requirements in solving problems.[8] This research addresses the threat of using illegal language in group discussions. Use deep neural networks, fuzzy communication, speech binarization, and CNN-based methods for accurate speech measurement. This study traces the evolution of the detection process and examines how the combination of fuzzy speech and CNN-based methods represents an advance in speech analysis[9].To demonstrate the implementation of secure transmission, this study uses CNN for histogram analysis and Deep Voice recognition. This research examines the evolution of detection mechanisms, examining how CNN-based methods and speech recognition can help prevent content contamination and improve overall performance. [10] Using ASVspoof data, this article provides advanced and integrated data. The proposed model uses LSTM in the background to perform pressure measurement. This research explores the evolution of detection techniques by examining how the combination of data augmentation, hybrid feature extraction, and LSTM represent advances in speech perception. [11]Motivated by competition This paper focuses on the quality score of similar products in XAI image distribution. Recognizing the difference between human and machine understanding, this study examines the development of search tools and examines how a focus on quality scores can contribute to progress in solving XAI image classification problems. [12] In audio analysis, this paper uses the Fourier transform to create human-interpretable spectrograms. This work traces the evolution of detection techniques by examining the use of Fourier transforms and feature-based spectrograms to represent advances in voice feature extraction for dynamic speech.[13]A proposed model for deep speech detection using spectrogram-based features, simplicity over realism. This research examines the evolution of the exploration process by comparing CNN and LSTM models, examining the importance of simplicity and the use of communication technologies for advances in spectrogram-based deep-sea exploration. [14] Introducing Explainable Artificial Intelligence (XAI) using Taylor decomposition, this article evaluates the validity of other methods. It uses gradient integration for evaluation and correlation analysis to achieve translation. This research traces the evolution of the search process by examining how the introduction of Taylor decomposition, gradient integration and redistribution relations represent progress in the XAI model for image classification. [15] The book describes the Griffin-Lim algorithm for generating speech from spectrogram scores; This article focuses on the simple and available ability to intercept speech without perfect segmentation. This research examines the evolution of the perception process by examining how the Griffin-Lim algorithm and the use of spectrogram scores represent advances in speech. [16] Inspired by the XAI program visualization function, this paper investigates the segmentation, detection and matching of complex words in comparison with spectrogram-based language. This study evaluates the evolution of the search process,

examining how the focus compares with the advancement of spectrogram-based voice support in complex word classification, detection, and understanding in the human voice and comprehension in general.

Discussion

[1]The first article uses machine learning and deep learning, specifically MFCC, to identify deepfakes. SVM is good at processing short sounds, gradient boosting on artifacts, and VGG-16 on artifacts. This approach shows that there are many types of models with different characteristics.[2]Providing an in-depth study inspired by Siamese architecture, the second article leads an integrated analysis and aims to explore the depth of similar sound and emotion. AUC scores are impressive on DFDC and DF-TIMIT data, highlighting the importance of combining hypotheses in multivariate analysis.[3]Concerning the sight and sound of threats, the third article recommends that joint search operations be more effective. Synchronization of vision and hearing enhances overall exploration and highlights the importance of considering both in deep exploration.[4]The fourth paper focuses on the FoR dataset and compares the visual performance of deep speech and image detection over CNN models. This approach demonstrates the value of using text-to-speech models to generate inaccurate or accurate information.[5]The fifth article introduces the 2022 ADD Competition solving real-life problems. The inclusion of LF, PF and FG pieces in the competition emphasizes the need for quality models in different sizes and deep tones.[6]The sixth article provides an overview of deep voice, highlights gaps in current systems, and solicits recommendations for research that finds AD to be more robust. This comprehensive review underscores the need for further research to address the challenges posed by real noise.[7]The seventh article evaluates various CNN architectures for deep sound detection; This shows that different models with different backgrounds are needed in deep voice search. The customized architecture of Malik et al. High sensitivity is found, paving the way for a good model for deep voice search.[8]To solve the problem of synthetic speech abuse, the eighth paper adopted the deep neural networks method, which includes negative communication, speech binarization, and CNN. This approach emphasizes the importance of separating speech from speech synthesis, allowing for accurate and effective speech analysis.[9]The ninth paper demonstrating the use of context immunity uses CNN for histogram analysis and deep speech learning. This approach demonstrates the versatility of CNNs in predicting the spread of content.[10]To solve the problem of deepfake detection, the tenth paper is applied to the ASVspoof dataset. Using LSTM, MFCC + GTCC and SMOTE, the proposed model completes the accuracy test and demonstrates robustness against deep speech.[11-16]The rest of the article explores the problems of the XAI method such as XAI image classification, feature extraction using Fourier transform, model comparison of deep speech, Taylor decomposition using Griffin-Lim calculus.

Conclusion and Future Scope

In summary, this research supports false or true (FoR) data and provides a powerful deep language search method. Feature engineering using Mel Cepstral Coefficients

(MFCC) has proven useful and the effectiveness of machine learning algorithms has been demonstrated. Support vector machines (SVM) and gradient boosting superiority show very good accuracy. Comparison of deep sound detection performance demonstrates the superiority of Temporal Convolutional Networks (TCN) compared to traditional CNN models. While future directions include investigating different MFCC window sizes in the extraction process, measuring the structure in real conditions remains a major challenge.

