

Review Article

# A Comprehensive Review of Deepfake and its Detection Techniques

Tatwadarshi P. Nagarhalli<sup>1</sup>, Ashwini Save<sup>2</sup>, Sanket Patil<sup>3</sup>, Uday Aswalekar<sup>4</sup>

<sup>1</sup>Department of Artificial Intelligence & Data Science, Vidyavardhini's College of Engineering and Technology, University of Mumbai, Maharashtra, India.

<sup>2</sup>Computer Engineering Department, VIVA Institute of Technology, University of Mumbai, Maharashtra, India.

<sup>3</sup>Department of Computer Engineering, Vidyavardhini's College of Engineering and Technology, University of Mumbai, Maharashtra, India.

<sup>4</sup>Department of Mechanical Engineering, Vidyavardhini's College of Engineering and Technology, University of Mumbai, Maharashtra, India.

<sup>1</sup>Corresponding Author : [tatwadarshipn@gmail.com](mailto:tatwadarshipn@gmail.com)

Received: 05 June 2024

Revised: 08 July 2024

Accepted: 06 August 2024

Published: 31 August 2024

**Abstract** - Deepfake technology has emerged as a significant concern in the era of digital media, posing threats to various sectors by enabling the creation of highly realistic synthetic content. This paper presents a comprehensive review of deepfake techniques and detection methods. It analyzes 14 research papers covering a range of approaches, including machine learning algorithms, computer vision techniques, and signal processing methods. Key aspects explored include face and voice manipulation, multimodal fusion, and the use of attention mechanisms. The review highlights the challenges in detecting deepfakes, such as dataset bias and the arms race between creators and detectors. Additionally, it discusses the limitations of current detection techniques and the need for robust, scalable solutions. Through a critical analysis of the literature, this review provides insights into the strengths and weaknesses of existing approaches and identifies areas for future research. The paper contributes to understanding deepfake technology and its implications for society, emphasizing the importance of developing effective detection mechanisms to combat the spread of synthetic media.

**Keywords** - Deepfake, Deepfake detection, Face swap, Audio-video manipulation, Deep Learning, Voice spoofing, Synthetic media.

## 1. Introduction

Artificial Intelligence (AI) stands at the forefront of transformative technologies, reshaping industries and revolutionizing human interaction with machines. AI encompasses a spectrum of methodologies, including Machine Learning (ML) [1], Deep Learning (DL) [2], and Natural Language Processing (NLP) [3], each contributing to the advancement of intelligent systems. Machine Learning involves algorithms that enable computers to learn from data and make predictions or decisions without explicit programming [4].

A subclass of machine learning called "Deep Learning" uses multi-layered neural networks to derive hierarchical representations from input, making feature learning and complicated pattern recognition possible [5]. Natural language processing aims to make it possible for computers to comprehend, interpret, and produce human language. This makes jobs like sentiment analysis, text summarization, and language translation easier [6].

Among its myriad applications, Generative AI, a subset of AI that focuses on creating data or content, has garnered substantial attention [7]. This branch of AI encompasses various techniques, including Generative Adversarial Networks (GANs) [8] and Variational Autoencoders (VAEs) [9], enabling machines to generate realistic images, videos, text, and even audio. Businesses have increasingly embraced Generative AI due to its potential to streamline processes, enhance creativity, and personalize user experiences, with projections suggesting significant contributions to the global economy by 2030 [7]. Amidst the rapid advancements in Generative AI, the emergence of deepfake technology has garnered both fascination and concern. Deepfake refers to synthetic media generated by AI algorithms, typically employing Generative AI techniques, to manipulate or replace existing content with fabricated images, videos, or audio [10]. While deepfake technology offers novel avenues for entertainment, artistic expression, and digital content creation, its proliferation raises significant ethical, social, and security concerns.



The allure of deepfake lies in its ability to alter reality, presenting a host of potential benefits convincingly. Content creators can leverage deepfake technology to enhance film visual effects, create lifelike animations, and personalize user experiences in virtual environments [11]. Moreover, deepfake applications extend to medical imaging, where AI-generated simulations aid disease diagnosis and treatment planning. However, the democratization of deepfake tools, coupled with their ease of use, has facilitated the spread of misinformation, cyberbullying, and illicit activities.

Despite the potential benefits, the rapid spread and increasing sophistication of deepfake technology pose significant threats. The menace of deepfake encompasses issues of privacy, identity theft, and trust in digital media [12]. Politicians, celebrities, and ordinary individuals alike face the threat of malicious manipulation, as deepfake videos can be weaponized to disseminate false narratives, incite discord, and undermine democratic processes. In a world increasingly reliant on digital communication and media consumption, the proliferation of deepfake content erodes trust in information sources and exacerbates societal polarization.

Given these challenges, developing robust deepfake detection techniques has become imperative [12]. As deepfake technology evolves, so must our capabilities to discern fabricated content from genuine media. While varied and innovative, current detection methods often face limitations in accuracy, generalizability, and scalability. This paper addresses this research gap by presenting a comprehensive review of existing deepfake detection techniques, shedding light on their strengths, limitations, and avenues for future research and development.

By advancing our understanding of deepfake detection, we aim to fortify defenses against the pernicious effects of synthetic media manipulation and uphold the integrity of digital communication in the AI era. This review is intended to provide a critical resource for researchers and practitioners seeking to develop more effective deepfake detection technologies, ultimately contributing to safeguarding individual integrity, privacy, and public trust.

## 2. Deepfake Generation: Techniques and Technologies

Deepfake generation has become a prominent area of research within the field of Generative AI, utilizing sophisticated algorithms to create highly realistic synthetic media. Among the various techniques employed for deepfake creation, Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs) are two of the most significant methodologies. This section thoroughly explores these techniques, highlighting their mechanisms and applications in generating deepfake content.

### 2.1. Generative Adversarial Networks (GANs)

Since its introduction in 2014 by Ian Goodfellow and associates, generative adversarial networks, or GANs, have completely transformed the area of generative modeling [8]. Two neural networks, a generator and a discriminator, comprise a GAN. They are trained concurrently using adversarial procedures. Whereas the discriminator seeks to discern between genuine and created data, the generator seeks to create realistic synthetic data. Because of this antagonistic connection, the generator keeps producing better data until the discriminator cannot distinguish between actual and phony data.

Because GANs can generate high-fidelity pictures and videos, they are often employed in deepfake production. The quality and control over the generated material have been further improved by techniques like Nvidia's StyleGAN [18], which allow for fine-grained customization of visual features. The versatility of GANs allows for various applications, including face swapping, attribute manipulation, and the creation of entirely synthetic personas.

