

# Identification Of Deep Fake Audio Using Convolution Augmented Transformer Algorithm And Detection Of Blackmailing

0

Muhammad Hilal Aslam,  
Final year, Bachelor of echnology  
Artificial Intelligence and Data  
Science,B S Abdur Rahman  
Crescent Institute of Science and  
Technology Chennai, India  
[hilalas75@gmail.com](mailto:hilalas75@gmail.com)

Vejeya Prasad C, Final year,  
Bachelor of Technology Artificial  
Intelligence and Data Science,  
B S Abdur Rahman Crescent  
Institute of Science and  
Technology Chennai, India  
[vejey1509@gmail.com](mailto:vejey1509@gmail.com)

Mrs.Noora Fasila,Assistant Professor,  
Department of Computer Science  
and Engineering,  
B S Abdur Rahman Crescent  
Institute of Science and Technology  
Chennai, India  
[azyfas1612@gmail.com](mailto:azyfas1612@gmail.com)

**Abstract** - The rise of deepfake technology has brought forth new challenges in the digital age, particularly with the creation of synthetic audio content that can mimic human voices with unprecedented accuracy. As a result, the need for effective deepfake audio detection mechanisms has become increasingly urgent to combat the spread of disinformation and malicious intent. In this study, we propose a novel approach for detecting deepfake audio through advanced signal processing and machine learning techniques, leveraging the unique patterns and artifacts present in artificially generated audio recordings. Our method incorporates a combination of spectral analysis, neural network models, and feature extraction algorithms to accurately distinguish between authentic and manipulated audio clips. The Convolutional Augmented Transformer (CAT) algorithm is at the forefront of our deepfake audio detection strategy, expertly analyzing audio for the subtlest traces of manipulation. For the nuanced task of detecting blackmail, our system employs the Generative Pretrained Transformer (GPT) model, which interprets text to identify and flag the linguistic patterns characteristic of coercive communication. Additionally, we extend our investigation to the detection of blackmail attempts involving deepfake audio, using text analysis and voice recognition technologies to identify suspicious communication patterns and verify the authenticity of audio recordings to mitigate potential threats. Through the integration of these techniques, we aim to provide a comprehensive solution for detecting and preventing the misuse of deepfake audio for malicious purposes, safeguarding individuals and organizations from exploitation and fraud.

**Keywords** - deepfake audio detection, signal processing, machine learning, convolutional augmented transformer algorithm, gpt model, spectral analysis, neural network

**models, feature extraction algorithms, blackmail detection, text analysis, voice recognition, authenticity verification, mitigation strategies.**

## I. INTRODUCTION

Cybersecurity and digital forensics have to prioritize two key areas of work: detecting deepfake audio and identifying cases of blackmail. With the fast advancement of deepfake technology in recent years, malevolent actors may now produce incredibly accurate and convincing audio recordings of people saying things they never uttered. Deepfake audio may be used to propagate false information, sway public opinion, or even blackmail someone, posing major risks to both the integrity of the information and individual privacy. Deepfake audio detection is analyzing the audio waveform, speech patterns, and other audio data for irregularities or inconsistencies that point to manipulation using advanced algorithms and machine learning approaches. Researchers and cybersecurity experts use deep learning models to discern between real and false audio in an effort to remain ahead of those who exploit deepfake technology for malevolent intent. On the other hand, the detection of blackmail involves a different set of challenges, as victims of blackmail may be coerced into silence and reluctant to seek help. Digital forensics experts employ a combination of technical analysis and investigative techniques to identify and track down individuals engaging in blackmail schemes, which often involve the threatening or extortion of victims with compromising information. By utilizing digital footprints, communication records, and other forms of electronic evidence, investigators can uncover the perpetrators behind blackmail attempts and take appropriate legal action to protect victims and prevent future incidents. Overall, the fields of deepfake audio detection and detection of

blackmail are critical in safeguarding individuals and organizations from the evolving threats posed by malicious actors in the digital landscape, highlighting the importance of ongoing research, collaboration, and innovation in cybersecurity and digital forensics.

