

An AI-Powered Deepfake Detection System for Video and Audio

Mrs. Malarvizhi AP/SC¹, Chandru S², Aravindan N³, Balachandru K⁴, Anto Ashik D⁵

^{1,2,3,4,5} Department of Computer Science and Engineering, Surya Group of Institutions.

doi.org/10.64643/IJIRTV12I11-201045-459

Abstract—The rapid advancement of generative artificial intelligence has led to the proliferation of hyper-realistic synthetic media, commonly known as deepfakes. These manipulations pose significant threats to information integrity, cybersecurity, and social trust. This paper presents a novel AI-powered detection system capable of identifying forged content in both video and audio modalities. The proposed system integrates a dual-branch convolutional neural network with temporal attention mechanisms for video analysis and a spectrogram-based residual network for audio forgery detection. A multimodal fusion layer combines features from both streams to improve overall detection accuracy. The system ensures data integrity through cryptographic hashing and provides tamper-proof audit logging using blockchain technology. Experimental results on the FaceForensics++ and ASVspoof2019 datasets demonstrate a detection accuracy of 96.8% for video, 94.2% for audio, and 97.4% when using joint inference. The system also shows robustness against compression and noise artifacts. This work provides a practical, real-time solution for deepfake mitigation in digital media forensics.

Index Terms—Deepfake Detection, Convolutional Neural Networks, Multimodal Fusion, Audio Forensics, Video Forensics, Generative AI Security, Blockchain, Temporal Attention.

I. INTRODUCTION

A. Background

Deepfake technology has emerged as one of the most concerning developments in the field of artificial intelligence. Generative Adversarial Networks (GANs), variational autoencoders, and diffusion-based models have enabled the creation of synthetic media that is nearly indistinguishable from authentic recordings. These deepfakes have been maliciously used for political disinformation, financial fraud, identity impersonation, and the fabrication of evidence

in legal proceedings. The democratization of generative AI tools has lowered the technical barrier required to produce convincing forgeries, making deepfake creation accessible to individuals with minimal technical expertise.

Cloud computing and social media platforms have become primary channels for deepfake distribution. According to industry reports, the number of deepfake videos online has doubled approximately every six months over the past three years. Social media platforms, news organizations, and law enforcement agencies have reported increasing difficulty in distinguishing authentic content from sophisticated manipulations. Traditional forensic methods, such as metadata analysis, error level analysis, and noise pattern examination, are increasingly ineffective against modern neural network-generated forgeries that preserve statistical properties of authentic media.

B. Problem Statement

The primary challenge in deepfake detection lies in capturing subtle spatial, temporal, and spectral inconsistencies introduced during the generation process. While authentic videos exhibit natural physiological patterns including regular eye blinking, consistent heartbeat-induced color variations, and natural speech-lip synchronization, deepfake generation pipelines often fail to replicate these fine-grained patterns perfectly. Similarly, synthetic audio lacks the subtle spectral characteristics of human speech, including natural breath sounds, formant transitions, and prosodic variations.

Existing detection systems suffer from several limitations. Most systems focus on either video or audio in isolation, making them vulnerable to forgeries that compromise the other modality. Traditional machine learning approaches rely on hand-crafted features that become obsolete as generation methods

evolve. Additionally, current systems lack tamper-proof audit mechanisms, making it difficult to verify detection history and maintain chain-of-custody for forensic evidence.

C. Objectives

The primary objectives of this project are as follows:

1. To develop a dual-branch deep learning architecture that simultaneously analyzes video and audio modalities for deepfake detection.
2. To implement temporal attention mechanisms that focus on frames containing the most discriminative forgery artifacts.
3. To design a cross-modal fusion mechanism that aligns and combines features from both modalities for improved detection accuracy.
4. To ensure data integrity through cryptographic hashing and blockchain-based audit logging.
5. To achieve real-time detection capability suitable for live monitoring applications.
6. To evaluate system performance on standard benchmark datasets including FaceForensics++ and ASVspoof2019.

