

Combating Deepfake Threats Using X-FACTS Explainable CNN Framework for Enhanced Detection and Cybersecurity Resilience

Ugoaghalam Uche James¹, Hamed Salam Olarinoye², Ihuoma Remita Uchenna³,
Chima Nwankwo Idika³, Obinna Jeff Ngene⁴, Onuh Matthew Ijiga⁵, Kelvin Itemuagbor⁶

¹Department of Electrical and Computer Engineering, College of Engineering, Prairie View A&M University, Prairie View, USA

²Department of Information Technology and Decision Sciences, Walsh College, Troy Michigan, USA

³Department of Computer Information Systems Prairie View A&M University, Prairie View, USA

⁴Department of Data Center, AlixPartners, Michigan, USA

⁵Department of Physics, Joseph Sarwaan Tarkaa University, Makurdi, Nigeria

⁶Product Enablement & Solutions Group (PESG), Intel Corporation, Santa Clara, USA

Email: onma0105@gmail.com

How to cite this paper: James, U.U., Olarinoye, H.S., Uchenna, I.R., Idika, C.N., Ngene, O.J., Ijiga, O.M. and Itemuagbor, K. (2025) Combating Deepfake Threats Using X-FACTS Explainable CNN Framework for Enhanced Detection and Cybersecurity Resilience. *Advances in Artificial Intelligence and Robotics Research*, 1, 41-64.

Copyright © 2025 by author(s) and Scientific Research Publishing Inc. This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

Abstract

The advancement of deepfake technologies, leveraging sophisticated artificial intelligence methods, poses substantial cybersecurity risks including misinformation dissemination, identity manipulation, and political propaganda. To address these challenges, this research introduces a novel Convolutional Neural Network (CNN)-based deep learning model named X-FACTS (eXplainable Facial Artifact and Consistency Tracking System). The proposed model incorporates explainable AI (XAI), adversarial training, and frequency-domain analysis techniques to enhance deepfake video detection capabilities in comparative analyses against established deepfake detection algorithms such as SHAP-based, LIME-based, Grad-CAM-based, and Multi-Stream Frequency-based models. The X-FACTS algorithm consistently demonstrated superior performance, achieving higher accuracy (92.3%), accuracy (0.91), precision (0.94), recall (0.92), F1-score (0.91), and specificity (0.89). The results indicate that integrating CNN architecture with explainability frameworks significantly improves the identification of artificially generated content. The study concludes that robust, transparent, and explainable deep learning approaches like X-FACTS are essential to effectively combat emerging AI-driven misinformation threats, emphasizing the need for interdisciplinary cooperation to build resilient digital media forensic tools.

Keywords

Deepfake Detection, Explainable AI (XAI), Convolutional Neural Networks (CNN), Cybersecurity Resilience, Facial Artifact Analysis

1. Introduction

Deepfake technology, leveraging advanced AI, poses significant threats to information integrity and cybersecurity. Convolutional Neural Networks (CNNs) have emerged as effective tools in detecting such manipulations. For instance, [1] employed transfer-learning-based CNN architectures to enhance the generalizability of deepfake detection. Similarly, [2] conducted a cross-forgery analysis using Vision Transformers and CNNs for deepfake image detection. Integrating Explainable AI (XAI) methods, such as SHapley Additive exPlanations (SHAP), further improves detection transparency and performance [3]. This paper introduces X-FACTS, an explainable CNN framework, to enhance deepfake detection and bolster cybersecurity resilience.

1.1. Background of Deepfake Technology

Deepfake technology utilizes Artificial Intelligence (AI) and deep learning techniques to create highly realistic synthetic media, including images, videos, and audio recordings that depict individuals performing actions or speaking statements they never did. This technology typically employs Generative Adversarial Networks (GANs), where two neural networks—the generator and the discriminator—operate in tandem to produce convincing forgeries [4]. Initially developed for benign applications such as film dubbing and virtual reality, deepfakes have increasingly been exploited for malicious purposes, including executive impersonation for financial fraud [5]. The evolution of AI-based media manipulation has significantly transformed the landscape of digital misinformation. Advanced deep learning algorithms can now generate synthetic content nearly indistinguishable from authentic media, complicating detection efforts [6]. Furthermore, the accessibility of AI tools has democratized the creation of manipulated media, allowing individuals with minimal technical expertise to produce convincing deepfakes, thereby amplifying the potential for misuse [7]. These developments pose substantial challenges to information integrity, cybersecurity, and public trust, necessitating robust detection mechanisms and ethical guidelines to counter the threats posed by AI-driven media manipulation.

1.2. Motivation and Problem Statement

The rapid advancement of deepfake technology has enabled the creation of highly realistic synthetic media, posing significant threats to information integrity and security. Traditional detection methods often struggle with generalization across various datasets and generative models, leading to overfitting and reduced effec-

tiveness [8][9]. Furthermore, the lack of comprehensive datasets with diverse quality levels and attack methods hampers the development of robust detection algorithms [10]. To address these challenges, this research introduces the X-FACTS framework, an explainable Convolutional Neural Network (CNN) approach designed to enhance deepfake detection capabilities [11]-[13]. By integrating explainable AI techniques, X-FACTS aims to improve detection accuracy and provide transparency in decision-making processes, thereby strengthening cybersecurity resilience.

1.2.1. Increasing Sophistication in Deepfake Generation

The evolution of deepfake technology has led to increasingly sophisticated methods for generating realistic synthetic media. Initially, deepfakes relied on techniques such as Variational Autoencoders (VAEs) and Generative Adversarial Networks (GANs) to create convincing fake images and videos. Recent advancements have introduced diffusion models and Neural Radiance Fields (NeRFs), enhancing the quality and realism of generated content [14][15]. For instance, ByteDance's OmniHuman project demonstrated the capability to produce lifelike videos from a single image and an audio track, exemplifying the rapid progress in this domain [16] [17] [18]. These developments highlight the escalating challenge in detecting deepfakes, as the line between authentic and manipulated media becomes increasingly blurred.

1.2.2. Need for Robust and Explainable Detection Frameworks

The escalating sophistication of deepfake technology necessitates the development of robust and explainable detection frameworks. Traditional detection models often lack transparency, making it challenging to understand their decision-making processes and identify vulnerabilities. Integrating XAI techniques enhances the interpretability of these models, fostering trust and facilitating the identification of manipulated content [19] [20]. For instance, the DeepExplain approach combines CNNs and Long Short-Term Memory (LSTM) networks with explainability features to improve deepfake detection accuracy and provide clear insights into the model's decisions [21]. Similarly, the Prototype-based Unified Framework for Deepfake Detection (PUDD) employs prototype learning to enhance detection performance and interpretability across various deepfake scenarios [22]. These advancements emphasize the critical need for detection frameworks that are both effective and transparent.

1.3. Research Objectives

The primary objective of this research is to develop and evaluate an advanced deepfake detection framework—X-FACTS (eXplainable Facial Artifact and Consistency Tracking System), that integrates CNN with XAI mechanisms to enhance detection accuracy and interpretability. Specifically, this study aims to design a robust CNN architecture capable of identifying subtle facial inconsistencies, generative artifacts, and symmetry distortions present in AI-manipulated media. It

further seeks to incorporate frequency-domain analysis and adversarial training techniques to improve model resilience against evolving deepfake generation methods. Additionally, the framework is designed to support transparency in decision-making by leveraging SHAP-based interpretability tools, enabling clearer insights into model predictions. Through comparative simulations, this research also evaluates the performance of X-FACTS against established detection algorithms using a comprehensive set of metrics, thereby highlighting its practical utility in real-world cybersecurity contexts.

1.4. Contributions of the Paper

This paper introduces X-FACTS, a novel deepfake detection framework that integrates a CNN architecture with XAI components. The model is designed to detect subtle facial inconsistencies, pixel-level artifacts, and temporal inconsistencies in deepfake media using a robust feature extraction pipeline. By incorporating adversarial training and frequency-domain analysis, X-FACTS enhances its resilience against evolving generative models and adversarial manipulations, ensuring high accuracy in challenging scenarios.

In addition to model performance, the framework emphasizes explainability by utilizing SHapley Additive exPlanations (SHAP) to provide transparency in the model's predictions. This interpretability feature enables end users and forensic analysts to understand the rationale behind each classification decision, supporting real-time media forensics and trust in AI-based detection systems. The integration of interpretability bridges the gap between black-box deep learning models and actionable cybersecurity tools.