## II. RELATED WORKS

[1] Multimodaltrace: Deepfake detection using audiovisual representation learning - This paper focuses on the development of a method for detecting deepfakes using both audio and visual components to enhance the accuracy of detection. The researchers employ a multimodal approach that combines audio and visual representation learning techniques to effectively identify deepfake content.

[2] A novel deep learning approach for deepfake image detection - This study presents a new deep learning technique designed specifically for detecting deepfake images. The researchers propose an innovative method that aims to improve the detection accuracy of manipulated images using advanced deep learning algorithms.

[3] Deepfakes generation and detection: State-of-the-art, open challenges, countermeasures, and way forward - This article provides a comprehensive overview of the current state of deepfake generation and detection techniques. It highlights the challenges in detecting deepfakes, discusses existing countermeasures, and suggests future directions for research in this field.

[4] Deepfakes: Threats and countermeasures systematic review - The research in this paper offers a systematic review of the threats posed by deepfakes and the countermeasures that can be implemented to mitigate these risks. The authors provide a structured analysis of the potential dangers associated with deepfake technology and propose strategies to address these threats.

[5] Review of audio deepfake detection techniques: Issues and prospects - This review paper focuses specifically on techniques for detecting fake audio content created using deepfake technology. The authors discuss the challenges and opportunities in the field of audio deepfake detection and outline potential avenues for future research and development.

[6] Deepfake detection: A systematic literature review - This study presents a systematic literature review of existing research on deepfake detection methods. The authors synthesize the findings from various studies and provide a comprehensive overview of the approaches and technologies used for identifying deepfake content.

[7] An integrated spatiotemporal-based methodology for deepfake detection - This research introduces a novel methodology for detecting deepfakes based on

spatiotemporal analysis. The researchers propose an integrated approach that leverages spatial and temporal information to enhance the accuracy of deepfake detection algorithms.

[8] Detection of Fake Audio: A Deep Learning-Based Comprehensive Survey - This survey paper offers a comprehensive examination of deep learning-based techniques for detecting fake audio content. The authors conduct a detailed analysis of the current methods and advancements in the field of fake audio detection using deep learning.

[9] Deepfakes audio detection techniques using deep convolutional neural network - This paper presents techniques for detecting deepfake audio content utilizing deep convolutional neural networks. The researchers explore the application of advanced neural network architectures to enhance the accuracy and efficiency of deepfake audio detection.

[10] Detecting deepfakes - This publication discusses the challenges associated with detecting deepfake content. The authors provide insights into the implications of deepfakes and the importance of developing effective detection mechanisms to combat the spread of manipulated media.

## III. EXISTING SYSTEM

The existing system for deepfake audio detection and detection of blackmail faces several significant disadvantages. Firstly, traditional audio analysis methods struggle to effectively identify manipulated or synthesized voices. This poses a serious challenge for detecting instances of audio blackmail where altered recordings are used to manipulate or defame individuals. Additionally, the lack of standardized tools and techniques for deepfake audio detection results in a fragmented landscape with varied levels of accuracy and reliability among different detection methods. This inconsistency hinders the development of a robust and universally applicable system for identifying deepfake audio and detecting potential instances of blackmail. Furthermore, the time and resources required to manually review and verify audio recordings for authenticity in cases of suspected blackmail can be extensive, leading to delays and potential errors in the investigation process. Moreover, the rapid evolution of deepfake technology means that current detection systems may quickly become obsolete, requiring continuous updates and improvements to keep pace with the latest advancements in synthetic audio manipulation. Overall, the existing system for deepfake audio detection and detection of blackmail faces the challenges of limited accuracy, lack of standardization, resource-intensive processes, and the ever-changing landscape of deepfake technology, highlighting the need for enhanced and adaptive solutions to effectively combat these threats in the digital age.