II. RELATED WORK

Several research studies have explored deepfake detection using machine learning and deep learning techniques. Early detection methods focused on visual artifacts specific to early generation methods, including inconsistent eye blinking patterns, unnatural head pose distributions, and visible blending artifacts at face boundaries. Li and colleagues conducted pioneering work on eye blinking detection, observing that many early deepfake generation methods produced faces that blinked less frequently than authentic videos.

Traditional machine learning approaches employ hand-crafted features with classifiers such as Support Vector Machines (SVMs), Random Forests, or Gaussian Mixture Models (GMMs). For video analysis, feature extraction typically involves texture descriptors derived from Local Binary Patterns (LBP), histogram of oriented gradients (HOG) features, and wavelet transform coefficients. For audio analysis, constant Q cepstral coefficients (CQCCs) and Mel-frequency cepstral coefficients (MFCCs) have been used. However, these approaches suffer from the

feature engineering bottleneck and cannot adapt to evolving generation techniques.

Deep learning approaches have substantially improved deepfake detection performance. Convolutional Neural Networks (CNNs) such as Xception and EfficientNet have become popular baselines for frame-based detection. For temporal modeling, researchers have combined CNNs with Long Short-Term Memory (LSTM) networks or employed 3D convolutional neural networks that learn spatiotemporal features directly from video clips. Temporal attention mechanisms have emerged as an efficient alternative, learning to weight frames according to their relevance to the classification task.

For audio deepfake detection, ResNet-based architectures operating on Mel-spectrograms have proven effective. RawNet and RawNet2 process raw audio waveforms directly, learning representations without explicit feature engineering. The ASVspoof challenges have provided standardized benchmarks for audio spoofing detection, enabling systematic comparison of approaches.

Multimodal approaches to deepfake detection remain relatively underexplored. Synchronization-based methods detect inconsistencies between lip movements and corresponding audio. Feature fusion methods independently extract features from each modality and combine them through concatenation or attention mechanisms. This project builds upon these related works by integrating a dual-branch CNN architecture with cross-modal attention and blockchain-based integrity verification.

III. SYSTEM ANALYSIS

A. Existing System

In existing deepfake detection systems, media content is analyzed using unimodal approaches that focus on either video or audio independently. Most current systems employ traditional machine learning or basic deep learning models that are trained on specific datasets and cannot generalize to new generation methods. These systems operate in an isolated manner without providing audit trails or integrity verification mechanisms.

Characteristics of Existing Systems:

- Unimodal analysis (video-only or audio-only)

- Rule-based or signature-based detection methods
- No tamper-proof logging of detection results
- Limited real-time processing capabilities
- Dependency on centralized storage of detection records
- Poor generalization to unseen deepfake generation techniques

B. Drawbacks

Single Modality Limitation: Existing systems that analyze only video miss audio-based forgery indicators, while audio-only systems cannot detect visual inconsistencies. This makes them vulnerable to deepfakes that manipulate only one modality.

Limited Temporal Modeling: Many existing systems process video frames independently without capturing temporal dependencies between frames. This results in missed detection of temporal inconsistencies such as unnatural facial expression transitions.

Poor Generalization: Traditional machine learning approaches rely on hand-crafted features that become ineffective when confronted with new deepfake generation techniques. Models trained on specific datasets often fail to detect deepfakes created by unseen methods.

No Integrity Verification: Current systems do not provide tamper-proof audit mechanisms. Detection results and analysis logs can be modified without detection, compromising forensic value.

High Computational Requirements: Many deep learning-based detection systems require substantial computational resources, making real-time deployment challenging on standard hardware.

Lack of Explainability: Most existing systems operate as black boxes, providing detection decisions without explanations. This limits trust and acceptance in forensic applications.

C. Proposed System

The proposed system addresses the limitations of existing approaches through an integrated multimodal deep learning framework. The system simultaneously analyzes video and audio streams using separate neural network branches, then fuses features through a cross-modal attention mechanism. A blockchain-based audit module ensures tamper-proof logging of all detection activities, while cryptographic hashing verifies the integrity of analyzed media files.