### 2.2. Variational Autoencoders (VAEs)

As proposed by Kingma and Welling in 2013, Variational Autoencoders (VAEs) are an additional potent method for producing deepfake material [9]. One kind of autoencoder that adds probabilistic components to the encoding and decoding procedures is called a VAE. They comprise a decoder that reconstructs the data from this latent representation and an encoder that maps input data to a latent space represented by a distribution.

Because of their probabilistic character, VAEs enable sampling from the latent distribution to provide unique and varied data samples. VAEs are particularly useful in generating images and videos with controlled variations, making them suitable for applications where variability and creativity are desired. While VAEs generally produce less sharp images than GANs, they offer advantages in interpretability and structured latent spaces.

### 2.3. Combining GANs and VAEs

Recent advancements have explored the combination of GANs and VAEs to leverage the strengths of both techniques. For instance, VAE-GAN models integrate the generative capabilities of VAEs with the adversarial training of GANs, resulting in improved image quality and better latent space representation [35]. These hybrid models enhance the robustness and flexibility of deepfake generation, enabling the creation of more realistic and diverse content. The applications of GANs and VAEs in deepfake generation are vast and varied. In the entertainment industry, these techniques are used to create lifelike visual effects, generate synthetic actors, and develop immersive virtual environments.

In healthcare, AI-generated simulations aid in medical training and diagnosis, providing realistic scenarios for practitioners. However, the misuse of these technologies for malicious purposes, such as spreading misinformation, cyberbullying, and identity theft, underscores the need for ethical considerations and robust detection methods.

So, deepfake generation techniques, particularly GANs and VAEs, have significantly advanced the capabilities of Generative AI. While these technologies offer exciting possibilities for innovation and creativity, they pose substantial ethical and security challenges. Understanding the mechanisms and applications of GANs and VAEs is crucial for developing effective countermeasures against the malicious use of deepfake technology.

### 3. Types of Deepfake

Deepfakes encompass a variety of techniques for creating synthetic media based on the type of data being manipulated, and the deepfakes can be categorized into the following types as well.

#### 3.1. Face Swaps

One common kind of deepfake is face swapping, which is when someone's face is seamlessly replaced with another's in a picture or video. Usually, advanced deep learning models—especially GANs, which are excellent at producing realistic synthetic data—are used to implement this strategy. GANs are made up of a generator and a discriminator neural network that have been trained adversarially to create excellent false pictures that are identical to genuine ones. [8].

The process of face swapping begins with collecting a dataset of images containing both source and target faces. These images are then fed into the GAN, which learns to map the facial features of the target individual onto the source individual's face while preserving other characteristics such as head pose, lighting conditions, and facial expressions. This mapping process is iteratively refined through training until the generated face appears convincingly realistic.

While face swap deepfakes can be used for entertainment purposes, such as inserting a person's face into a movie scene or creating humorous videos, they also raise significant concerns about identity theft, privacy infringement, and misinformation. Malicious actors could exploit face swap technology to impersonate individuals, spread false information, or manipulate visual content for fraud, highlighting the urgent need for robust detection and mitigation strategies in this domain.

#### 3.2. Voice Cloning

Voice cloning deepfakes represent a significant advancement in synthetic media generation, utilizing deep learning algorithms to replicate a person's voice with

remarkable accuracy. These algorithms are typically based on techniques such as WaveNet [13] or Tacotron [14], which can generate highly realistic audio waveforms. By training on a large dataset of audio samples containing the target individual's voice, these algorithms learn to capture the subtle nuances of speech patterns, intonation, and timbre. This enables them to synthesize speech that closely resembles the original speaker.

One of the key challenges in voice cloning is capturing the speaker's unique characteristics, including accent, pitch, and cadence. Research by Google's DeepMind team introduced WaveNet, a generative model for raw audio waveforms, demonstrating impressive results in generating natural-sounding speech [13]. Similarly, Tacotron, developed by researchers at Google, is a sequence-to-sequence model capable of directly synthesizing speech from text inputs, achieving high-quality speech synthesis with natural-sounding prosody [14].

Voice cloning deepfakes have significant implications for various applications, including impersonation, fraud, and generating fake audio recordings. For instance, malicious actors could exploit voice cloning technology to impersonate individuals for fraud, such as phishing scams or social engineering attacks. Furthermore, because voice cloning deepfakes may be used to produce false or misleading audio recordings, they raise questions about the veracity and authenticity of audio information.

#### 3.3. Body Swaps

Body swaps represent a significant extension of the deepfake technology beyond facial manipulation, allowing for the replacement of an entire person's body in a video with that of another individual. This process typically involves employing techniques similar to those used in face swapping but applied to the entire body instead. By leveraging deep learning algorithms and computer vision methodologies, body swap deepfakes can generate highly realistic simulations of individuals performing actions or engaging in activities that they did not actually partake in.

The development of body swap deepfake technology has been facilitated by advancements in deep learning and GANs. For instance, research by Nirkin et al. [15] introduced a method for generating high-quality full-body deepfakes by leveraging a GAN-based architecture capable of synthesizing realistic human poses and appearances. Additionally, Li et al. [16] proposed a deep learning framework for seamless body swapping in videos, enabling the transfer of body movements and expressions from a source actor to a target actor while preserving visual realism. Furthermore, the proliferation of body swap deepfakes has raised concerns about their potential misuse and implications for various applications, including entertainment, advertising, and misinformation.

For example, malicious actors could exploit body swap deepfake technology to create convincing simulations of individuals engaging in inappropriate or unethical behaviour, leading to reputational damage or misinformation dissemination. Therefore, there is a growing need for research and development of detection techniques to identify and mitigate the spread of body swap deepfakes in online platforms and media.

### 3.4. Text-Based Deepfake

Text-based deepfakes, a burgeoning facet of synthetic media, entail the generation of artificial text that emulates the writing style of a specific individual. This process relies on sophisticated Natural Language Processing (NLP) techniques and deep learning models trained on extensive corpora of text samples attributed to the target individual. By analyzing and understanding the nuances of the individual's writing style, including vocabulary, grammar, syntax, and tone, these models can generate synthetic text closely resembling authentic writings.

Advances in deep learning architectures, namely in transformer models and Recurrent Neural Networks (RNNs), have been the driving force behind the creation of text-based deepfake technologies. For example, OpenAI's Generative Pre-trained Transformer (GPT) model series has shown an impressive capacity to produce logical and contextually appropriate text through extensive pre-training on various text sources [17]. These models excel at capturing the stylistic nuances of individual writers, allowing for the creation of convincing text-based deepfakes.