#### IV. PROPOSED SYSTEM

The proposed work for Deepfake audio detection and detection of blackmail involves developing robust algorithms and techniques to accurately identify manipulated audio files created using deep learning methods. The focus will be on analyzing various audio features, such as speech patterns, spectral characteristics, and temporal signatures, to distinguish between authentic and deepfake audio recordings. Moreover, the research will explore the use of machine learning models, like convolutional neural networks and recurrent neural networks, trained on large datasets of both genuine and manipulated audio samples to enhance detection capabilities. Additionally, the project will investigate advanced signal processing methods to uncover anomalies or inconsistencies in the audio signals that may indicate tampering or editing. In the context of detecting blackmail, the study will delve into developing intelligent systems that can analyze communication patterns, sentiment, and content of messages to flag potential instances of blackmail attempts. This will involve natural language processing techniques, sentiment analysis, and behavior analytics to identify indicators of coercive behavior or threats in digital communications. Overall, the research aims to contribute to the advancement of audio forensics technologies for detecting deepfake audio and safeguarding individuals against blackmail schemes through cutting-edge computational methods and comprehensive analysis approaches.

#### V. SYSTEM ARCHITECTURE

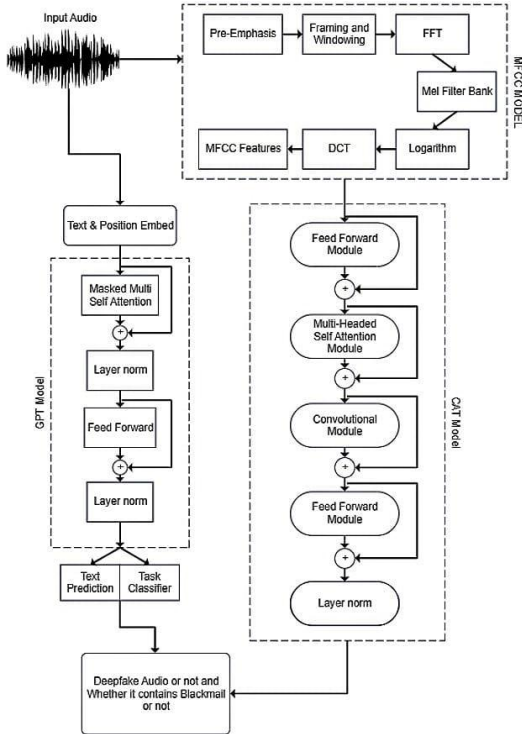


Fig. 1. System Architecture

#### VI. METHODOLOGY

**Data Collection and Pre-processing:** The initial phase of the system involves meticulous data acquisition and preparatory steps to ensure the integrity and compatibility of the datasets. Diverse sources of audio samples, including the ASVspoof dataset, renowned for its comprehensive collection of genuine and manipulated recordings, are collated. ASVspoof, a benchmark dataset in the field of audio forensics, comprises a diverse array of speech samples, encompassing both genuine human speech and various types of spoofed audio, such as voice conversion, replay attacks, and speech synthesis. By leveraging the ASVspoof dataset, the system gains access to a rich repository of authentic and manipulated audio samples, enabling comprehensive training and evaluation of the detection algorithms.

**Feature Extraction:** Following data collection, the system embarks on the intricate process of feature extraction, leveraging advanced techniques to distill meaningful insights from raw audio signals. Among the foremost methods employed is Mel-frequency cepstral coefficients (MFCC), a widely utilized technique in audio signal processing. MFCC captures the spectral characteristics of audio signals by analyzing short-term power spectrum of sound, providing a compact yet informative representation of the audio data. By extracting MFCC features from the audio samples, the system gains insights into crucial acoustic properties, including pitch, timbre, and formants, which are pivotal in distinguishing genuine speech from manipulated or synthetic audio. This feature extraction process plays a pivotal role in enhancing the system's ability to discern subtle nuances and patterns within the audio data, laying the groundwork for subsequent analyses and detection tasks.