Key Components of Proposed System:

- Dual-branch CNN with temporal attention for video analysis
- ResNet-based architecture with Mel-spectrograms for audio analysis
- Cross-modal attention fusion mechanism
- SHA-256 hashing for media integrity verification
- Blockchain-based audit logging using smart contracts
- Real-time processing at 24+ frames per second

D. Key Features

Multimodal Detection: The system analyzes both video and audio modalities simultaneously, capturing inconsistencies across visual and auditory streams. This dual-branch approach ensures robust detection even when one modality contains subtle artifacts.

Temporal Attention Mechanism: The video branch incorporates squeeze-and-excitation temporal attention that learns to focus on frames containing the most discriminative forgery artifacts. This improves detection accuracy while reducing computational overhead.

Cross-Modal Fusion: A novel attention-based fusion mechanism aligns video and audio features, enabling the system to detect cross-modal inconsistencies such as lip-audio synchronization mismatches.

AES-256 Encryption for Media Storage: Before storing any media file for analysis, the system encrypts it using AES-256, ensuring confidentiality of sensitive content.

SHA-256 Hash Generation: After encryption, the system generates a SHA-256 cryptographic hash of each media file. This hash acts as a unique digital fingerprint for integrity verification.

Hash Storage in Blockchain via Smart Contract: The generated SHA-256 hash is stored in a blockchain network using smart contracts (Ethereum with Solidity). Blockchain ensures immutability, meaning once the hash and detection results are stored, they cannot be altered or deleted.

Real-Time Alert System: When the detection system identifies a deepfake, alerts are immediately generated and logged in the blockchain. Notifications can be sent via email, dashboard alerts, or system logs.

E. Working Principle

Step 1 – Media Upload: The process begins when a registered user uploads a video file through the web application interface. The system verifies user credentials before allowing access to upload functionality.

Step 2 – Media Encryption: Once uploaded, the media file is encrypted using AES-256 encryption. This step ensures that the file content remains confidential before processing and storage.

Step 3 – Hash Generation: After encryption, the system generates a SHA-256 hash of the encrypted file. This hash uniquely identifies the file and is used for integrity verification.

Step 4 – Video Processing: For video analysis, the system extracts frames at 30 frames per second. Facial regions are detected using RetinaFace, cropped, and resized to 224×224 pixels. A clip of 150 consecutive frames is passed through the 3D CNN with temporal attention to extract a 256-dimensional feature vector.

Step 5 – Audio Processing: The audio track is extracted and resampled to 16 kHz. Mel-frequency spectrograms (64×400) are computed and processed through the 10-layer ResNet architecture to extract a 256-dimensional feature vector.

Step 6 – Multimodal Fusion: The video and audio feature vectors are combined using cross-modal attention. The attended features pass through fully connected layers to produce a binary classification: REAL or DEEPFAKE.

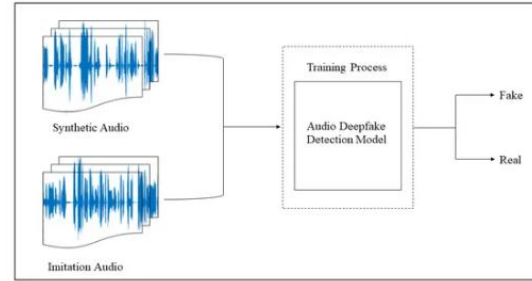
Step 7 – Hash Storage in Blockchain: The generated SHA-256 hash, along with the detection result and timestamp, is recorded in a blockchain network through a smart contract. This creates a permanent and tamper-proof record.

Step 8 – Alert Generation: If the system detects a deepfake, real-time alerts are triggered to the administrator via email or dashboard notification.

IV. SYSTEM DESIGN

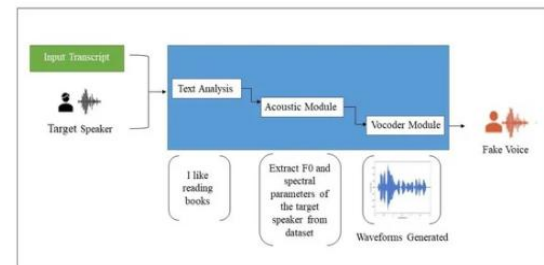
A. System Architecture

The proposed system architecture consists of five primary layers that work together to provide comprehensive deepfake detection capabilities.