Moreover, the proliferation of text-based deepfakes has raised significant concerns regarding their potential for misuse, particularly in disseminating fake news, misinformation, and social engineering attacks. For example, malicious actors could exploit text-based deepfake technology to fabricate news articles, social media posts, or emails that appear to originate from reputable sources or influential individuals, thereby manipulating public opinion or deceiving individuals for nefarious purposes. Consequently, there is a growing need for research and development of detection techniques to identify and combat the spread of text-based deepfakes in online platforms and communication channels.

### 3.5. Image Generation

Image generation using deep learning models, particularly GANs, has emerged as a powerful tool for synthesizing realistic images of people, animals, objects, and scenes that do not exist. GANs consist of two neural networks, a generator and a discriminator, which are trained adversarially to generate increasingly indistinguishable images from real photographs. This process involves the generator network creating candidate images while the discriminator network attempts to differentiate between real and fake images. Through iterative training, GANs learn to

produce high-quality and realistic images, capturing intricate details and visual characteristics.

The seminal work by Goodfellow et al. [8] introduced GANs as a novel approach to generative modeling, demonstrating their ability to generate compelling images across various domains. Subsequent research has further advanced GAN-based image generation techniques, leading to the development of models capable of producing highly realistic and diverse visual content. For instance, Karras et al. [18] proposed StyleGAN, a progressive growing GAN architecture that synthesises high-resolution images with unprecedented realism and variation.

However, the proliferation of image generation technology has raised concerns about its potential for misuse, particularly in creating deceptive or misleading visual content. Malicious actors could exploit GANs to fabricate images for various purposes, including propaganda, disinformation, and digital manipulation. Therefore, research and development of detection techniques are desperately needed to recognize and stop the propagation of deceptive visual material produced by deep learning models.

### 3.6. Video Manipulation

Video manipulation deepfakes represent a sophisticated application of deep learning techniques to alter videos and manipulate the facial expressions, gestures, or actions of individuals depicted in the footage. This process involves analyzing and modifying video frames using advanced computer vision and deep learning algorithms to achieve the desired effects. Deep learning models, such as Convolutional Neural Networks (CNNs) and RNNs, are trained on large datasets of video sequences to learn the temporal and spatial dependencies within the data, enabling them to generate realistic and coherent video manipulations.

Research by Suwajanakorn et al. [19] introduced a method for generating highly realistic video manipulations using a deep learning approach known as "Synthesizing Obama." The method leveraged a generative neural network to learn the facial movements and expressions of former President Barack Obama from existing video footage and then applied these learned characteristics to manipulate the facial expressions of Obama in target videos. This groundbreaking work demonstrated the potential of deep learning techniques to produce convincing video manipulations that are difficult to distinguish from authentic footage.

Furthermore, the proliferation of video manipulation deepfakes has raised concerns about their potential for misuse in various contexts, including entertainment, satire, and propaganda. Deepfakes for video manipulation can be used for good or benign purposes, such as making humorous satire or entertaining material. However, they also carry much danger when it comes to disinformation, fraud, and public opinion

manipulation. Therefore, to detect and stop the spread of harmful or deceptive video modifications in online platforms and media, there is an urgent need for the study and development of detection systems.

### 3.7. Audio-Visual Deepfakes

Audio-visual deepfakes represent a sophisticated integration of techniques from both voice cloning and face swaps, aiming to synchronize manipulated facial movements with cloned speech, thereby creating highly convincing simulations of individuals speaking. This process involves leveraging deep learning algorithms to analyze and modify both the facial expressions captured in video frames and the corresponding audio content to ensure coherence and synchronization between the two modalities.

By seamlessly combining synthesized audio with manipulated facial movements, audio-visual deepfakes produce multimedia content where the speaker's lip movements and facial expressions closely match the synthesized speech, resulting in persuasive simulations of individuals delivering speeches, interviews, or other multimedia content.

Research by Zhou et al. [20] introduced a method for generating audio-visual deepfakes using a conditional GAN architecture, enabling the synthesis of synchronized facial expressions and cloned speech. The proposed approach leveraged paired audio-visual data to train the cGAN model, enabling it to learn the intricate correlations between speech content and facial movements. By conditioning the generator network on input audio features, the model could generate realistic facial expressions that align with the synthesized speech, achieving highly convincing audio-visual synchronization.

However, the widespread availability of audio-visual deepfake technology has raised significant concerns about its potential for misuse, particularly in creating fake interviews, speeches, or other multimedia content to deceive or manipulate viewers. While audio-visual deepfakes can have legitimate applications in entertainment and visual effects, their misuse poses risks to trust, authenticity, and the integrity of multimedia content. Therefore, there is an urgent need for research and development of robust detection techniques to identify and mitigate the spread of misleading or malicious audio-visual deepfakes in online platforms and media.

## 4. Deepfake Benchmark Datasets

Deepfake detection model development and assessment heavily rely on benchmark datasets. These datasets offer the information required to assess algorithms' performance, train new techniques, and compare them to industry standards. This section examines and highlights the qualities, contributions, and limits of some of the most popular datasets in deepfake detection.

### 4.1. FaceForensics++

One of the most extensive and popular datasets for deepfake detection is FaceForensics++. It is made up of more than a thousand films that have been altered using four distinct methods: NeuralTextures, Face2Face, Deepfakes, and FaceSwap [21]. The dataset offers a strong basis for training and assessing detection algorithms, containing both the altered films and the matching original videos. Furthermore, FaceForensics++ provides movies at varying degrees of compression, accurately capturing the diverse quality of videos seen in real-world situations.

### 4.2. Deepfake Detection Challenge Dataset (DFDC)

Facebook launched the Deepfake Detection Challenge (DFDC) dataset in association with many business partners and academic organizations [36]. It has over 100,000 video clips, combining authentic and deepfake films using cutting-edge generating techniques. This dataset is noteworthy for its size and diversity since it offers substantial data required for building reliable models. In order to facilitate the thorough examination and benchmarking of detection methods, the DFDC dataset additionally contains metadata and annotations.

### 4.3. Celeb-DF

A dataset called Celeb-DF solves some of the shortcomings noted in previous datasets, including the authenticity and caliber of deepfake films [37]. It includes 5,639 deepfake films that correlate to 590 genuine celebrity videos, produced with an enhanced Deepfake synthesis technique. Celeb-DF is a tough dataset for detection algorithms and offers a realistic baseline for assessing their performance since it emphasizes good visual quality and minimal artifacts.

### 4.4. DeeperForensics-1.0

DeeperForensics-1.0 is another significant dataset containing 50,000 videos generated using various techniques and tools [38]. This dataset is notable for its focus on diverse perturbations, such as different lighting conditions, compression levels, and occlusions, reflecting the complex nature of real-world deepfakes. DeeperForensics-1.0 also includes many subjects and scenes, enhancing its applicability for training generalizable detection models.