1. Improved Feature Extraction Techniques: Current research mostly uses MFCC and spectrogram-based features. Future research may explore optimal extraction techniques such as wavelet transform or hybrid methods to better preserve patterns in deep sounds.

2. Integration of different methods: Although current research focuses on deep sound perception, the integration of visual perception can be improved throughout reality. Future studies could explore the combination of audio and video using this combination to create a more powerful search engine.

3. Adversarial Robustness: Adversarial attacks pose a threat to deep search models. Future work should focus on improving the robustness of the model against counterattacks to ensure performance in real-world situations.

4. Real world testing: Testing the model in the real world, for example manipulating ambient noise, reverberation and different recording equipment, is crucial to clearly ensure the validity of deep exploration. system. This involves testing the model in different settings to test its effectiveness and generalizability.

5. Continuous dataset development: Continuous efforts should be directed towards creating diverse and complex datasets. These datasets should reflect real-world scenarios, encompassing a wide range of accents, languages, and environmental conditions to improve the robustness and generalization of models.

6. Examination of Ethical Implications: As deepfake technology evolves, ethical considerations become increasingly important. Future research should delve into the ethical implications of deepfake detection, including issues related to privacy, consent, and responsible use of detection systems.

7. Exploration of Explainability Techniques: Given the complexity of deep learning models, developing explainability techniques is crucial for gaining insights into model decisions. Future work should explore explainable AI methods, ensuring transparency and interpretability in deepfake detection systems.

8. Standardization and Benchmarking: Creating benchmarks and benchmarks for deep learning systems will facilitate fair comparison of different methods. This involves defining common criteria to measure the effectiveness of the model.

9. Human Systems in the Loop: Involving human intelligence in the review process can improve performance. Future studies may explore human-machine circulatory systems where the combination of intelligence and human judgment helps produce accurate and reliable results.

10. Continuous collaboration and knowledge sharing: Collaboration between researchers, industry experts and policymakers is crucial to staying ahead of the curve. Establishing platforms for continuous information sharing and collaboration can lead to a more unified and effective response to the challenges posed by deepfake technology.

References

- [1] A. Hamza et al., "Deepfake Audio Detection via MFCC Features Using Machine Learning," *IEEE Access*, vol. 10, pp. 134018-134028, 2022, doi: 10.1109/ACCESS.2022.3231480.
- [2] T. Mittal et al., "Emotions Don't Lie: An Audio-Visual Deepfake Detection Method using Affective Cues," in *Proceedings of the 28th ACM International Conference on Multimedia (MM '20)*, ACM, New York, NY, USA, 2020, pp. 2823-2832, doi: 10.1145/3394171.3413570.
- [3] Y. Zhou, S.-N. Lim, "Joint Audio-Visual Deepfake Detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 14800-14809.
- [4] J. Khochare et al., "A Deep Learning Framework for Audio Deepfake Detection," in *Arabian Journal for Science and Engineering*, vol. 47, no. 3, 2022, pp. 3447, doi: [journal DOI].
- [5] J. Yi et al., "ADD 2022: the first Audio Deep Synthesis Detection Challenge," in *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Singapore, 2022, pp. 9216-9220, doi: 10.1109/ICASSP43922.2022.9746939.
- [6] "Challenges A Review of Modern Audio Deepfake Detection Methods and Future Directions," <https://doi.org/10.3390/a15050155>.
- [7] M. Mcuba et al., "The Effect of Deep Learning Methods on Deepfake Audio Detection for Digital Investigation," in *Procedia Computer Science*, vol. 219, 2023, pp. 211-219, doi: 10.1016/j.procs.2023.01.283.
- [8] R.L.M.A.P.C.Wijethunga et al., "Deepfake Audio Detection: A Deep Learning Based Solution for Group Conversations," in *2020 2nd International Conference on Advancements in Computing (ICAC)*, Malabe, Sri Lanka, 2020, pp. 192-197, doi: 10.1109/ICAC51239.2020.9357161.
- [9] D.M. Ballesteros et al., "Deep4SNet: deep learning for fake speech classification," in *Expert Systems with Applications*, vol. 184, 2021, 115465, doi: 10.1016/j.eswa.2021.115465.
- [10] N. Chakravarty and M. Dua, "Data augmentation and hybrid feature amalgamation to detect audio deep fake attacks," in *Physica Scripta*, vol. 98, no. 9, 2023.
- [11] S.-Y. Lim et al., "Detecting Deepfake Voice Using Explainable Deep Learning Techniques."
- [12] "A Review of Deep Learning Based Speech Synthesis," in *Appl. Sci.*, 2019, 9, 4050, doi: [DOI].
- [13] Y. Ren et al., "FastSpeech 2: Fast and High-Quality End-to-End Text to Speech," arXiv, 2020, arXiv:2006.04558.
- [14] J. Shen et al., "Natural Tts Synthesis by Conditioning Wavenet on Mel Spectrogram Predictions," in *IEEE*, 2018, pp. 4779-4783.
- [15] W. Ping et al., "Deep voice 3: Scaling text-to-speech with convolutional sequence learning," arXiv, 2017, arXiv:1710.07654.
- [16] Z. Khanjani et al., "How deep are the fakes? Focusing on audio deepfake: A survey," arXiv, 2021, arXiv:2111.14203.