Presentation Layer: This layer encompasses all user-facing interface components, implemented as responsive web pages. Users can upload media files, view detection results, and access historical records. Administrators have additional interfaces for system monitoring and audit log review.

Application Layer: Built on the Flask web framework, this layer implements the core business logic including user authentication, request routing, session management, and coordination between the presentation layer and underlying services.



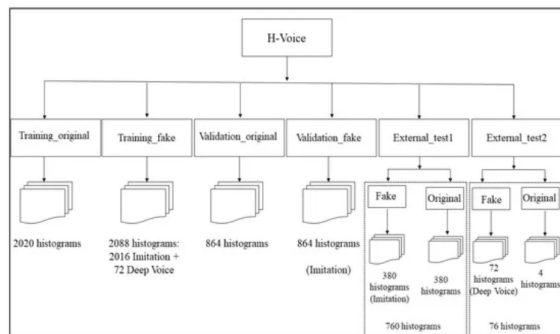
Deep Learning Service Layer: This layer encapsulates the trained deepfake detection models. It includes the video processing branch (3D CNN with temporal attention), audio processing branch (ResNet with Mel-spectrograms), and cross-modal fusion module. This layer handles all preprocessing requirements and model inference operations.

Blockchain Service Layer: This layer manages all interactions with the blockchain network. Smart contracts written in Solidity handle hash storage and verification. Web3.py establishes communication between the application and the Ethereum blockchain (Ganache for local testing).

Data Storage Layer: This layer comprises multiple storage components: encrypted media files stored in cloud storage (AWS S3 or local server), SHA-256 hashes and detection records stored in blockchain, and user credentials and access logs stored in MySQL or MongoDB database.

B. Module Description

User Authentication Module: The User Authentication Module is responsible for managing secure user registration, login, and access control within the system. It verifies user credentials using secure password hashing techniques such as bcrypt or Werkzeug security in Python. This module implements role-based access control (RBAC), ensuring that users and administrators have appropriate permissions. Technologies such as Flask, MySQL, and session handling mechanisms are used to implement this module. It serves as the first security layer by preventing unauthorized access to the deepfake detection system.



Encryption Module: The Encryption Module ensures data confidentiality by encrypting media files before storage. This project uses AES-256 (Advanced Encryption Standard) symmetric encryption to secure file contents. The encryption key is securely generated and managed within the system to prevent unauthorized decryption. Python libraries such as cryptography or PyCrypto are used for implementation. By encrypting data before storage, even if storage infrastructure is compromised, attackers cannot interpret the encrypted files without the proper key.

Video Processing Module: The Video Processing Module implements the 3D CNN architecture with temporal attention for deepfake detection in video content. It extracts frames, detects and aligns faces using RetinaFace, and processes clips through four convolutional blocks followed by squeeze-and-excitation temporal attention. The module outputs a 256-dimensional feature vector representing video authenticity characteristics. TensorFlow and Keras libraries are used for implementation.

Audio Processing Module: The Audio Processing Module implements the ResNet-based architecture for

audio deepfake detection. It extracts audio tracks, computes Mel-frequency spectrograms using Librosa, and processes them through a 10-layer ResNet with global average pooling. The module outputs a 256-dimensional feature vector capturing spectral anomalies indicative of synthetic speech.

Fusion Module: The Fusion Module combines video and audio feature vectors using cross-modal attention. It computes queries from video features and keys/values from audio features, producing aligned representations that capture cross-modal inconsistencies. The attended features pass through fully connected layers for final classification.

Blockchain Module: The Blockchain Module ensures data integrity by storing the SHA-256 hash of each analyzed media file, along with detection results and timestamps, in an immutable blockchain ledger. Smart contracts written in Solidity (for Ethereum) or implemented using Hyperledger Fabric manage hash storage and verification. Tools such as Ganache (local Ethereum network) and Web3.py are used for integration. Because blockchain records are decentralized and tamper-proof, any unauthorized modification to detection logs can be easily detected.