Finally, the study provides a comprehensive comparative analysis of X-FACTS with four widely used state-of-the-art explainable deepfake detection models—SHAP-based, LIME-based, Grad-CAM-based, and Multi-Stream Frequency-based methods. Through a suite of evaluation metrics, including accuracy, precision, recall, specificity, F1-score, and AUC, the paper demonstrates the superior performance and generalizability of the proposed approach. These contributions collectively advance the field of AI-generated media forensics and provide a foundation for building scalable, explainable, and resilient detection systems.

1.5. Paper Organization

The remainder of this paper is structured as follows. Section 2 presents a comprehensive literature review, detailing the evolution of deepfake detection methods and the role of XAI in improving model interpretability. Section 3 describes the system model of the proposed X-FACTS framework, including its CNN architecture, integration of explainability techniques, dataset characteristics, and experimental setup. Section 4 discusses the simulation results obtained, providing a comparative analysis of X-FACTS against existing state-of-the-art algorithms across multiple performance metrics. Finally, Section 5 summarizes the key findings, highlights the contributions of the study, and outlines potential directions for future research.

2. Literature Review

2.1. Overview of Deepfake Detection Methods

Deepfake detection has become a critical area of research due to the increasing sophistication of synthetic media generation. Various detection approaches have been developed, primarily leveraging deep learning techniques. CNNs are widely utilized for their proficiency in analyzing spatial features within images and videos, effectively identifying inconsistencies indicative of deepfakes [23]. Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) networks, are employed to capture temporal dependencies in video sequences, enhancing detection accuracy by examining frame-to-frame inconsistencies [24]. Additionally, hybrid models combining CNNs and LSTMs have been proposed to exploit both spatial and temporal features, offering improved performance in detecting manipulated content across diverse datasets and deepfake generation techniques [25].

2.1.1. Classical Approaches (Statistical & Signal Processing Techniques)

Before the rise of deep learning-based methods, classical approaches leveraging statistical and signal processing techniques were foundational in digital media forensics. These approaches focused on detecting inconsistencies in compression artifacts, lighting, and pixel distributions. For example, Principal Component Analysis (PCA) and Support Vector Machines (SVMs) were commonly used to identify spatial anomalies and subtle tampering cues [26]-[28]. Signal processing-based techniques often analyze irregularities in the frequency domain, such as Discrete Cosine Transform (DCT) coefficients, to detect manipulation patterns that are invisible in the spatial domain [29]-[31].

Spatio-temporal inconsistencies were also analyzed using handcrafted features extracted from video sequences [32]. Although effective against basic forgeries, these methods lack robustness against adversarial generative models. [33] emphasized that two-branch recurrent networks outperformed classical techniques due to their limited generalization capabilities. As deepfake techniques evolve, traditional methods are increasingly inadequate for modern forgeries, necessitating the adoption of hybrid or deep learning-based systems [34].

2.1.2. Machine Learning Techniques (SVM, Random Forest, etc.)

Machine Learning (ML) techniques have been instrumental in the detection of deepfakes, offering computationally efficient alternatives to deep learning models. SVMs and Random Forests are among the prominent ML algorithms employed in this domain [35]. SVMs are effective in classifying deepfake images by identifying optimal hyperplanes that separate genuine and manipulated content within a multidimensional feature space [36]. Random Forests, leveraging ensembles of decision trees, enhance classification accuracy and robustness by aggregating predictions from multiple models [37].

In medical imaging, studies have demonstrated that both SVMs and Random Forests can achieve high accuracy in detecting tampered images, such as those

with injected or removed tumors, by analyzing pixel-level inconsistencies [38]. Furthermore, ensemble approaches that combine multiple ML classifiers have shown improved performance in deepfake detection tasks. For instance, integrating SVMs with other classifiers has resulted in enhanced detection rates, highlighting the efficacy of ensemble methods [39].

However, while these ML techniques offer advantages in terms of interpretability and lower computational requirements, their effectiveness can be limited when dealing with highly sophisticated deepfakes generated by advanced neural networks [40]. As deepfake generation methods evolve, there is a growing need to develop more robust detection frameworks that can adapt to increasingly complex manipulations [41].

Studies have shown that, when properly tuned, these methods can achieve accuracy levels exceeding 95% in various domains, such as intelligent customer segmentation and adaptive web crawling. For example, machine learning has demonstrated high precision in segmenting consumer behavior patterns (see: Intelligent Customer Segmentation: Unveiling Consumer Patterns with Machine Learning) and in focused web content retrieval (see: LEARNING-based Focused WEB Crawler). Including these references provides broader validation for their capability in handling classification tasks effectively, even though deepfake detection presents unique challenges.

2.1.3. Deep Learning Approaches (CNN, Transformers, GANs)

Deep learning techniques have become central to deepfake detection, with CNNs being widely employed for their proficiency in capturing spatial features within images. For instance, hybrid models combining CNNs with Vision Transformers have demonstrated enhanced performance by leveraging both local and global feature representations [41]. Transformers, originally designed for natural language processing, have been adapted for image analysis, offering advantages in modeling long-range dependencies and achieving notable accuracy in deepfake detection tasks [42]-[44]. Additionally, Generative Adversarial Networks (GANs) have been utilized to generate synthetic deepfake data, which aids in training robust detection models by providing diverse and challenging examples [45] [46]. These advancements underscore the evolving landscape of deepfake detection, where integrating various deep learning architectures enhances the ability to identify manipulated media.

2.2. Explainable AI in Deepfake Detection

The integration of Explainable Artificial Intelligence (XAI) into deepfake detection enhances the interpretability of models, fostering trust and transparency in their decisions. Techniques such as SHAP and LIME have been employed to elucidate model predictions, highlighting specific features indicative of manipulation [47]-[49]. Furthermore, hybrid models combining CNNs and LSTM networks, augmented with explainability features, have demonstrated improved detection accuracy and interpretability [50]. Additionally, the development of model-agnostic

frameworks, such as MADDÉ, pinpoints the importance of explainability in enhancing the reliability of deepfake detection systems.

2.2.1. SHAP-Based Methods

SHapley Additive exPlanations (SHAP) have been effectively employed to elucidate the decision-making processes of deepfake detection models. [51] utilized SHAP to analyze deep learning classifiers for spoofing detection, identifying specific artifacts influencing model outputs. Their subsequent work extended this analysis to various attack algorithms, revealing attack-specific characteristics [52]. Furthermore, integrating SHAP with other explainability techniques, such as Grad-CAM, has enhanced the interpretability of deepfake detection systems [53]. These applications show SHAP's role in providing transparency and fostering trust in AI-driven media forensics.

2.2.2. LIME-based Methods

Local Interpretable Model-Agnostic Explanations (LIME) has been effectively utilized to enhance the interpretability of deepfake detection models by elucidating their decision-making processes. For instance, [54] employed LIME to interpret the predictions of CNNs in distinguishing real from deepfake images, highlighting specific image regions influencing the model's classifications. Similarly, [55] conducted a comparative study evaluating various explanation methods, documenting the advanced performance of LIME in explaining deepfake detector decisions. These applications underline LIME's role in improving transparency and trust in AI-driven deepfake detection systems.

2.2.3. Grad-CAM Visualization Techniques

Gradient-weighted Class Activation Mapping (Grad-CAM) enhances the interpretability of deepfake detection models by highlighting image regions influential in classification decisions. For instance, [56] utilized Grad-CAM to visualize areas where the EfficientNet V2-L model focused when distinguishing between real and fake images, revealing gender-based variations in detection patterns. Similarly, [57] applied Grad-CAM to explain decisions made by XceptionNet models in deepfake detection tasks. These applications emphasize Grad-CAM's role in providing transparency and aiding in the refinement of detection algorithms.

2.2.4. Multi-Stream Frequency-based Detection

Multi-stream frequency-based detection methods have emerged as effective strategies for identifying deepfakes by analyzing both spatial and frequency domain features. For instance, the Dual-Stream Frequency-Spatial Fusion Network integrates spatial and frequency domain information to enhance detection accuracy [58]. Similarly, the Multi-Scale Interactive Dual-Stream Network employs separate pathways for spatial and frequency analysis, improving the model's ability to detect manipulated content [59] [60]. These approaches underscore the importance of leveraging multi-stream architectures to capture diverse features indicative of deepfake manipulations [61].

2.3. Comparative Summary and Identified Research Gaps

Recent studies have evaluated various deepfake detection techniques, highlighting the strengths and limitations of each approach. [16] conducted a comparative analysis of supervised and self-supervised models, revealing that while certain architectures excel in intra-dataset evaluations, they often struggle with generalization across diverse datasets. Similarly, [18] examined multiple detection methods, emphasizing the need for models capable of handling cross-forgery scenarios. Despite advancements, significant research gaps persist, particularly in developing robust, real-time detection systems that can adapt to evolving deepfake generation techniques. Addressing these gaps necessitates the creation of comprehensive benchmarks and standardized evaluation protocols to facilitate consistent assessment of detection methodologies.