### 4.5. Google/Jigsaw Deepfake Detection

Google and Jigsaw released a deepfake detection dataset to support research in this area. The dataset includes thousands of videos created using various face-swapping algorithms [39]. This dataset was designed to provide researchers with a resource for training and evaluating models capable of detecting AI-manipulated content, contributing to improving detection technologies.

### 4.6. WildDeepfake

WildDeepfake is a dataset specifically curated to include deepfake videos in the wild, reflecting real-world conditions

more accurately than synthetic datasets [40]. It comprises 7,314 real videos and 3,509 deepfake videos, with annotations for each. This dataset challenges detection models with videos that have undergone diverse post-processing steps, such as compression and noise addition, which are common in videos distributed online.

## 5. Deepfake Detection Techniques

Among the different deepfake techniques, research on detection, face swap, voice spoofing, and a combination of all types of deepfakes have been done extensively. Research by Rossler et al. [21] introduced FaceForensics++, a deep learning-based approach for detecting manipulated facial images, including face swaps. Their method utilizes CNNs to extract subtle artifacts and inconsistencies introduced during face-swapping, such as misaligned facial features, unnatural lighting, or inconsistent skin tones. Their model achieved state-of-the-art performance in detecting manipulated faces by training on a large dataset of both real and manipulated images.

The growing problem of digital face modification, which casts doubt on the veracity of media output, is discussed in the study [22]. The suggested framework classifies modified face photos using deep learning methods, namely the EfficientNet learning model. Using several modification techniques, including Face-Swap, Face-2-Face, Deepfakes, and neural textures, the framework offers an all-encompassing method for detecting altered areas in facial image data. Furthermore, by using BlazeFace tracking, it becomes easier to locate the pixel coordinates and face characteristics, which improves the precision of manipulation detection.

In terms of accuracy and efficiency, the article promises greater performance when compared to existing methodologies. However, more empirical validation and benchmarking against well-established datasets are required to support these findings. Although the suggested approach shows promise in identifying digital face video tampering in forensics and security, its implementation may be impacted by computing resources and scalability issues. Overall, the study makes a substantial addition to the subject, although more investigation is required to resolve any issues and guarantee the framework's applicability in actual situations.

The paper [23] addresses the critical issue of fake face forgery, which poses significant security concerns across various domains, such as fake news dissemination, fraud, and impersonation. While existing face forgery detection methods have demonstrated success within specific domains, they often lack generalization capability and experience significant performance degradation when applied to new, unforeseen domains. To overcome this challenge, the paper proposes a novel approach based on the Vision Transformer (ViT) architecture to enhance the generalizability of fake face detection models.

The suggested approach uses pretrained ViT weights, updating the Low-Rank Adaptation (LoRA) modules as it goes through training. By using this technique, the model may leverage the information embedded in the pretrained weights and adapt to new domains. Additionally, by offering more supervision during training, the adoption of Single Center Loss (SCL) substantially improves the model's capacity for generalization. Therefore, the suggested approach successfully addresses the generalization problem by achieving state-of-the-art detection performance in both cross-manipulation and cross-dataset assessments.

The study [24] investigates the possibilities of crowd-based decision fusion techniques in order to tackle the problem of digital face alteration detection. The study tries to improve detection accuracy by utilizing the choices made by human examiners and combining variables like decision time, expertise level, and decision confidence. In psychophysical evaluation trials, 223 participants' choices were blended to imitate crowds of up to seven human examiners using various modification methods, including face morphing, swapping, and retouching.

The findings show that decision fusion considerably increases accuracy, whereas individual human examiners may only achieve mediocre detection performance. In particular, the most competitive detection performance is obtained using a weighted fusion technique that considers the examiners' judgment confidence. This method shows how crowd-powered detection may help overcome the limits of the skills of individual examiners.

Implementing a crowd-powered system in real-world circumstances may present certain obstacles, such as scalability, dependability, and resource needs. Furthermore, subjectivity and variability are introduced by using human examiners, which may impact the consistency and dependability of detection results.

The study [25] discusses the growing threat posed by face-forging techniques, which have allowed fake face recordings to be widely shared online and jeopardize the security and reliability of digital assets. This paper suggests using 3D CNNs for video-level face forgery detection, whereas previous approaches mostly depend on 2D CNNs for identifying forged frames inside movies. The suggested 3D attention network can extract spatial and temporal data from video sequences by integrating a lightweight attention module into the network designs.

A noteworthy feature of the suggested method is the attention module's creation of attention maps, which highlight fabricated areas inside fictitious faces. Furthermore, implementing a global attention pool improves detection accuracy by reducing disparities across various locations. The model outperforms previous techniques in experimental

evaluations using the FaceForensics++ dataset, demonstrating high accuracy in forgery detection. Additionally, the model's great transferability and generalization ability across several datasets is confirmed by cross-dataset assessment on the Celeb-DF dataset. However, there is little talk of pragmatic issues like model deployment and difficulties with real-world

implementation. Table 1 compares the five papers based on the problem identified, techniques and technology used, accuracy, and practicality. Each paper's strengths and weaknesses are highlighted, allowing for a comprehensive understanding of their contributions and limitations in face forgery detection.

**Table 1. Analysis of face swap detection techniques**

Paper	Problem Identified	Techniques & Technology Used	Practicality	Pros	Cons
[21]	Threat of manipulated facial images and face swaps	CNNs for artifact extraction	Challenges related to scalability and real-world implementation	State-of-the-art performance in detecting manipulated faces, utilization of CNNs for artifact extraction.	Potential challenges in scalability and deployment.
[22]	Lack of generalization capability in detection methods	Crowd-powered decision fusion, considering confidence levels, experience, and decision time	Potential challenges in scalability and reliability	Exploration of crowd-powered detection, improved accuracy through fusion.	Reliance on human examiners introduces subjectivity and variability.
[23]	Lack of generalization capability in existing methods	Vision Transformer (ViT) architecture, Low-Rank Adaptation (LoRA) modules, Single Center Loss (SCL)	Challenges in computational resources and scalability	Utilization of ViT architecture, incorporating LoRA modules and SCL, insights into leveraging human examiners' decisions for detection enhancement.	Limited empirical validation in real-world scenarios.
[24]	Threat of widespread dissemination of synthetic face videos	3D CNNs, attention module, global attention pool	Potential computational resource requirements and scalability challenges	Introduction of 3D CNNs for video-level detection, incorporation of attention mechanism, strong transferability and generalization ability.	Limited discussion on practical considerations such as deployment and real-world implementation challenges.
[25]	Threat to security and trustworthiness of digital content	3D CNN, attention module, decision fusion, spatial and temporal feature extraction	Challenges related to computational resources and scalability	A comprehensive approach to manipulation detection, utilization of 3D CNNs and attention mechanism, and potential for forensic and security applications.	Limited empirical validation of claimed performance.