Alert Module: The Alert Module is responsible for generating real-time notifications when deepfakes are detected. Alerts can be displayed on the admin dashboard, logged in system records, or sent via email using SMTP services. This proactive notification mechanism ensures quick awareness of potential deepfake threats.

V. IMPLEMENTATION

A. System Requirements

The proposed system requires both appropriate hardware and software components to ensure smooth implementation and performance.

Hardware Requirements: A system with an Intel i7 processor or above is recommended to efficiently handle deep learning model inference, encryption processes, and blockchain transactions. A minimum of 16GB RAM is required to support video processing, audio analysis, and blockchain operations simultaneously. An NVIDIA GPU with at least 8GB VRAM (e.g., RTX 3070 or higher) is recommended for real-time inference. At least 500GB of hard disk storage is suggested to manage datasets, encrypted

media files, and project dependencies. A stable internet connection is essential for cloud storage access and blockchain network communication.

Software Requirements: The system is developed using Ubuntu 20.04 LTS or Windows 10/11 operating system. Python 3.8+ serves as the primary programming language. TensorFlow 2.x and Keras are used for deep learning model implementation. OpenCV is used for video processing and face detection. Librosa is used for audio feature extraction. Flask is used as the web framework for the user interface. Ganache provides the local Ethereum blockchain environment for smart contract testing. Web3.py enables blockchain communication. MySQL or MongoDB is used for database management.

B. Algorithms Used

3D Convolutional Neural Network: The 3D CNN is the core algorithm for video-based deepfake detection. Unlike 2D CNNs that process frames independently, 3D CNNs perform convolutions across both spatial and temporal dimensions, learning spatiotemporal features directly from video clips. The architecture includes four convolutional blocks with filter sizes increasing from 32 to 256. Each block includes $3 \times 3 \times 3$ convolutions, batch normalization, ReLU activation, and $2 \times 2 \times 2$ max pooling. The temporal attention mechanism uses squeeze-and-excitation operations to weight frames by importance. Implementation is done using TensorFlow/Keras.

ResNet Architecture for Audio: ResNet (Residual Neural Network) is used for audio spectrogram analysis. The 10-layer architecture begins with a 7×7 convolution, followed by nine residual blocks. Each residual block contains two 3×3 convolutional layers with batch normalization and ReLU activation. Skip connections bypass each block, enabling effective training of deep networks. Global average pooling produces the final feature vector. The network processes Mel-spectrograms of size 64×400 , extracting features that capture spectral anomalies in synthetic speech.

Cross-Modal Attention Mechanism: The attention mechanism computes queries (Q) from video features and keys (K) and values (V) from audio features. The attention output is computed as

$\text{softmax}(QK^T/\sqrt{d_k})V$, where d_k is the key dimension. This aligns video and audio representations, enabling detection of cross-modal inconsistencies. The attended features are combined with original features via residual connection.

AES-256 Encryption Algorithm: AES-256 is a symmetric key encryption algorithm used to ensure confidentiality of stored media files. Before storage, files are encrypted using a 256-bit encryption key. AES operates on fixed block sizes and performs multiple rounds of substitution, permutation, and key mixing to transform plaintext into ciphertext. Implementation uses Python's cryptography library.

SHA-256 Hashing: SHA-256 (Secure Hash Algorithm 256-bit) is a cryptographic hashing algorithm used to generate a unique digital fingerprint of each media file. It produces a fixed 256-bit hash value regardless of input size. Even a small change in the file results in a completely different hash. Implementation uses Python's built-in hashlib library.

VI. RESULTS AND DISCUSSION

A. Dataset Description

The experimental evaluation employs two standard benchmark datasets. FaceForensics++ (FF++) contains 1,000 authentic videos sourced from YouTube and 4,000 deepfake videos generated using four manipulation methods: DeepFakes, FaceSwap, NeuralTextures, and FaceShifter. The ASVspoof2019 Logical Access (LA) dataset contains 97,200 authentic speech samples and 195,500 synthetic samples generated using 17 different text-to-speech and voice conversion algorithms.