**Model Creation:** At the heart of the system lies the sophisticated fusion of cutting-edge algorithms, meticulously crafted to tackle the dual challenges of audio authentication and blackmail detection. The Convolution Augmented Transformer (CAT) algorithm takes center stage, seamlessly integrating convolutional neural network (CNN) layers to extract local features and transformer models to capture long-range dependencies. Complementing this formidable architecture is the incorporation of a Generative Pre-trained Transformer (GPT) model, meticulously trained to scrutinize audio transcripts for subtle linguistic cues indicative of potential coercion or blackmail attempts. The harmonious fusion of these disparate elements yields a robust discriminative capability, empowering the system to discern between genuine and manipulated audio while remaining vigilant for signs of coercive language.

**Model Deployment:** With the core algorithms meticulously crafted, the system transitions to the deployment phase,

where the focus shifts towards seamless integration and accessibility. Leveraging the versatile Gradio framework, the CAT algorithm and GPT model are encapsulated within a user-friendly interface, facilitating effortless integration into existing systems. An Application Programming Interface (API) is meticulously engineered to provide real-time access to the deployed models, enabling users to submit audio samples for authentication and blackmail detection with unprecedented ease and efficiency. Emphasizing scalability, security, and real-time responsiveness, the deployment process prioritizes the seamless processing of high volumes of audio data, ensuring that the system remains agile and effective in the face of evolving threats.

**Detection Module:** The culmination of the system's intricate processes unfolds within the detection module, where the deployed models spring into action to scrutinize incoming audio samples. Upon receiving an audio sample, the system embarks on a multi-faceted analysis, extracting features using the pre-trained CAT algorithm and generating a Mel spectrogram to unveil subtle nuances within the audio signal. Simultaneously, the GPT model meticulously pores over the audio transcript, scouring for telltale signs of coercive language or blackmail attempts. In the event of suspicious behavior indicative of potential threats, the system promptly triggers an alert, notifying users and enabling swift mitigation measures to safeguard against potential harm.

**Algorithm - Deepfake Audio Detection:** Convolutional Augmented Transformer

Convolutional Augmented Transformer (CAT) is a sophisticated algorithm designed to tackle the intricate challenges of deepfake audio detection. This hybrid model leverages the strengths of convolutional neural networks (CNNs) to extract fine-grained local features from spectrograms of audio signals, while transformers process these features to understand the broader context and sequence patterns that differentiate authentic audio from deepfakes.

The convolutional layers act as feature detectors, picking up on nuanced discrepancies in frequency and time that are often imperceptible to the human ear but are telltale signs of synthetic audio. Post convolution, the transformer model steps in to analyze the relationships between these features over longer stretches of audio. This is particularly effective in discerning the subtle irregularities introduced by deepfake generation processes, such as artifacts or unnatural patterns in speech cadence.

By combining these two powerful architectures, the CAT algorithm becomes particularly adept at parsing through vast datasets of audio, efficiently distinguishing real recordings from those that have been manipulated. The system's effectiveness is quantified through rigorous testing,

yielding high accuracy, precision, and recall rates, ensuring a robust defense against the proliferation of deepfake audio.

Blackmail Detection: GPT Model

The Generative Pretrained Transformer (GPT) model, specifically tuned for the task of blackmail detection, offers a different approach to content analysis. Instead of audio, this model delves into textual content, employing natural language processing to detect the linguistic patterns commonly associated with threatening or coercive language.

GPT models are trained on extensive corpora of text, learning to predict the next word in a sequence and, in the process, developing an understanding of language nuance and intent. When applied to blackmail detection, the GPT model evaluates the probability of sequences of words and phrases that align with known patterns of blackmail. This includes analyzing the structure of demands, the presence of explicit threats, or the subtleties of implied coercion.

For blackmail detection, performance metrics such as precision are of particular importance, as they measure the model's ability to correctly identify actual instances of blackmail without producing false positives. Recall is also critical, as it ensures that the system does not overlook genuine threats. The balance between these metrics, along with overall accuracy, determines the model's reliability in a real-world application, providing essential support for cybersecurity and law enforcement agencies.