The study [26] addresses the pressing issue of voice spoofing attacks, which have increased tenfold during the past few years. These attacks employ Automatic Speaker Verification (ASV) systems to deceive them into allowing unauthorized users to access confidential information, such as bank account and home control access, using fake voice.

The availability of sophisticated tools has made it easier for attackers to carry out voice spoofing attacks, which poses major security concerns. In order to mitigate the challenges caused by synthetic speech that circumvents ASV system

security, the research proposes an effective synthetic speech detector that utilizes a fusion of spectral properties. The proposed architecture especially employs a fused feature vector composed of Spectral Flux, Spectral Centroid, Mel-Frequency Cepstral Coefficients (MFCC), and Gammatone Cepstral Coefficients (GTCC) to describe audio signals. These features search for both the computational aberrations characteristic of artificial signals and the natural speech fluctuation characteristics of genuine signals. The framework trains a Bidirectional Long Short-Term Memory (BiLSTM) network to discriminate between real and fraudulent signals

by employing these features. This enables the system to detect attacks, including voice conversion and synthetic speech.

The efficacy of the suggested approach in identifying logical access attacks, such as voice conversion and cloned/synthetic voice assaults, is assessed using the ASVspoof 2019 LA dataset. Although the article offers a promising countermeasure to the increasingly dangerous voice spoofing issue, more empirical validation and benchmarking against a wider range of datasets are required to evaluate its accuracy fully. Other practical factors like scalability and deployment feasibility in actual ASV systems should be considered to guarantee the framework's practical viability.

Hafiz Malik discusses the rising threat that voice cloning technologies represent to speech-based access control systems and voice-driven interfaces [27]. Because these technologies may produce virtually indistinguishable speech from real audio samples, security and privacy issues have been highlighted despite being useful for tailored voice interfaces and other applications. Modern speech synthesis methods make use of trained generative models, which can introduce distinctive distortions into synthetic speech and make it difficult to tell the difference between real and fake sounds.

The research suggests a unique method for capturing characteristics that distinguish real audio from cloned audio by using higher-order spectrum analysis to reduce these security threats. In particular, the technique finds generative model artifacts in the cloned speech by using Quadrature Phase Coupling (QPC) in the estimated bicoherence, Gaussianity test statistics, and linearity test statistics. The efficiency of the suggested strategy is demonstrated by experimental assessments on a dataset that includes both real and cloned speech samples with near-perfect detection rates.

Independent analysis of voice cloning is one important area of deepfake and deepfake detection, seen in active research recently. Some significant deepfake detection techniques with the potential for largescale, industrywide implementation have been reviewed.

Aakriti Aggarwal et al. [28] address the pressing challenge of identifying and mitigating the spread of falsified information, particularly through manipulating videos using deep fake technology. The paper highlights the exponential growth of digital media consumption and the potential for malicious actors to exploit this medium to spread misinformation, especially during sensitive periods such as elections. The paper proposes a model for detecting deep fake videos using a combination of XceptionNet and ResNet50 architectures, augmented with multiple hidden layers of neural networks. These deep learning techniques reflect the contemporary trend in leveraging artificial intelligence to address complex problems like deep fake detection.

One notable strength of the proposed approach is its reported performance in outperforming other state-of-the-art methods in terms of precision and recall rates. By employing advanced neural network architectures and leveraging multiple layers of abstraction, the model demonstrates promising accuracy in distinguishing between authentic and manipulated videos. Moreover, the practicality of the proposed method is underscored by its ability to efficiently predict the authenticity of a given video, which is crucial in the context of combating the proliferation of deep fake content.

However, while the paper presents a promising solution to the challenge of deep fake detection, there are potential limitations to consider. Firstly, the abstract lacks specific details regarding the dataset used for training and evaluation and the methodology for collecting and annotating deep fake videos. Additionally, the practical deployment of the proposed model may face challenges related to computational resources and scalability, particularly in real-time or large-scale video analysis scenarios. Nonetheless, with further refinement and validation, the IsSwap model holds potential as a valuable tool in the fight against deep fake proliferation.

The study [29] discusses the increasing demand for reliable deepfake detection algorithms due to the spread of face synthesis technologies and related security issues. The abstract draws attention to the shortcomings of current methods, especially regarding their inability to efficiently capture spatially important frequency aspects crucial for identifying more realistic counterfeit material. The research suggests a unique Spatial-Frequency Dynamic Graph technique to address this issue by combining multiscale attention map learning, dynamic graph spatial-frequency feature fusion, and content-guided adaptive frequency extraction.

The complete methodology of the suggested solution, which integrates several cutting-edge elements to utilize relation-aware characteristics in both the spatial and frequency domains, is one of its most significant strengths. By applying dynamic graph learning approaches, the model improves generalization performance across different types of deepfakes by strengthening the relationship between frequency characteristics and picture content. Furthermore, the published experimental findings highlight the usefulness of the suggested strategy in deepfake detection by showing notable performance advantages over state-of-the-art techniques across many benchmark datasets.

Although the research offers a promising solution to the stated issue, there are some possible drawbacks. Specifically, the efficacy of the suggested strategy may differ based on the variety and intricacy of deepfake techniques that are actually used.



The prevalent problem of false information and fake media on the internet, especially in textual and visual forms, is discussed in the study [30]. The study draws attention to the shortcomings of current detection techniques, which mainly concentrate on binary classification and single-modality forgeries, missing minute manipulation traces across many modalities. In order to detect the authenticity of multi-modal media while grounding manipulated content—such as image bounding boxes and text tokens—through deeper reasoning of multi-modal manipulation, the paper introduces a novel research problem called Detecting and Grounding Multi-Modal Media Manipulation (DGM4).

Creating the initial DGM4 dataset, which comprises altered image-text pairings annotated with various alterations, is a noteworthy strength of the suggested methodology that allows for extensive research and benchmarking. Furthermore, the study presents the Hierarchical Multi-modal Manipulation Reasoning Transformer (HAMMER), which captures fine-grained interaction across several modalities by utilizing modality-aware cross-attention mechanisms and manipulation-aware contrastive learning. Extensive investigations show that HAMMER performs better in detecting and grounding multi-modal media manipulation by merging manipulation detection and grounding heads at both superficial and deep levels based on interacting multi-modal information.

However, while the paper presents a promising solution to the identified problem, potential limitations exist. The effectiveness of DGM4 and HAMMER may vary depending on the diversity and complexity of multi-modal manipulation techniques encountered in practice. Nonetheless, with further validation, the proposed approach holds promise as a valuable tool in combating multi-modal media manipulation and advancing research in this area.