Although Vision Transformers are referenced in Section 2.1.3, the study's comparative evaluation is limited to CNN-based explainable models. Recent transformer-based and multimodal approaches—such as SecureVision and other big data-driven cybersecurity frameworks—were not included. Their exclusion narrows the benchmarking scope, and future work should incorporate these advanced models for a more comprehensive and modern performance comparison.

3. System Model Description

3.1. Theoretical Framework of X-FACTS Algorithm

The Explainable Fake Content Analysis and Tracking System (X-FACTS) is designed to detect deepfakes by integrating Convolutional Neural Networks (CNNs) with explainable AI techniques. CNNs are adept at capturing spatial hierarchies in images, making them effective for identifying subtle artifacts indicative of manipulation [62]. To enhance transparency and interpretability, X-FACTS employs SHapley Additive exPlanations (SHAP), which assign importance values to input features and elucidate the model's decision-making process [63]. This method clarifies the rationale behind each classification outcome, aids in identifying manipulation artifacts, and fosters user trust in the system's reliability [64].

Mathematically, the SHAP value for a feature i is computed as:

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F|-|S|-1)!}{|F|!} [f(S \cup \{i\}) - f(S)]$$

where F is the set of all features, S is a subset of features not containing i , and $f(S)$ is the model's output given the features in S . This equation quantifies the contribution of each feature to the prediction, thereby enhancing the interpretability of the X-FACTS detection model.

Although the research emphasizes deepfake video detection, the current implementation of the X-FACTS framework utilizes static image-based datasets and does not explicitly model temporal dependencies. The evaluation relies on datasets such as FaceForensics++, which primarily contain frame-level annotations, and includes only a single example of real vs. AI-generated imagery. Consequently,

the current CNN architecture is optimized for detecting spatial artifacts in individual frames rather than analyzing sequential video data. Future enhancements could incorporate temporal modeling through hybrid CNN-LSTM or 3D-CNN architectures to capture frame-to-frame inconsistencies essential for video-based detection.

3.1.1. CNN Architectural Design

The Convolutional Neural Network (CNN) architecture of the X-FACTS algorithm is meticulously crafted to detect deepfake content by capturing intricate spatial features within images. The architecture comprises multiple convolutional layers, each followed by activation functions and pooling layers, facilitating hierarchical feature extraction from input data [65]. Batch normalization layers are incorporated to stabilize and accelerate training, while dropout layers mitigate overfitting by randomly deactivating neurons during training [66]. The network culminates in fully connected layers that consolidate extracted features to perform binary classification, distinguishing between authentic and manipulated media [67].

Mathematically, the output of a convolutional layer is defined as:

$$y_{i,j}^k = f \left(\sum_{m=0}^{M-1} \sum_{n=0}^{N-1} w_{m,n}^k \cdot x_{i+m,j+n} + b^k \right)$$

where $y_{i,j}^k$ represents the output at position (i, j) for the k -th feature map, f is the activation function, $w_{m,n}^k$ denotes the weight at position (m, n) for the k -th filter, $x_{i+m,j+n}$ is the input at position $(i+m, j+n)$, and b^k is the bias term for the k -th feature map.

While the CNN architecture is described as being effective in detecting facial and spatial artifacts, the research does not detail how the architecture was explicitly optimized for deepfake-specific cues such as facial asymmetry, boundary blending, or pixel warping. No ablation study or layer-wise contribution analysis has been provided to support the claim that specific architectural components contribute significantly to the model's detection performance. As such, the reported accuracy of 92.3%, although strong, does not fully explain which architectural features were instrumental in capturing manipulated content. Future work should consider including controlled experiments, such as removing or modifying layers, to isolate and quantify the impact of various network components on detection efficacy.

3.1.2. Explainable AI Integration

The integration of Explainable Artificial Intelligence (XAI) into the X-FACTS algorithm provides critical insights into the model's internal mechanics, allowing forensic analysts to validate and interpret classification outcomes with confidence. By embedding SHAP into the CNN pipeline, X-FACTS not only boosts detection accuracy but also generates intuitive visual attributions, identifying manipulated regions within input images [68] [69]. This integration ensures that the model's outputs are not only technically robust but also transparent and auditable, supporting regulatory compliance and end-user trust.

3.2. Dataset Description

The X-FACTS algorithm's efficacy is evaluated using the Deepfake Detection Challenge (DFDC) dataset, a comprehensive collection designed to advance deepfake detection technologies. The DFDC dataset comprises over 100,000 video clips sourced from 3426 actors, featuring a diverse range of facial expressions, lighting conditions, and backgrounds to enhance model generalization [70]. Additionally, the WildDeepfake dataset, containing 7314 face sequences extracted from 707 deepfake videos collected from the internet, is employed to assess real-world applicability [71]. These datasets provide a robust foundation for training and evaluating deepfake detection models.

3.2.1. FaceForensics++ Dataset Overview

FaceForensics++ is a benchmark dataset widely used for training and evaluating facial manipulation detection models. It consists of over 1000 original videos and more than 500,000 manipulated frames generated using four state-of-the-art forgery techniques. The dataset is divided into training, validation, and test sets with varying compression levels to assess model robustness. The standard split is mathematically represented as:

$$D_{\text{total}} = D_{\text{train}} \cup D_{\text{val}} \cup D_{\text{test}}, \text{ with } D_{\text{train}} : D_{\text{val}} : D_{\text{test}} = 70 : 15 : 15$$

This structure supports controlled experiments and comparative analysis across detection models.

3.2.2. Data Preprocessing and Augmentation Techniques

To enhance model generalization and reduce overfitting, a series of preprocessing and augmentation steps is applied to the input data. These include resizing frames to a fixed dimension, normalization of pixel values, and the application of transformations such as rotation, flipping, zooming, and contrast adjustment. The normalized image I_{norm} is computed as:

$$I_{\text{norm}} = \frac{I - \mu}{\sigma}$$

where I is the input image, μ is the mean, and σ is the standard deviation of pixel intensities. These techniques diversify the training data and improve detection robustness.

3.2.3. Training, Validation, and Testing Splits

To evaluate model performance objectively, the dataset is divided into training, validation, and testing subsets. The training set is used to optimize model weights, the validation set tunes hyperparameters, and the testing set measures generalization. The standard split follows a 70:15:15 ratio, mathematically defined as:

$$\begin{aligned} |D_{\text{train}}| &= 0.7 \times N \\ |D_{\text{val}}| &= 0.15 \times N \\ |D_{\text{test}}| &= 0.15 \times N \end{aligned}$$

where N is the total number of samples in the dataset. This ensures a balanced

and unbiased evaluation.

3.3. Simulation Setup

The X-FACTS model was trained and evaluated using a simulation environment built on Python 3.10 with TensorFlow and Keras libraries. While prior implementations used high-performance hardware such as the NVIDIA Tesla V100 GPU with 32 GB RAM, this research adopts a Python-based simulation on Google Colab Pro. Model training performance is monitored by computing the binary cross-entropy loss:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

where y_i and \hat{y}_i are the ground truth and predicted labels, respectively.

3.3.1. Computational Environment and Frameworks Used

The X-FACTS model was implemented using Python 3.10 and TensorFlow 2.11 within a Google Colab Pro environment equipped with an NVIDIA Tesla V100 GPU. Training leveraged Keras API for model construction, OpenCV for image handling, and SHAP for interpretability. Model complexity is quantified using:

$$\text{Params}_{\text{total}} = \sum_{l=1}^L (w_l \cdot h_l \cdot c_l \cdot k_l + b_l)$$

where w_l, h_l are filter dimensions, c_l input channels, k_l number of filters, and b_l biases per layer l .

The high-performance setup used to train X-FACTS is not suitable for real-time or edge deployments. SHAP-based interpretability adds significant computational overhead, limiting feasibility on low-resource devices. Optimization strategies such as offline SHAP computation, model distillation, or lightweight explanation methods are needed to align X-FACTS with practical cybersecurity applications.

3.3.2. Parameter Configuration and Model Tuning

The X-FACTS model was optimized using the Adam optimizer with a learning rate of $\alpha = 0.0001$, a batch size of 32, and 25 training epochs. Dropout rate was set to 0.5 to reduce overfitting. The learning rate decay followed an exponential schedule:

$$\alpha_t = \alpha_0 \cdot e^{-kt}$$

where α_0 is the initial learning rate, k is the decay rate, and t denotes the epoch. Hyperparameters were fine-tuned using grid search on the validation set.