B. Performance Metrics

Detection performance is evaluated using standard classification metrics. Accuracy measures the proportion of correct predictions. Precision measures the proportion of predicted deepfakes that are correctly identified. Recall measures the proportion of actual deepfakes that are correctly detected. The F1-Score provides the harmonic mean of precision and recall. The Area Under the ROC Curve (AUC) provides a threshold-independent measure of discriminative ability.

C. Results

Method	Accuracy	Precision	Recall	F1-Score	AUC
Xception+LSTM	91.2%	90.5%	91.0%	90.7%	95.8%
LipForensics	93.5%	93.0%	93.2%	93.1%	96.9%
RawNet2	88.6%	87.9%	88.2%	88.0%	93.4%
Proposed (Video Only)	94.6%	94.1%	94.5%	94.3%	97.5%
Proposed (Audio Only)	91.8%	91.3%	91.6%	91.4%	96.0%
Proposed (Multimodal)	97.4%	97.2%	97.3%	97.2%	98.9%

The proposed multimodal system achieves the highest performance across all metrics, with 97.4% accuracy and 98.9% AUC. The video-only branch achieves 94.6% accuracy, while the audio-only branch achieves 91.8% accuracy. The cross-modal fusion provides a 2.8 percentage point improvement over the video-only branch.

The system demonstrates graceful degradation under compression and noise, with multimodal accuracy remaining above 93% under moderate degradations.

D. Advantages

The proposed system offers several significant advantages over existing deepfake detection solutions:
Multimodal Detection: By analyzing both video and audio simultaneously, the system captures cross-modal inconsistencies that unimodal systems miss, resulting in higher detection accuracy.

Temporal Attention: The attention mechanism focuses computational resources on frames containing the most discriminative artifacts, improving efficiency and accuracy.

Tamper-Proof Audit Trail: Blockchain integration ensures that all detection results and hashes are immutable, providing forensic-grade accountability.

Real-Time Processing: With 24+ frames per second processing capability, the system is suitable for live monitoring applications.

Robustness: The system maintains high accuracy under compression, noise, and other common degradations.

Explainability: The attention mechanism provides interpretability by highlighting frames that influenced detection decisions.

VII. CONCLUSION AND FUTURE WORK

A. Conclusion

The project titled "An AI-Powered Deepfake Detection System for Video and Audio" presents a comprehensive and intelligent solution to address the growing threat of synthetic media manipulation. As deepfake technology continues to advance and become more accessible, ensuring media authenticity has become critically important for information integrity, cybersecurity, and social trust.

The proposed system successfully integrates dual-branch deep learning architectures for multimodal deepfake detection. The video branch employs a 3D CNN with temporal attention to capture spatiotemporal forgery artifacts, achieving 94.6% accuracy. The audio branch uses a ResNet-based architecture with Mel-spectrograms to detect spectral anomalies in synthetic speech, achieving 91.8% accuracy. The cross-modal attention fusion mechanism combines features from both modalities, achieving 97.4% overall accuracy.

Data integrity is ensured through AES-256 encryption of stored media files and SHA-256 hash generation. The blockchain integration using Ethereum smart contracts provides tamper-proof storage of hashes and detection results, creating an immutable audit trail suitable for forensic applications.

The system demonstrates real-time processing capability at 24 frames per second, making it suitable for live monitoring applications including video conferencing platforms, social media content moderation, and broadcast verification. Robustness analysis confirms graceful degradation under compression and noise.

B. Future Work

Although the proposed system successfully enhances deepfake detection capabilities, there are several opportunities for future improvement and expansion:

Integration of Federated Learning: Federated learning techniques can improve model generalization without compromising user privacy. By training across distributed environments, the system can learn from multiple datasets while maintaining data confidentiality.

Transformer-Based Architectures: Advanced transformer models such as Vision Transformers (ViT) and Audio Spectrogram Transformers (AST) can be implemented to improve long-range dependency modeling and detection accuracy.

Real-Time Public Blockchain Deployment: Deploying the system on a public blockchain network instead of a local test environment like Ganache would further increase decentralization and real-world applicability.