By utilizing Domain Generalization (DG) approaches, the study [31] tackles the problem of enhancing Face Anti-Spoofing (FAS) systems' generalization in unknown settings. The study draws attention to the shortcomings of earlier approaches that align domain distributions using imprecise and arbitrary domain labels, which results in an insufficient depiction of actual domain distributions. Furthermore, these approaches neglect fine-grained instance-level distinctions that may impact generalization in favor of domain-level alignment. The study suggests a unique method that aligns features at the instance level without requiring domain labels to get around these problems.

The novel viewpoint on DG FAS that prioritizes instance-level feature alignment to reduce susceptibility to instance-specific styles is one of the approach's standout strengths. The proposed Instance-Aware Domain Generalization framework introduces techniques like Asymmetric Instance Adaptive Whitening, Dynamic Kernel Generator, and Categorical Style

Assembly to adaptively remove style-sensitive feature correlations and produce style-diversified features with significant style shifts. The suggested approach outperforms its cutting-edge rivals, as evidenced by the experimental findings, which also show enhanced accuracy and resilience in the face of spoofing challenges.

While the paper presents a promising solution to the identified problem, potential limitations exist. The paper does not give any specific details regarding the proposed framework's computational complexity and resource requirements, which may impact its practical deployment in real-world scenarios.

The study [32] discusses the problem of deepfake films, which is becoming increasingly common, and the critical need for sophisticated systems for identification and prevention. The authors suggest a novel method for processing audio-visual data in real-time via a web interface—more precisely, by using a browser plugin—by applying Deep Learning technology. Using a multimodal neural network that combines visual and auditory information from videos, a thorough approach to deepfake prediction is presented, which may improve the detection of altered material.

The suggested system's ability to attain a maximum validation accuracy of 90%, which indicates strong performance in identifying deepfake films, is one of its main strengths. Furthermore, practicality and usability are improved by the use of JavaScript to construct an end-to-end solution as a Chrome extension and by creating a responsive, low-latency API. However, potential limitations may arise concerning the scalability and generalizability of the system across different types of deepfake manipulations and varying environmental conditions. Additionally, the paper would benefit from further elucidation on the methodology employed for feature extraction and model training and the validation process used to assess accuracy.

The proposed Audio-Visual Deepfake Detection System presents a promising contribution to advancing deepfake detection and mitigation technologies. While its achievements in accuracy and practicality are commendable, continued refinement and validation are necessary to address potential limitations and ensure effectiveness across diverse scenarios and environments.

The difficulty of identifying edited videos is addressed in the research [33] by taking advantage of minute discrepancies between the visual and auditory inputs. The suggested approach uses anomaly detection techniques and only utilizes genuine, unlabeled data to train an autoregressive model. By creating audio-visual feature sequences that record the temporal synchronization between sound and video frames, the model can identify films that have low probability assignments due to possible tampering.

One notable strength of the approach is its self-supervised learning framework, which eliminates the need for labelled data, making it more scalable and adaptable to various scenarios. Additionally, the emphasis on detecting manipulated speech videos highlights the relevance of the method in addressing specific types of deepfake content. However, potential limitations may arise concerning the model's generalizability across different manipulations beyond speech videos. Also, the paper could benefit from further elaboration on the specific anomaly detection techniques employed and their effectiveness in capturing subtle inconsistencies between audio and visual signals.

So, the proposed self-supervised video forensics method presents a promising approach to detecting manipulated videos by leveraging audio-visual anomaly detection. While its self-supervised learning framework and focus on specific types of manipulations are commendable, further research and validation are needed to assess its effectiveness across a broader range of manipulation techniques and real-world scenarios.

The paper [34] addresses the critical challenge of identifying deepfake videos, which pose a significant threat due to their potential for spreading misinformation and manipulating public opinion. While existing research has

predominantly focused on visual deepfake detection methods, this work highlights the overlooked aspect of audio manipulation, such as synthetic speech, in creating deepfakes.

By introducing a novel joint detection task that considers both visual and auditory cues, the proposed framework offers a holistic approach to deepfake detection. Leveraging the inherent synchronization between visual and auditory modalities, the model demonstrates improved performance compared to independently trained models. This approach enhances detection accuracy and exhibits superior generalization capability, particularly in identifying previously unseen types of deepfakes.

While the framework shows promise in addressing the multifaceted nature of deepfake manipulation, potential challenges may arise regarding computational complexity and resource requirements for simultaneous analysis of both modalities. Nonetheless, the paper sheds light on the importance of integrating audio-visual cues for more robust deepfake detection, contributing valuable insights to the ongoing efforts to combat deceptive media content proliferation. Further research to refine the approach and mitigate computational constraints will be crucial for practical implementation in real-world scenarios. Table 2 shows the analysis of the deepfake detection techniques.

**Table 2. Analysis of deepfake detection techniques**

Paper	Problem Identified	Techniques & Technology Used	Accuracy	Practicality
[26]	Addresses the challenge of detecting synthetic speech attacks	Presents a synthetic speech detector that combines spectrum parameters such as spectral flux, spectral centroid, MFCC, and GTCC.	Effective detection of voice conversion and synthetic speech attacks.	Offers a reliable method for detecting voice spoofing attacks with practical applications.
[27]	Highlights the necessity of protecting voice-driven interfaces from attacks using cloned audio.	Suggests using higher-order spectrum analysis as the basis for recording generative model artifacts from cloned audio.	Ability to identify real and copied audio with nearly flawless detection rates.	Provides a robust approach for detecting cloned audio attacks, reducing false alarms and enabling large-scale analysis.
[28]	Addresses the challenge of detecting deepfake videos	The proposed method combines XceptionNet, ResNet50, and multiple hidden layers of neural networks.	Demonstrated superior performance in terms of precision and recall rate	Provides a fast and reliable method for determining video authenticity, suitable for practical applications.
[29]	Addresses the need for improved deepfake detection methods	Introduces the Multiple Domains Attention Map Learning, Content-guided Adaptive Frequency Extraction, and Dynamic Graph Spatial-Frequency Feature Fusion Network in the Spatial-Frequency Dynamic Graph technique context.	Sustainedly exceeds state-of-the-art performance in experiments.	Enhances detection accuracy and is suitable for various scenarios.
[30]	Identifies the challenge of detecting and	Proposes Hierarchical Multi-Modal Manipulation Reasoning	Demonstrates superiority over	Offers robust detection and localization of