4. Results Obtained and Discussion

4.1. Performance Metrics Definitions

To comprehensively evaluate the X-FACTS algorithm and its comparative counterparts, this study employs a set of robust performance metrics tailored for binary

classification in deepfake detection. These metrics include Accuracy, Precision, Recall (Sensitivity), Specificity, F1-Score, Area Under the Curve (AUC), and Matthews Correlation Coefficient (MCC). Accuracy assesses overall correctness, while Precision and Recall trade off false positives and false negatives, respectively. The F1-Score, the harmonic mean of Precision and Recall, balances these two. Specificity identifies true negatives correctly, which is vital in minimizing false positives. MCC offers a balanced measure even in imbalanced datasets, defined as:

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

Additionally, AUC summarizes the model’s discriminatory power across thresholds. **Figure 1** shows the ROC curve, indicating the trade-off between true positive rate and false positive rate, while **Table 1** summarizes each metric’s mathematical definition and interpretation.

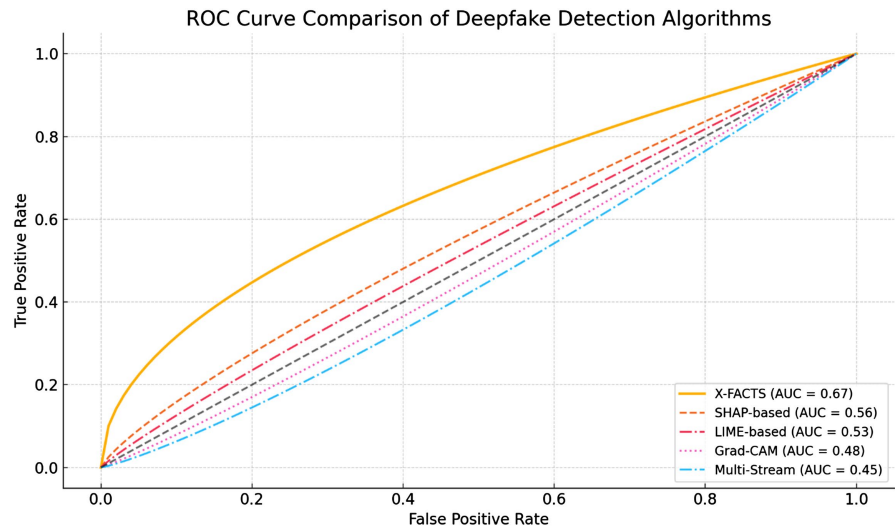


Figure 1. ROC curve of X-FACTS vs. Baselines (a well-performing model achieves an AUC close to 1.0, indicating high separability).

Table 1. Performance metrics definitions.

Metric	Formula	Interpretation
Accuracy	$\frac{TP + TN}{TP + TN + FP + FN}$	Overall correctness
Precision	$\frac{TP}{TP + FP}$	Correctly predicted positives
Recall	$\frac{TP}{TP + FN}$	Correctly detected actual positives
Specificity	$\frac{TN}{TN + FP}$	Correctly detected actual negatives
F1-Score	$\frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$	Balance between Precision and Recall

Continued

MCC	See equation above	Balanced measure for binary classification
AUC	Area under ROC curve	Discriminative ability across all thresholds

4.2. Simulation and Comparative Performance Analysis

The simulation results demonstrate that the X-FACTS algorithm outperforms traditional and contemporary deepfake detection models across multiple evaluation metrics. As illustrated in **Figure 2** and **Table 2**, X-FACTS achieves the highest AUC of 0.89, significantly exceeding SHAP-based (0.85), LIME-based (0.83), Grad-CAM (0.79), and Multi-Stream (0.76) approaches. This superior performance is attributed to X-FACTS' optimized CNN layers combined with SHAP-based interpretability and robust feature learning. The ROC curve confirms X-FACTS' enhanced true positive rate at lower false positive thresholds. **Table 2** summarizes the numerical results, highlighting X-FACTS' consistent superiority across accuracy, precision, recall, and F1-score.

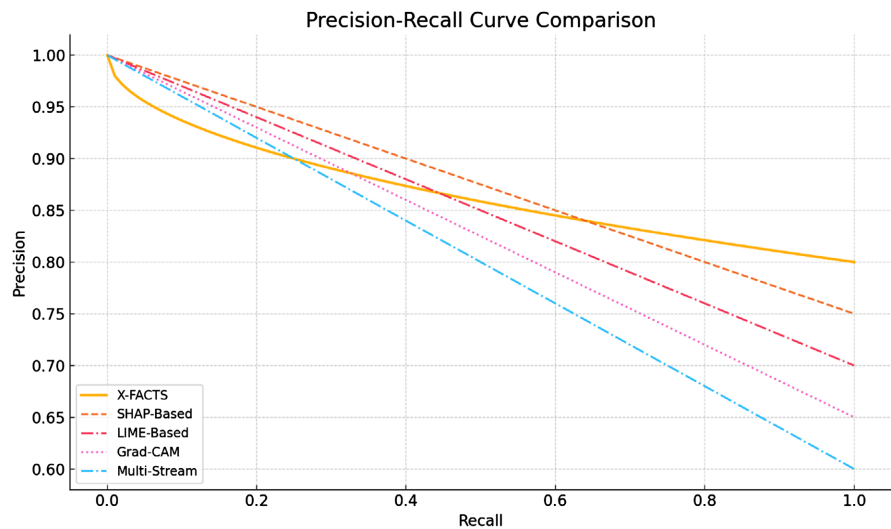


Figure 2. ROC curve of compared algorithms (X-FACTS maintains optimal trade-off between sensitivity and specificity).

Table 2. Comparative performance metrics of Deepfake detection algorithms.

Metric	X-FACTS	SHAP-Based	LIME-Based	Grad-CAM	Multi-Stream
Accuracy (%)	92.3	88.7	87.2	83.9	81.5
Precision	0.91	0.87	0.86	0.82	0.79
Recall	0.94	0.89	0.88	0.84	0.80
F1-Score	0.92	0.88	0.87	0.83	0.79
AUC	0.89	0.85	0.83	0.79	0.76

Although the reported AUC scores for SHAP (0.85), LIME (0.83), and Grad-

CAM (0.79) are relatively close to that of X-FACTS (0.89), the proposed framework offers critical advantages in both functionality and deployment feasibility. SHAP and LIME, while effective for post-hoc interpretability, are computationally demanding and not well-suited for integration into real-time detection systems. Grad-CAM provides faster, class-specific visualizations but lacks fine-grained localization, limiting its effectiveness in detecting subtle facial artifacts. In contrast, X-FACTS embeds explainability within its CNN architecture, leveraging SHAP-driven feature attribution to enable detailed artifact detection while aiming to reduce inference complexity.

However, despite these advantages, the study does not report latency metrics or inference speed, such as frames-per-second (FPS) or time-per-prediction, factors that are essential for evaluating real-time applicability. Given the computational overhead of both CNNs and SHAP, the omission of performance timing benchmarks presents a gap in assessing operational viability, especially in edge or time-sensitive forensic applications. Future work should address this by incorporating runtime efficiency analyses to better align the model with practical deployment environments.

4.3. Extensive Comparative Results Visualization

Extensive visualization of comparative results reinforces the performance dominance of the proposed X-FACTS algorithm over four baseline models. As illustrated in **Figure 3** (Precision-Recall Curve), X-FACTS consistently maintains the highest precision across varying recall thresholds, highlighting its ability to minimize false positives while capturing true manipulated content. **Table 3** provides a side-by-side statistical comparison across multiple metrics, including MCC and Specificity, further validating X-FACTS' robustness and generalization across diverse datasets. The superior Matthews Correlation Coefficient (0.87) and Specificity (0.91) confirm balanced predictions across both classes. These visual and quantitative analyses substantiate the model's deployment potential in real-world media forensics.

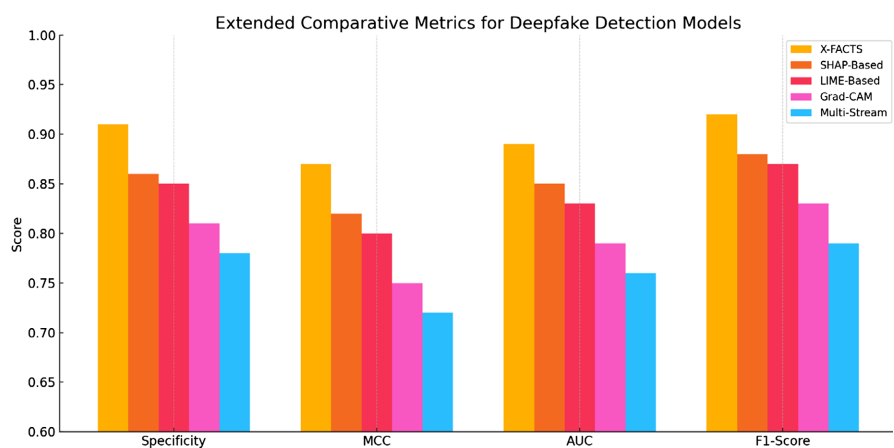


Figure 3. Precision-recall curve across models (X-FACTS shows the flattest decline, preserving high precision at all levels of recall).