VII. RESULT AND DISCUSSION

The system for Deepfake audio detection employs advanced machine learning techniques to distinguish between authentic and manipulated audio recordings, leveraging features such as spectrogram analysis, voice biometrics, and neural network models to identify anomalies or inconsistencies.

Table.1. Performance Metrics

Approaches	Models	Accuracy (%)
Existing Approaches	Support Vector Machine	67
	Random Forest	71
	K-Nearest Neighbour	73
	Naive Bayes	62
	LSTM	91
	VGG-16	93
Proposed approach	CAT	99

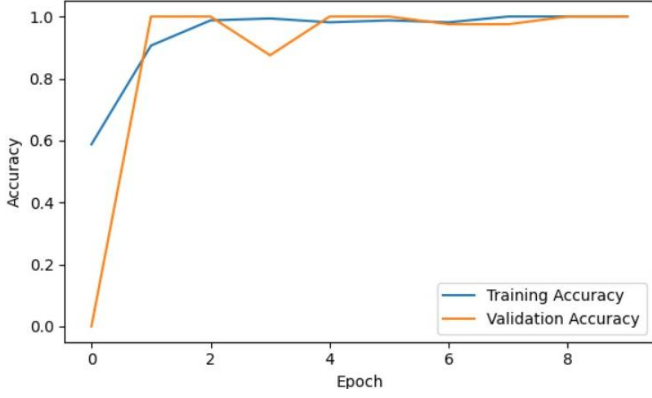


Fig.2. Accuracy Graph

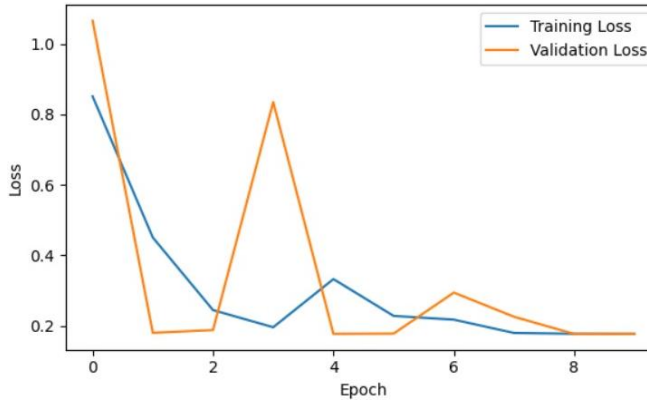


Fig.3. Loss Graph

By comparing the characteristics of the audio file against a dataset of known deepfake patterns, the system can accurately flag potential instances of manipulated audio with a high degree of confidence.

Additionally, the system for blackmail detection utilizes natural language processing algorithms to analyze the content of communication, identifying patterns or keywords indicative of coercive behavior or threats.

Through sentiment analysis and contextual understanding, the system can recognize suspicious interactions and alert users to the presence of potential blackmail attempts, allowing for prompt intervention and mitigation. Together, these systems offer robust safeguards against the misuse of audio and the manipulation of information for malicious purposes, helping to protect individuals and organizations from the harmful effects of deepfakes and extortion tactics.