	grounding multi-modal media manipulation	tRansformer (HAMMER) framework, focusing on manipulation detection and grounding. Includes creation of a benchmark dataset.	state-of-the-art methods.	manipulations, facilitating large-scale analysis.
[31]	Addresses the need for improved face anti-spoofing methods	Presents the Instance-Aware Domain Generalization system, which includes the Dynamic Kernel Generator, Categorical Style Assembly, and Asymmetric Instance Adaptive Whitening.	Outperforms state-of-the-art competitors	Provides a robust method for face anti-spoofing with strong performance and generalization capabilities.
[32]	Tackles the challenge of detecting audio-visual deepfakes	Uses a multimodal neural network that is fed visual and auditory information in order to forecast deepfakes.	90% validation accuracy	It offers a responsive, low-latency solution implemented as a browser plugin, contributing to the fight against misinformation.
[33]	Addresses the issue of detecting manipulated speech videos	Proposes an autoregressive model trained on real, unlabeled data for audio-visual anomaly detection.	Demonstrates strong performance on the task	Offers a self-supervised approach for video forensics with potential for real-world application.
[34]	Identifies the need for joint detection of visual and auditory deepfakes	Introduces a joint detection task exploiting synchronization between visual and auditory modalities.	Outperforms independently trained models	Enhances deepfake detection accuracy by considering both visual and auditory cues simultaneously.

The fourteen papers comprehensively explore techniques and methodologies to detect various forms of digital media manipulation, ranging from deepfake videos to synthetic speech attacks. Leveraging advanced deep learning architectures such as Vision Transformers (ViT), CNNs, and multimodal neural networks, these studies propose innovative solutions to address the growing threat of manipulated media. Techniques like content-guided spatial-frequency relation reasoning, instance-aware domain generalization, and joint audio-visual deepfake detection demonstrate a nuanced understanding of the intricacies of identifying sophisticated manipulations across different modalities.

Moreover, the incorporation of self-supervised video forensics, crowd-powered decision fusion, and anomaly detection in audio-visual data showcase the field's interdisciplinary nature, where insights from computer vision, signal processing, and human-computer interaction converge to tackle the challenges posed by digital media manipulation. By exploring methods for both detection and localization of manipulations, these papers contribute significantly to developing robust and reliable forensic tools essential for maintaining the integrity of digital content.

However, despite the strides made in algorithmic advancements and empirical evaluations demonstrating promising results, several challenges remain. Issues related to scalability, computational resource requirements, and real-world deployment pose practical hurdles to the widespread

adoption of these techniques. Additionally, the reliance on labeled datasets and the potential for adversarial attacks highlight the need for further research to enhance detection systems' resilience and generalization capabilities. Overall, the collective efforts showcased in these papers underscore the urgency and complexity of the task at hand and provide valuable insights into the ongoing efforts to combat digital media manipulation effectively.

## 6. Conclusion

In conclusion, this comprehensive review has delved into the multifaceted realm of deepfake technology and the evolving landscape of detection methodologies. Deepfakes, propelled by advanced machine learning algorithms such as Generative Adversarial Networks (GANs) and Convolutional Neural Networks (CNNs), pose a significant challenge in today's digital era, with the potential to disseminate misinformation, violate privacy, and compromise security. However, the concerted efforts of researchers have yielded diverse, innovative techniques aimed at identifying and mitigating the adverse effects of deepfakes.

The papers examined in this review exemplify the interdisciplinary nature of deepfake detection research, drawing upon techniques from computer vision, signal processing, and machine learning. From utilising attention mechanisms and dynamic graph learning to explore multimodal fusion strategies and crowd-powered decision fusion, these methodologies showcase the breadth and depth

of approaches employed in the fight against deepfake proliferation.

Looking ahead, several critical challenges and future directions emerge on the horizon of deepfake detection research. Firstly, there is an urgent need to develop more robust and scalable detection methodologies capable of adapting to the evolving landscape of deepfake generation techniques. Efforts should focus on enhancing the generalization capabilities of detection systems to combat adversarial manipulations effectively. This entails exploring novel architectures, such as graph neural networks and reinforcement learning-based approaches, to improve detection accuracy and resilience.

Furthermore, addressing the issue of data scarcity and dataset bias is paramount for advancing the field of deepfake detection. Researchers should explore techniques for generating synthetic data and developing transfer learning frameworks to mitigate the reliance on annotated datasets. Moreover, establishing standardized evaluation benchmarks and fostering interdisciplinary collaboration will facilitate

knowledge sharing and benchmarking performance across different detection methodologies.

Beyond technical advancements, education and awareness initiatives are essential for empowering users to identify and mitigate the risks posed by deepfake technology. By promoting media literacy and critical thinking skills, individuals can become more adept at discerning between authentic and manipulated content, thereby reducing the potential impact of deepfake dissemination.

In summary, while the proliferation of deepfake technology presents formidable challenges, the collective efforts of the research community offer promise for effective detection and mitigation strategies. By embracing interdisciplinary collaboration, fostering innovation, and prioritizing education and awareness, we can work towards safeguarding the integrity of digital content and preserving trust in the digital ecosystem. As we navigate the complexities of the deepfake landscape, continued vigilance, collaboration, and innovation will be indispensable in addressing this evolving threat.

## References

- [1] Manju M. et al., "Smart Fields: Enhancing Agriculture with Machine Learning," *2024 2<sup>nd</sup> International Conference on Artificial Intelligence and Machine Learning Applications Theme: Healthcare and Internet of Things (AIMLA)*, Namakkal, India, pp. 1-5, 2024. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [2] Harini B. et al., "Advanced Sound Detection and Behavior Examination for Real-Time Intruder Detection Using Deep Learning: A Comprehensive Security Framework," *2024 2<sup>nd</sup> International Conference on Artificial Intelligence and Machine Learning Applications Theme: Healthcare and Internet of Things (AIMLA)*, Namakkal, India, pp. 1-6, 2024. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [3] Smriti Sett, and Ajay Vikram Singh, "Applying Natural Language Processing in Healthcare Using Data Science," *2024 11<sup>th</sup> International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions) (ICRITO)*, Noida, India, pp. 1-6, 2024. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [4] Tom M. Mitchel, *Machine Learning*, McGraw-Hill Science / Engineering / Math, 1997. [[Google Scholar](#)]
- [5] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton, "Deep Learning," *Review Articles*, vol. 521, pp. 436-444, 2015. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [6] Daniel Jurafsky, and James H. Martin, "Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics and Speech Recognition," *Computational Linguistics*, vol. 26, no. 4, pp. 638-641, 2000. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [7] Jacques Bughin et al., "Artificial Intelligence The Next Digital Frontier?," McKinsey & Company, 2017. [[Google Scholar](#)] [[Publisher Link](#)]
- [8] Ian J. Goodfellow et al., "Generative Adversarial Networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139-144, 2020. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [9] Diederik P. Kingma, and Max Welling, "Auto-Encoding Variational Bayes," *arXiv*, 2022. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [10] Gregory Barber, Deepfakes Are Getting Better, But They're Still Easy to Spot, WIRED, 2019. [Online]. Available: <https://www.wired.com/story/deepfakes-getting-better-theyre-easy-spot/>
- [11] Diya Garg, and Rupali Gill, "Deepfake Generation and Detection - An Exploratory Study," *2023 10<sup>th</sup> IEEE Uttar Pradesh Section International Conference on Electrical, Electronics and Computer Engineering (UPCON)*, Gautam Buddha Nagar, India, pp. 888-893, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [12] Chandrasekaran, "The Rise of Deepfake Technology: A Threat to the Future of Truth?," *Journal of Cybersecurity*, vol. 6, no. 1, 2020.
- [13] Aaron van den Oord et al., "WaveNet: A Generative Model for Raw Audio," *arXiv*, 2016. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [14] Yuxuan Wang et al., "Tacotron: Towards End-to-End Speech Synthesis," *arXiv*, 2017. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [15] Yuval Nirkin, Yosi Keller, and Tal Hassner, "FSGANv2: Improved Subject Agnostic Face Swapping and Reenactment," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 1, pp. 560-575, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]