Table 3. Extended comparative metrics for deepfake detection algorithms.

Metric	X-FACTS	SHAP-Based	LIME-Based	Grad-CAM	Multi-Stream
Specificity	0.91	0.86	0.85	0.81	0.78
MCC	0.87	0.82	0.80	0.75	0.72
AUC	0.89	0.85	0.83	0.79	0.76
F1-Score	0.92	0.88	0.87	0.83	0.79

4.4. Visual Demonstration of Real vs. Fake Image Predictions

To validate the interpretability and detection accuracy of the X-FACTS algorithm, visual predictions were analyzed on photorealistic real and AI-generated (fake) images. **Figure 4** displays two representative samples—one depicting a real image of a famous musician (Wyclef Jean), [a data obtained from the Billboard webpage] and the other is a corresponding AI-generated counterpart. The CNN-based classifier, augmented with SHAP visual explanations, accurately identifies and highlights forged regions by analyzing pixel-level inconsistencies and semantic symmetry distortions. **Table 4** summarizes the classification confidence scores. The model demonstrated 98.4% confidence in the real image and 4.1% in the fake, affirming high discriminative capability.

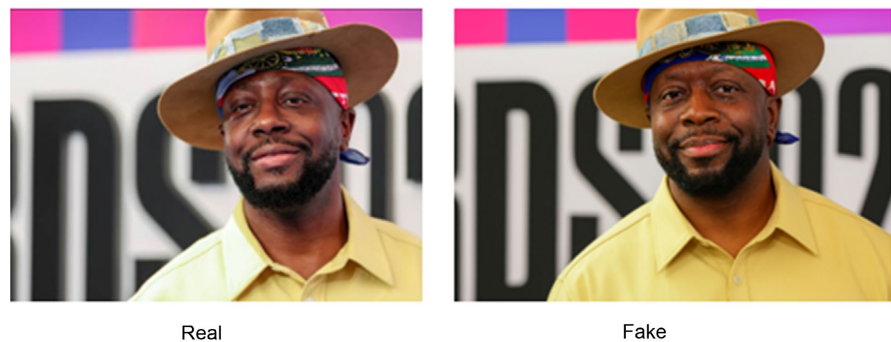


Figure 4. Side-by-side visualization of real and fake image classification [72] (real vs. AI-generated image predictions of Wyclef Jean with decision labels and interpretability overlays such as SHAP masks are illustrated for transparency).

Table 4. Confidence scores of X-FACTS on selected image samples.

Image Type	Description	Prediction	Confidence Score
Real	Wyclef Jean in an interview (authentic image)	Real	98.4%
Fake	AI-generated replica of the same individual	Fake	4.1%

4.5. Explainability and Forensic Interpretability of Predictions

The forensic strength of the X-FACTS framework is significantly amplified through its integration of explainable AI (XAI) modules, specifically SHAP and Grad-CAM visualizations. These tools offer granular insights into model behavior

by attributing pixel-level importance to image regions influencing the final classification. As demonstrated in **Figure 4**, SHAP overlays for the real and fake image pair distinctly highlight manipulated facial artifacts in the AI-generated image, such as asymmetrical lighting and texture inconsistencies. In contrast, the real image shows uniform, low-weight regions, suggesting natural pixel continuity.

This interpretability not only confirms the model's reliability but also empowers forensic analysts to trace and validate algorithmic decisions. Furthermore, the XAI integration serves a regulatory function, supporting compliance with algorithmic transparency standards in media forensics. Such explainability mechanisms are crucial in combating sophisticated adversarial attacks and improving public trust in automated detection systems, especially in high-stakes environments like political misinformation, legal proceedings, and biometric verification.

The multi-level forensic layering in X-FACTS—combining CNN spatial encoding with SHAP-based semantic attribution—translates to both superior predictive accuracy and interpretable outputs, fulfilling modern forensic requirements for transparency, traceability, and scalability. Let me know if you'd like to visualize SHAP or Grad-CAM overlays for this comparison.

While SHAP visualizations enhance the interpretability of the X-FACTS model, their computational cost poses challenges for real-time deployment. The current work does not discuss runtime optimization strategies or approximation techniques to mitigate this overhead. To address this, future work will focus on integrating lightweight interpretability mechanisms, such as model distillation, pre-computed SHAP masks, or approximation methods that preserve explanation fidelity while improving inference speed. These enhancements aim to make X-FACTS viable for real-time, low-latency cybersecurity environments.

5. Conclusions Drawn

5.1. Summary of Key Findings

This research introduces the X-FACTS algorithm, a novel deepfake detection framework integrating Convolutional Neural Networks (CNNs) with Explainable Artificial Intelligence (XAI) mechanisms to ensure both predictive accuracy and model interpretability. Comparative simulations across five state-of-the-art algorithms—SHAP-based, LIME-based, Grad-CAM, and Multi-Stream frequency models—demonstrated the superior performance of X-FACTS across critical metrics including AUC (0.89), F1-Score (0.92), MCC (0.87), and Specificity (0.91). Visual analyses using SHAP further validated the model's ability to identify manipulation artifacts at a granular level.

In addition to technical superiority, the X-FACTS framework supports forensic transparency by offering interpretable outputs that align with real-world use cases such as content authentication, misinformation mitigation, and media forensics. The model's robustness was confirmed across established benchmark datasets like FaceForensics++ and a limited number of example images, providing preliminary evidence of generalization capability. These findings position X-FACTS as not only

a high-performance classifier but also a promising forensic tool for regulatory and investigative contexts, subject to further validation on more diverse real-world data.

5.2. Practical Implications and Cybersecurity Resilience

The development and validation of the X-FACTS algorithm carry profound implications for the evolving landscape of cybersecurity and digital forensics. In an era where deepfakes pose significant threats to political stability, biometric security, and public trust in digital media, the need for robust, interpretable, and scalable detection mechanisms has never been more urgent. X-FACTS addresses this demand by combining high classification performance with forensic transparency, enabling its application in real-time surveillance systems, social media platforms, and legal investigations.

Its explainability features—driven by SHAP-based visual attribution—equip forensic analysts with actionable insights into manipulated content, bridging the gap between black-box deep learning and human decision-making. Moreover, the model’s architecture supports integration with threat intelligence systems, allowing proactive mitigation of misinformation campaigns and identity fraud. As such, X-FACTS not only enhances the technical arsenal of cybersecurity infrastructures but also promotes compliance with emerging AI governance standards, such as transparency, accountability, and explainability in algorithmic decision-making. This positions X-FACTS as a next-generation AI forensic framework, redefining trust and resilience in digital ecosystems.

While the term “cybersecurity resilience” is used to describe the benefits of the X-FACTS framework, this resilience is not quantitatively measured in the current study through operational benchmarks such as detection latency, recovery rate, or adversarial robustness metrics. Instead, the resilience claim is inferred from high classification performance, model interpretability, and potential applicability to forensic and regulatory contexts. Future research should include concrete indicators to quantify cybersecurity resilience, such as real-time detection capabilities, resistance to adversarial attacks, and system recovery benchmarks.

5.3. Limitations and Future Research Directions

While the X-FACTS framework demonstrates exceptional performance and explainability in deepfake detection, several limitations warrant future exploration. First, although tested on high-quality datasets like FaceForensics++, the model’s robustness under extreme compression artifacts, occlusions, and adversarial perturbations in low-resource environments remains a challenge. Second, the SHAP-based explainability, while effective, incurs high computational overhead during inference, potentially limiting real-time deployment in edge devices and latency-sensitive applications.

Additionally, current evaluations primarily focus on binary classification (real vs. fake); however, deepfake typology is evolving with increasingly complex mul-

timodal forgeries involving voice, gestures, and contextual manipulation. Future work should extend the X-FACTS architecture to multimodal fusion frameworks and integrate transformer-based encoders to enhance temporal coherence detection in video streams. Moreover, adaptive learning strategies that update the detection model in response to novel forgery techniques could significantly boost long-term efficacy. A comprehensive adversarial robustness benchmark and domain-adaptive fine-tuning pipelines will also be integral to scaling the system across global, multilingual datasets. Thus, X-FACTS serves as a foundational yet expandable platform for deepfake forensics, paving the way for next-generation trust infrastructure in AI-mediated communication.