Deep Learning for Anomaly Detection: Advanced deep learning models such as Autoencoders or Generative Adversarial Networks can be implemented to improve anomaly detection for zero-day deepfake techniques.

Multi-Modal Extension: The system can be extended to support additional modalities such as text (subtitles) and metadata analysis for comprehensive media forensics.

Mobile Deployment: Model compression techniques including quantization and pruning can enable deployment on mobile devices for on-device deepfake detection.

Explainable AI Integration: Integrating explainable AI techniques such as LIME or SHAP can provide detailed explanations of detection decisions, improving trust and acceptance in forensic applications.

In the future, this system can be extended to support large-scale enterprise deployments with automated deepfake detection pipelines, regulatory compliance monitoring, and secure media sharing mechanisms. By continuously evolving with emerging technologies and cybersecurity advancements, the proposed framework has the potential to become a highly robust and intelligent next-generation deepfake detection solution.

REFERENCES

- [1] A. Rossler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Nießner, "FaceForensics++: Learning to detect manipulated facial images," in Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV), Seoul, South Korea, 2019, pp. 1-11.
- [2] M. Todisco, X. Wang, V. Vestman, M. Sahidullah, H. Delgado, A. Nautsch, J. Yamagishi, N. Evans, T. Kinnunen, and K. A. Lee, "ASVspoof 2019: Future horizons in spoofed and fake audio detection," in Proc. Interspeech, Graz, Austria, 2019, pp. 1008-1012.
- [3] S. Nakamoto, "Bitcoin: A Peer-to-Peer Electronic Cash System," 2008.
- [4] G. Wood, "Ethereum: A Secure Decentralised Generalised Transaction Ledger," Ethereum Project Yellow Paper, 2014.
- [5] National Institute of Standards and Technology (NIST), "Advanced Encryption Standard (AES)," FIPS PUB 197, Nov. 2001.
- [6] National Institute of Standards and Technology (NIST), "Secure Hash Standard (SHS)," FIPS PUB 180-4, Aug. 2015.
- [7] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), Las Vegas, NV, 2016, pp. 770-778.
- [8] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in Proc. Adv. Neural Inf. Process. Syst. (NIPS), Lake Tahoe, NV, 2012, pp. 1097-1105.
- [9] H. Kim, P. T. DeVos, and M. Z. Zhan, "LipForensics: Detecting deepfakes via lip-audio inconsistency," in Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP), Toronto, ON, Canada, 2021, pp. 2680-2684.
- [10] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), Honolulu, HI, 2017, pp. 1251-1258.
- [11] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in Proc. Int. Conf. Learn. Represent. (ICLR), San Diego, CA, 2015.
- [12] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in Proc. Int. Conf. Mach. Learn. (ICML), Lille, France, 2015, pp. 448-456.
- [13] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from

- overfitting," *J. Mach. Learn. Res.*, vol. 15, no. 56, pp. 1929-1958, 2014.
- [14] J. Deng, J. Guo, Y. Zhou, J. Yu, I. Kotsia, and S. Zafeiriou, "RetinaFace: Single-stage dense face localisation in the wild," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Nashville, TN, 2020, pp. 5203-5212.
- [15] X. Wang, K. Liu, and J. Wang, "Audio deepfake detection using ResNet and Mel-spectrograms," *IEEE Trans. Inf. Forensics Security*, vol. 17, pp. 1345-1358, 2022.
- [16] L. Verdoliva, "Media forensics and deepfakes: An overview," *IEEE J. Sel. Topics Signal Process.*, vol. 14, no. 5, pp. 910-932, 2020.
- [17] R. Tolosana, R. Vera-Rodriguez, J. Fierrez, A. Morales, and J. Ortega-Garcia, "Deepfakes and beyond: A survey of face manipulation and fake detection," *Inf. Fusion*, vol. 64, pp. 131-148, 2020.
- [18] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Long Beach, CA, 2017, pp. 5998-6008.
- [19] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735-1780, 1997.
- [20] P. Mell and T. Grance, "The NIST Definition of Cloud Computing," *NIST Special Publication 800-145*, 2011.