- [16] Caroline Chan et al., “Everybody Dance Now,” *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, Seoul, Korea (South), pp. 5932-5941, 2019. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [17] Alec Radford et al., “Improving Language Understanding by Generative Pre-Training,” *Preprint*, 2018. [[Google Scholar](#)]
- [18] Tero Karras, Samuli Laine, and Timo Aila, “A Style-Based Generator Architecture for Generative Adversarial Networks,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 12, pp. 4217-4228, 2021. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [19] Supasorn Suwajanakorn, Steven M. Seitz, and Ira Kemelmacher-Shlizerman, “Synthesizing Obama: Learning Lip Sync from Audio,” *ACM Transactions on Graphics*, vol. 36, no. 4, pp 1-13, 2017. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [20] H. Zhou, J. Zhang, and J. Zhang, “Audio-Visual Deepfakes,” *arXiv*, 2020.
- [21] Andreas Rössler et al., “FaceForensics++: Learning to Detect Manipulated Facial Images,” *arXiv*, 2019. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [22] Sanskriti Chandra et al., “A Novel Framework for Detection of Digital Face Video Manipulation Using Deep Learning,” *2023 3<sup>rd</sup> International Conference on Computing and Information Technology (ICCIT)*, Tabuk, Saudi Arabia, pp. 348-352, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [23] Chenqi Kong, Haoliang Li, and Shiqi Wang, “Enhancing General Face Forgery Detection via Vision Transformer with Low-Rank Adaptation,” *2023 IEEE 6th International Conference on Multimedia Information Processing and Retrieval (MIPR)*, Singapore, pp. 102-107, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [24] C. Rathgeb et al., “Crowd-Powered Face Manipulation Detection: Fusing Human Examiner Decisions,” *2022 IEEE International Conference on Image Processing (ICIP)*, Bordeaux, France, pp. 181-185, 2022. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [25] Zhiyuan Ma, Xue Mei, and Jie Shen, “3D Attention Network for Face Forgery Detection,” *2023 4<sup>th</sup> Information Communication Technologies Conference (ICTC)*, Nanjing, China, pp. 396-401, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [26] Farman Hassan, and Ali Javed, “Voice Spoofing Countermeasure for Synthetic Speech Detection,” *2021 IEEE International Conference on Artificial Intelligence (ICAI)*, Islamabad, Pakistan, pp. 209-212, 2021. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [27] Hafiz Mallik, “Securing Voice-Driven Interfaces against Fake (Cloned) Audio Attacks,” *2019 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR)*, San Jose, CA, USA, pp. 512-517, 2019. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [28] Aakriti Aggarwal et al., “IsSwap: Deep Fake Detection,” *2021 7<sup>th</sup> International Conference on Signal Processing and Communication (ICSC)*, Noida, India, pp. 194-199, 2021. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [29] Yuan Wang et al., “Dynamic Graph Learning with Content-guided Spatial-Frequency Relation Reasoning for Deepfake Detection,” *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Vancouver, BC, Canada, pp. 7278-7287, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [30] Rui Shao, Tianxing Wu, and Ziwei Liu, “Detecting and Grounding Multi-Modal Media Manipulation,” *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Vancouver, BC, Canada, pp. 6904-6913, 2023. [[Google Scholar](#)] [[Publisher Link](#)]
- [31] Qianyu Zhou et al., “Instance-Aware Domain Generalization for Face Anti-Spoofing,” *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Vancouver, BC, Canada, pp. 20453-20463, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [32] Aman Parikh et al., “Audio-Visual Deepfake Detection System Using Multimodal Deep Learning,” *2023 IEEE 3<sup>rd</sup> International Conference on Intelligent Technologies (CONIT)*, Hubli, India, pp. 1-6, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [33] Chao Feng, Ziyang Chen, and Andrew Owens, “Self-Supervised Video Forensics by Audio-Visual Anomaly Detection,” *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10491-10503, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [34] Yipin Zhou, and Ser-Nam Lim, “Joint Audio-Visual Deepfake Detection,” *IEEE/CVF International Conference on Computer Vision (ICCV)*, Montreal, QC, Canada, pp. 14780-14789, 2021. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [35] Anders Boesen Lindbo Larsen et al., “Autoencoding Beyond Pixels Using a Learned Similarity Metric,” *ICML'16: Proceedings of the 33<sup>rd</sup> International Conference on International Conference on Machine Learning*, vol. 48, pp. 1558-1566, 2016. [[Google Scholar](#)] [[Publisher Link](#)]
- [36] Kaggle, Deepfake Detection Challenge, Identify Videos with Facial or Voice Manipulations, 2020. [Online]. Available: <https://www.kaggle.com/c/deepfake-detection-challenge>
- [37] Kaggle, Celeb DF (v2). [Online]. Available: <https://www.kaggle.com/datasets/reubensuju/celeb-df-v2>
- [38] EndlessSora/DeeperForensics-1.0, [CVPR 2020] A Large-Scale Dataset for Real-World Face Forgery Detection, 2020. [Online]. Available: <https://github.com/EndlessSora/DeeperForensics-1.0>
- [39] Nick Dufour, Contributing Data to Deepfake Detection Research, Google AI Blog, 2019. [Online]. Available: <https://ai.googleblog.com/2019/09/contributing-data-to-deepfake-detection.html>
- [40] Bojia Zi et al., “WildDeepfake: A Challenging Real-World Dataset for Deepfake Detection,” *Proceedings of the 28<sup>th</sup> ACM International Conference on Multimedia*, pp. 2382-2390, 2020. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]