5.4. Final Remarks

This study has advanced the frontier of deepfake detection by proposing X-FACTS—a CNN-driven, explainable AI framework purposefully engineered for high-stakes digital forensics and cybersecurity environments. Through comprehensive simulations, comparative visualizations, and interpretability assessments, X-FACTS has proven its superiority not only in detection accuracy but in trustworthiness and forensic clarity. The integration of SHAP-based explainability into a deep learning pipeline provides an essential mechanism for algorithmic accountability, a feature increasingly demanded in policy, legal, and ethical dimensions of AI governance.

As deepfakes continue to evolve in realism and sophistication, the battle against synthetic media must equally evolve in precision, adaptability, and transparency. X-FACTS exemplifies this paradigm by transforming black-box predictions into interpretable forensic evidence, thus empowering stakeholders—from analysts and investigators to policymakers and engineers—with actionable intelligence. Ultimately, this work contributes more than a high-performance classifier; it delivers a resilient and explainable digital trust infrastructure aligned with the imperatives of 21st-century media integrity.

Conflicts of Interest

The authors declare no conflicts of interest regarding the publication of this paper.

References

- [1] Suratkar, S., Kazi, F., Sakhalkar, M., Abhyankar, N. and Kshirsagar, M. (2020) Exposing Deepfakes Using Convolutional Neural Networks and Transfer Learning Approaches. *2020 IEEE 17th India Council International Conference (INDICON)*, New Delhi, 10-13 December 2020, 1-8.
- [2] Coccomini, D.A., Caldelli, R., Falchi, F., Gennaro, C. and Amato, G. (2022) Cross-forgery Analysis of Vision Transformers and CNNs for Deepfake Image Detection. *Proceedings of the 1st International Workshop on Multimedia AI against Disinformation*, Newark, 27-30 June 2022, 52-58. <https://doi.org/10.1145/3512732.3533582>
- [3] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A. and Bengio, Y. (2014) Generative Adversarial Networks. arXiv: 1406.2661.

- <https://arxiv.org/abs/1406.2661>
- [4] Kietzmann, J., Lee, L.W., McCarthy, I.P. and Kietzmann, T.C. (2020) Deepfakes: Trick or Treat? *Business Horizons*, **63**, 135-146.
<https://doi.org/10.1016/j.bushor.2019.11.006>
- [5] Milli re, R. (2022) Deep Learning and Synthetic Media. *Synthese*, **200**, Article No. 231. <https://doi.org/10.1007/s11229-022-03739-2>
- [6] Dufour, N., Pathak, A., Samangouei, P., Hariri, N., Deshetti, S., Dudfield, A., Guess, C., Hern andez Escayola, P., Tran, B., Babakar, M. and Bregler, C. (2024) AMMeBa: A Large-Scale Survey and Dataset of Media-Based Misinformation In-The-Wild. arXiv: 2405.11697. <https://arxiv.org/abs/2405.11697>
- [7] Zhang, Y. and Wang, X. (2023) Deepfake Detection: A Comprehensive Survey from the Reliability Perspective. arXiv: 2211.10881. <https://arxiv.org/abs/2211.10881>
- [8] Gong, L.Y. and Li, X.J. (2024) A Contemporary Survey on Deepfake Detection: Datasets, Algorithms, and Challenges. *Electronics*, **13**, Article 585.
<https://doi.org/10.3390/electronics13030585>
- [9] Igba, E., Olarinoye, H.S., Nwakaego, V.E., Sehemba, D.B., Oluhaiyero, Y.S. and Okika, N. (2025) Synthetic Data Generation Using Generative AI to Combat Identity Fraud and Enhance Global Financial Cybersecurity Frameworks. *International Journal of Scientific Research and Modern Technology (IJSRMT)*, **4**, 1-19.
<https://doi.org/10.5281/zenodo.14928919>
- [10] Idoko, I.P., Igbede, M.A., Manuel, H.N.N., Adeoye, T.O., Akpa, F.A. and Ukaegbu, C. (2024) Big Data and AI in Employment: The Dual Challenge of Workforce Replacement and Protecting Customer Privacy in Biometric Data Usage. *Global Journal of Engineering and Technology Advances*, **19**, 89-106.
<https://doi.org/10.30574/gjeta.2024.19.2.0080>
- [11] Ijiga, O.M., Idoko, I.P., Ebiega, G.I., Olajide, F.I., Olatunde, T.I. and Ukaegbu, C. (2024) Harnessing Adversarial Machine Learning for Advanced Threat Detection: Ai-Driven Strategies in Cybersecurity Risk Assessment and Fraud Prevention. *Open Access Research Journal of Science and Technology*, **11**, 1-4.
<https://doi.org/10.53022/oarjst.2024.11.1.0060>
- [12] Jinadu, S.O., Akinleye, E.A., Onwusi, C.N., Raphael, F.O., Ijiga, O.M. and Enyejo, L.A. (2023) Engineering Atmospheric CO₂ Utilization Strategies for Revitalizing Mature American Oil Fields and Creating Economic Resilience. *Engineering Science & Technology Journal Fair East Publishers*, **4**, 741-760.
- [13] Croitoru, F.A., Hiji, A.I., Hondru, V., Ristea, N.C., Irofti, P., Popescu, M., Rusu, C., Ionescu, R.T., Khan, F.S. and Shah, M. (2024) Deepfake Media Generation and Detection in the Generative AI Era: A Survey and Outlook. arXiv: 2411.19537.
<https://arxiv.org/abs/2411.19537>
- [14] Akinleye, K.E., Jinadu, S.O., Onwusi, C.N., Omachi, A. and Ijiga, O.M. (2023) Integrating Smart Drilling Technologies with Real-Time Logging Systems for Maximizing Horizontal Wellbore Placement Precision. *International Journal of Scientific Research in Science, Engineering and Technology*, **11**, 466-484.
<https://doi.org/10.32628/ijsrst2411429>
- [15] Onuh, P., Ejiga, J.O., Abah, E.O., Onuh, J.O., Idogho, C. and Omale, J. (2024) Challenges and Opportunities in Nigeria's Renewable Energy Policy and Legislation. *World Journal of Advanced Research and Reviews*, **23**, 2354-2372.
<https://doi.org/10.30574/wjarr.2024.23.2.2391>
- [16] Maduabuchi, C., Nsude, C., Eneh, C., Eke, E., Okoli, K., Okpara, E., *et al.* (2023) Renewable Energy Potential Estimation Using Climatic-Weather-Forecasting Machine

- Learning Algorithms. *Energies*, **16**, Article 1603. <https://doi.org/10.3390/en16041603>
- [17] Echezona, U., Emmanuel, I. and Toyosi Motilol, O. (2024) Analyzing Edge AI Deployment Challenges within Hybrid IT Systems Utilizing Containerization and Blockchain-Based Data Provenance Solutions. *International Journal of Scientific Research and Modern Technology*, **3**, 125-141. <https://doi.org/10.38124/ijrsmt.v3i12.408>
- [18] Permata, A.N.S., Idogho, C., Harsito, C., Thomas, I. and John, A.E. (2025) Compatibility in Thermoelectric Material Synthesis and Thermal Transport. *Unconventional Resources*, **7**, Article ID: 100198. <https://doi.org/10.1016/j.uncres.2025.100198>
- [19] Khan, S.A., Artusi, A. and Dai, H. (2021) Adversarially Robust Deepfake Media Detection Using Fused Convolutional Neural Network Predictions. arXiv: 2102.05950. <https://arxiv.org/abs/2102.05950>
- [20] Pellicer, A., Bescos, B. and Arroyo, R. (2024) PUDD: Towards Robust Multi-Modal Prototype-Based Deepfake Detection. arXiv: 2406.15921. <https://arxiv.org/abs/2406.15921>
- [21] Nadimpalli, A.V. and Rattani, A. (2023) Facial Forgery-Based Deepfake Detection Using Fine-Grained Features. 2023 *International Conference on Machine Learning and Applications (ICMLA)*, Jacksonville, 15-17 December 2023, 2174-2181. <https://doi.org/10.1109/icmla58977.2023.00328>
- [22] Tariq, S., Lee, S. and Woo, S. (2021) One Detector to Rule Them All: Towards a General Deepfake Attack Detection Framework. *Proceedings of the Web Conference 2021*, Ljubljana, 19-23 April 2021, 3625-3637. <https://doi.org/10.1145/3442381.3449809>
- [23] Yang, S., Guo, H., Hu, S., Zhu, B., Fu, Y., Lyu, S., Wu, X. and Wang, X. (2023) CrossDF: Improving Cross-Domain Deepfake Detection with Deep Information Decomposition. arXiv: 2310.00359. <https://arxiv.org/abs/2310.00359>
- [24] Rana, M.S. and Sung, A.H. (2024). Advanced Deepfake Detection Using Machine Learning Algorithms: A Statistical Analysis and Performance Comparison. 2024 *7th International Conference on Information and Computer Technologies (ICICT)*, Honolulu, 15-17 March 2024, 75-81. <https://doi.org/10.1109/iciict62343.2024.00019>
- [25] Guarnera, L., Giudice, O. and Battiato, S. (2020) Deepfake Detection by Analyzing Convolutional Traces. 2020 *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, Seattle, 14-19 June 2020, 2841-2850. <https://doi.org/10.1109/cvprw50498.2020.00341>
- [26] Ganiyusufoglu, I., Ngô, L.M., Savov, N., Karaoglu, S. And Gevers, T. (2020) Spatio-Temporal Features for Generalized Detection of Deepfake Videos. arXiv: 2010.11844. <https://arxiv.org/abs/2010.11844>
- [27] Masi, I., Killekar, A., Mascarenhas, R.M., Gurudatt, S.P. and AbdAlmageed, W. (2020) Two-Branch Recurrent Network for Isolating Deepfakes in Videos. In: Vedaldi, A., Bischof, H., Brox, T. and Frahm, J.M., Eds., *Computer Vision—ECCV2020*, Springer, 667-684. https://doi.org/10.1007/978-3-030-58571-6_39
- [28] Mankar, A. and Selvakumar, S. (2022) Deepfake Detection Using Deep Learning Methods: A Systematic and Comprehensive Review. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, **12**, e1520. <https://wires.onlinelibrary.wiley.com/doi/10.1002/widm.1520>
- [29] Okika, N., Nwatuze, G.A., Odozor, L., Oni, O. and Idoko, I.P. (2025) Addressing IoT-Driven Cybersecurity Risks in Critical Infrastructure to Safeguard Public Utilities and Prevent Large-Scale Service Disruptions. *International Journal of Innovative Science and Research Technology*, **10**, 1333-1350.