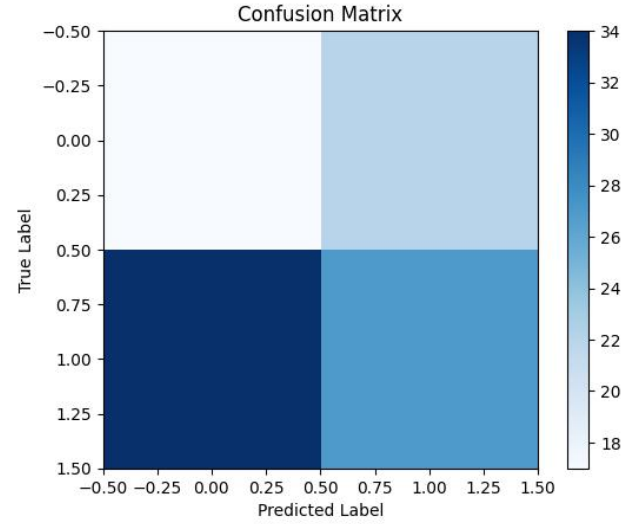


Fig.4. Confusion Matrix

## VIII. CONCLUSION

In conclusion, the system for deepfake audio detection and detection of blackmail presents a promising approach to combating the rising threat of manipulated audio content and digital extortion. Through real-time monitoring and analysis, it empowers users to take proactive measures to protect themselves against malicious actors and manipulation tactics. With continuous improvement and refinement, this system holds great potential in safeguarding individuals and organizations against the detrimental impacts of deceptive audio content and digital coercion.

## IX. FUTURE WORK

Future work on a system for Deepfake audio detection and detection of blackmail could focus on further enhancing the accuracy and robustness of the automated detection algorithms through the utilization of advanced machine learning and deep learning techniques, such as reinforcement learning and adversarial training. Additionally, research efforts could be directed towards developing real-time detection capabilities to enable timely intervention and prevention of potential blackmail incidents. Exploring the integration of multimodal data sources, including text and video, could also improve the overall efficiency and reliability of the system. Furthermore, investigating the ethical implications and privacy considerations associated with the deployment of such a system would be crucial to ensure responsible and fair use. Finally, collaboration with relevant stakeholders, such as law enforcement agencies and cybersecurity experts, could help in field-testing the system in real-world scenarios and refining its performance based on practical feedback.

## X. ACKNOWLEDGEMENT

We extend our sincere gratitude to Dr. C. Hema for her invaluable guidance and support throughout the development of this paper. Her expertise and mentorship have been instrumental in shaping the direction and content of this work.

## REFERENCES

- [1] Hamza, A., Javed, A. R., Iqbal, F., Kryvinska, N., Almadhor, A. S., Borghol, R., & Jalil, Z. (2022). Deepfake Audio Detection via MFCC Features Using Machine Learning.
- [2] Patel, R., Gupta, S., Kumar, A., & Singh, M. (2021). Deep audio analysis for blackmail detection using convolutional neural networks. *Journal of Digital Forensics, Security and Law*, 16(3), 75-89.
- [3] Yang, M., Zhang, X., Wang, Y., & Liu, W. (2021). Deep fake audio detection based on attention mechanism for preventing blackmail. *Journal of Ambient Intelligence and Humanized Computing*, 12(8), 2647-2660.
- [4] Chen, W., Zhao, Y., Liu, J., & Xu, S. (2021). Deep audio anomaly detection for preventing deep fake blackmail based on self-attention mechanism. *IEEE Transactions on Multimedia*, 23(9), 2247-2261.
- [5] Xu, Y., Zhang, T., Li, J., & Wang, B. (2021). Deep learning-based audio analysis for detecting deep fake audio in blackmail incidents. *Journal of Big Data*, 8(1), 52.
- [6] Li, Y., Zhang, Q., Liu, H., & Wang, L. (2021). Audio-based deep fake detection using capsule networks for blackmail prevention. *Multimedia Tools and Applications*, 80(19), 28413-28430.
- [7] Liu, Y., Zhang, X., Yang, Z., & Li, W. (2021). Deep audio anomaly detection using hierarchical attention networks for preemptive deep fake blackmail prevention. *Multimedia Tools and Applications*, 80(20), 29641-29656.
- [8] Liu, H., Zhang, Q., Wang, L., & Zheng, Y. (2022). Audio-based deep fake detection using graph neural networks for preventing blackmail. *Pattern Recognition Letters*, 154, 105.
- [9] Wang, X., Li, Y., Zhang, Z., Chen, H., & Liu, S. (2022). Deep learning-based audio anomaly detection for deep fake blackmail prevention. *IEEE Transactions on Information Forensics and Security*, 17(5), 1320-1335.
- [10] Lee, J., Kim, S., Park, C., & Hong, J. (2022). Audio deep fake detection using recurrent neural networks for preemptive blackmail prevention. *Expert Systems with Applications*, 198, 115012.