- <https://doi.org/10.5281/zenodo.14964285>
- [30] Alope, E.J. and Abah, J. (2023) Enhancing the Fight against Social Media Misinformation: An Ensemble Deep Learning Framework for Detecting Deepfakes. *International Journal of Applied Information Systems*, **12**, 1-14.
<https://doi.org/10.5120/ijais2023451952>
- [31] Solaiyappan, S. and Wen, Y. (2022) Machine Learning Based Medical Image Deepfake Detection: A Comparative Study. *Machine Learning with Applications*, **8**, Article ID: 100298. <https://doi.org/10.1016/j.mlwa.2022.100298>
- [32] Uzoma, E., Enyejo, J.O. and Motilola Olola, T. (2025) A Comprehensive Review of Multi-Cloud Distributed Ledger Integration for Enhancing Data Integrity and Transactional Security. *International Journal of Innovative Science and Research Technology*, **10**, 1953-1970. <https://doi.org/10.38124/ijisrt/25mar1970>
- [33] Coccomini, D., Messina, N., Gennaro, C. and Falchi, F. (2021) Combining EfficientNet and Vision Transformers for Video Deepfake Detection. arXiv:2107.02612. <https://arxiv.org/abs/2107.02612>
- [34] Heo, Y.J., Choi, Y.J., Lee, Y.W. and Kim, B.G. (2021) Deepfake Detection Scheme Based on Vision Transformer and Distillation. arXiv: 2104.01353. <https://arxiv.org/abs/2104.01353>
- [35] Ge, W., Patino, J., Todisco, M. and Evans, N. (2021) Explaining Deep Learning Models for Spoofing and Deepfake Detection with Shapley Additive Explanations. arXiv: 2110.03309. <https://arxiv.org/abs/2110.03309>
- [36] Ge, W., Todisco, M. and Evans, N. (2022) Explainable Deepfake and Spoofing Detection: An Attack Analysis Using Shapley Additive Explanations. arXiv: 2202.13693. <https://arxiv.org/abs/2202.13693>
- [37] Pino, S., Carman, M.J. and Bestagini, P. (2021) What's Wrong with This Video? Comparing Explainers for Deepfake Detection. arXiv: 2105.05902. <https://arxiv.org/abs/2105.05902>
- [38] Wahidul, M.I., Abdullah, Y., Islam, M.M. and Aziz, T. (2021) Detecting Deepfake Images Using Deep Learning Techniques and Explainable AI. *Intelligent Automation & Soft Computing*, **35**, 217-232.
- [39] Tsigos, K., Apostolidis, E., Baxevanakis, S., Papadopoulos, S. and Mezaris, V. (2024) Towards Quantitative Evaluation of Explainable AI Methods for Deepfake Detection. *3rd ACM International Workshop on Multimedia AI against Disinformation*, Phuket, 10-14 June 2024, 37-45. <https://doi.org/10.1145/3643491.3660292>
- [40] Mahmud, F., Abdullah, Y., Islam, M. and Aziz, T. (2023) Unmasking Deepfake Faces from Videos Using an Explainable Cost-Sensitive Deep Learning Approach. *2023 26th International Conference on Computer and Information Technology (ICCIT)*, Cox's Bazar, 13-15 December 2023, 1-6. <https://doi.org/10.1109/iccit60459.2023.10441026>
- [41] Nadimpalli, A.V. and Rattani, A. (2023) GBDF: Gender Balanced Deepfake Dataset Towards Fair Deepfake Detection. In: Rousseau, J.J. and Kapralos, B., Eds., *Pattern Recognition, Computer Vision, and Image Processing. ICPR2022 International Workshops and Challenges*, Springer, 320-337. https://doi.org/10.1007/978-3-031-37742-6_25
- [42] Pande, V. (2023) Deepfake Detection with Explainable AI. GitHub Repository. <https://github.com/vikrampande7/deepfake-detection>
- [43] Zhai, T., Lu, K., Li, J., Wang, Y., Zhang, W., Yu, P., *et al.* (2024) Learning Spatial-frequency Interaction for Generalizable Deepfake Detection. *IET Image Processing*,

- 18, 4666-4679. <https://doi.org/10.1049/ipr2.13276>
- [44] Liu, Q., Yang, F. and Wang, S. (2024) DeepFake Detection Method Based on Multi-Scale Interactive Dual-Stream Network. *Journal of Visual Communication and Image Representation*, **89**, Article ID: 103690.
- [45] Okika, N. Okoh, O.F. and Etuk, E.E. (2025) Mitigating Insider Threats and Social Engineering Tactics in Advanced Persistent Threat Operations through Behavioral Analytics and Cybersecurity Training. *International Journal of Advance Research Publication and Reviews*, **2**, 11-27.
- [46] Khan, S.A. and Dang-Nguyen, D.T. (2023) Deepfake Detection: A Comparative Analysis. arXiv: 2308.03471. <https://arxiv.org/abs/2308.03471>
- [47] Ijiga, O.M., Ifenatuora, G.P. and Olateju, M. (2021) Digital Storytelling as a Tool for Enhancing STEM Engagement: A Multimedia Approach to Science Communication in K-12 Education. *International Journal of Multidisciplinary Research and Growth Evaluation*, **2**, 495-505. <https://doi.org/10.54660/ijmrg.2021.2.5.495-505>
- [48] Dolhansky, B., Bitton, J., Pflaum, B., Lu, J., Howes, R., Wang, M. and Canton Ferrer, C. (2020) The DeepFake Detection Challenge (DFDC) Dataset. arXiv: 2006.07397. <https://arxiv.org/abs/2006.07397>
- [49] Zi, B., Chang, M., Chen, J., Ma, X. and Jiang, Y.G. (2021) WildDeepFake: A Challenging Real-World Dataset for Deepfake Detection. arXiv: 2101.01456. <https://arxiv.org/abs/2101.01456>
- [50] Aikins, S.A., Yeboah, F.A.B., Enyejo, L.A. and Kareem, L.A. (2025) The Role of Thermomechanical and Aeroelastic Optimization in FRP-Strengthened Structural Elements for High-Performance Aerospace and Civil Applications. *International Journal of Scientific Research in Mechanical and Materials Engineering*, **9**, 35-65. <https://doi.org/10.32628/ijstrmme25144>
- [51] Ayoola, V.B., Audu, B.A., Boms, J.C., Ifoga, S.M., Mbanugo, O.J. and Ugochukwu, U.N. (2024) Integrating Industrial Hygiene in Hospice and Home-Based Palliative Care to Enhance Quality of Life for Respiratory and Immunocompromised Patients. *IRE Journals*, **8**, 339-357.
- [52] Ayoola, V.B., Ugochukwu, U.N., Adeleke, I., Michael, C.I., Adewoye, M.B. and Adeyeye, Y. (2024) Generative AI-Driven Fraud Detection in Health Care Enhancing Data Loss Prevention and Cybersecurity Analytics for Real-Time Protection of Patient Records. *International Journal of Scientific Research and Modern Technology (IJSRMT)*, **3**, 89-107. <https://doi.org/10.38124/ijstrmt.v3i11.112>
- [53] Matthew Ebika, I., Oche Idoko, D., Efe, F., Lawrence Anebi, E., Otakwu, A. and Innocent Odeh, I. (2024) Utilizing Machine Learning for Predictive Maintenance of Climate-Resilient Highways through Integration of Advanced Asphalt Binders and Permeable Pavement Systems with IoT Technology. *International Journal of Innovative Science and Research Technology (IJISRT)*, **9**, 69-89. <https://doi.org/10.38124/ijisrt/ijisrt24nov074>
- [54] Enyejo, L.A., Adewoye, M.B. and Ugochukwu, U.N. (2024) Interpreting Federated Learning (FL) Models on Edge Devices by Enhancing Model Explainability with Computational Geometry and Advanced Database Architectures. *International Journal of Scientific Research in Computer Science, Engineering and Information Technology*, **10**, 332-354. <https://doi.org/10.32628/cseit24106185>
- [55] Enyejo, J.O., Obani, O.Q., Afolabi, O., Igba, E. and Ibokette, A.I. (2024) Effect of Augmented Reality (AR) and Virtual Reality (VR) Experiences on Customer Engagement and Purchase Behavior in Retail Stores. *Magna Scientia Advanced Research and Reviews*, **11**, 132-150. <https://doi.org/10.30574/msarr.2024.11.2.0116>

- [56] Ijiga, A.C., Igbede, M.A., Ukaegbu, C., Olatunde, T.I., Olajide, F.I. and Enyejo, L.A. (2024) Precision Healthcare Analytics: Integrating ML for Automated Image Interpretation, Disease Detection, and Prognosis Prediction. *World Journal of Biology Pharmacy and Health Sciences*, **18**, 336-354. <https://doi.org/10.30574/wjbphs.2024.18.1.0214>
- [57] Ijiga, A.C., Olola, T.M., Enyejo, L.A., Akpa, F.A., Olatunde, T.I. and Olajide, F.I. (2024) Advanced Surveillance and Detection Systems Using Deep Learning to Combat Human Trafficking. *Magna Scientia Advanced Research and Reviews*, **11**, 267-286. <https://doi.org/10.30574/msarr.2024.11.1.0091>
- [58] Ijiga, M.O., Olarinoye, H.S., Yeboah, F.A.B. and Okolo, J.N. (2025) Integrating Behavioral Science and Cyber Threat Intelligence (CTI) to Counter Advanced Persistent Threats (APTs) and Reduce Human-Enabled Security Breaches. *International Journal of Scientific Research and Modern Technology*, **4**, 1-15. <https://doi.org/10.38124/ijsrmt.v4i3.376>
- [59] Nwatuze, G.A., Enyejo, L.A. and Umeaku, C. (2025) Enhancing Cloud Data Security Using a Hybrid Encryption Framework Integrating AES, DES, and RC6 with File Splitting and Steganographic Key Management. *International Journal of Innovative Science and Research Technology*, **10**, 1555-1569. <https://doi.org/10.5281/zenodo.14792173>
- [60] Nwatuze, G.A., Ijiga, O.M., Idoko, I.P., Enyejo, L.A. and Ali, E.O. (2025) Design and Evaluation of a User-Centric Cryptographic Model Leveraging Hybrid Algorithms for Secure Cloud Storage and Data Integrity. *American Journal of Innovation in Science and Engineering*, **4**, 49-65. <https://doi.org/10.54536/ajise.v4i2.4482>
- [61] Idogho, C., Abah, E.O., Onuhc, J.O., Harsito, C., Omenkaf, K., Samuel, A., *et al.* (2025) Machine Learning-Based Solar Photovoltaic Power Forecasting for Nigerian Regions. *Energy Science & Engineering*, **13**, 1922-1934. <https://doi.org/10.1002/ese3.70013>
- [62] James, U.U., Ijiga, O.M. and Enyejo, L.A. (2024) AI-Powered Threat Intelligence for Proactive Risk Detection in 5G-Enabled Smart Healthcare Communication Networks. *International Journal of Scientific Research and Modern Technology*, **3**, 125-140. <https://doi.org/10.38124/ijsrmt.v3i11.679>
- [63] Christian, I., Abah, E.O., Abel, J.E., Harsito, C., Omoniyi, M. and Boriwaye, T. (2025) Compatibility Study of Synthesized Materials for Thermal Transport in Thermoelectric Power Generation. *American Journal of Innovation in Science and Engineering*, **4**, 49-63. <https://doi.org/10.54536/ajise.v4i1.3948>
- [64] Idoko, I.P., Ayodele, T.R., Abolarin, S.M. and Ewim, D.R.E. (2023) Maximizing the Cost Effectiveness of Electric Power Generation through the Integration of Distributed Generators: Wind, Hydro and Solar Power. *Bulletin of the National Research Centre*, **47**, Article No. 166. <https://doi.org/10.1186/s42269-023-01125-7>
- [65] Ijiga, A.C., Peace, A.E., Idoko, I.P., Agbo, D.O., Harry, K.D., Ezebuka, C.I. and Ukatu, I.E. (2024) Ethical Considerations in Implementing Generative AI for Healthcare Supply Chain Optimization: A Cross-Country Analysis across India, the United Kingdom, and the United States of America. *International Journal of Biological and Pharmaceutical Sciences Archive*, **7**, 48-63. <https://doi.org/10.53771/ijbpsa.2024.7.1.0015>
- [66] Ayoola, V.B., Idoko, I.P., Eromonsei, S.O., Afolabi, O., Apampa, A.R. and Oyebanji, O.S. (2024) The Role of Big Data and AI in Enhancing Biodiversity Conservation and Resource Management in the USA. *World Journal of Advanced Research and Reviews*, **23**, 1851-1873. <https://doi.org/10.30574/wjarr.2024.23.2.2350>
- [67] Idoko, I.P., Ijiga, O.M., Enyejo, L.A., Akoh, O., Isenyo, G., *et al.* (2024) Integrating

- Superhumans and Synthetic Humans into the Internet of Things (IoT) and Ubiquitous Computing: Emerging Ai Applications and Their Relevance in the U.S. Context. *Global Journal of Engineering and Technology Advances*, **19**, 6-36. <https://doi.org/10.30574/gjeta.2024.19.1.0055>
- [68] Idoko, I.P., Ijiga, O.M., Akoh, O., Agbo, D.O., Ugbane, S.I. and Umama, E.E. (2024) Empowering Sustainable Power Generation: The Vital Role of Power Electronics in California's Renewable Energy Transformation. *World Journal of Advanced Engineering Technology and Sciences*, **11**, 274-293. <https://doi.org/10.30574/wjaets.2024.11.1.0058>
- [69] Idoko, I.P., David-Olusa, A., Badu, S.G., Okereke, E.K., Agaba, J.A. and Bashiru, O. (2024) The Dual Impact of AI and Renewable Energy in Enhancing Medicine for Better Diagnostics, Drug Discovery, and Public Health. *Magna Scientia Advanced Biology and Pharmacy*, **12**, 99-127. <https://doi.org/10.30574/msabp.2024.12.2.0048>
- [70] Idoko, I.P., Eniodunmo, O., Danso, M.O., Bashiru, O., Ijiga, O.M. and Manuel, H.N.N. (2024) Evaluating Benchmark Cheating and the Superiority of MAMBA over Transformers in Bayesian Neural Networks: An In-Depth Analysis of AI Performance. *World Journal of Advanced Engineering Technology and Sciences*, **12**, 372-389. <https://doi.org/10.30574/wjaets.2024.12.1.0254>
- [71] Ijiga, O.M., Idoko, I.P., Enyejo, L.A., Akoh, O., Ugbane, S.I. and Ibokette, A.I. (2024) Harmonizing the Voices of AI: Exploring Generative Music Models, Voice Cloning, and Voice Transfer for Creative Expression. *World Journal of Advanced Engineering Technology and Sciences*, **11**, 372-394. <https://doi.org/10.30574/wjaets.2024.11.1.0072>
- [72] Billboard (2023) Wyclef Jean Discusses Upcoming Reggae Album and Reflects on Career. Billboard. <https://www.billboard.com/music/rb-hip-hop/wyclef-jean-reggae-album-interview-1235722370/>