

Research Article

# Deepfake Detection and Classification of Images from Video: A Review of Features, Techniques, and Challenges

Dennis Lucky Tuanwi Bale, Laud Charles Ochei\*, Chidiebere Ugwu

Department of Computer Science, University of Port Harcourt, Port Harcourt, Nigeria

## Abstract

The proliferation of deepfake technology poses significant challenges to the integrity and authenticity of visual content in videos, raising concerns about misinformation and deceptive practices. In this paper, we present a comprehensive review of features, techniques, and challenges related to the detection and classification of deepfake images extracted from videos. Existing literature has explored various approaches, including feature-based methods, machine learning algorithms, and deep learning techniques, to mitigate the adverse effects of deepfake content. However, challenges persist, such as the evolution of deepfake generation methods and the scarcity of diverse datasets for training detection models. To address these issues, this paper reviews related work on approaches for deepfake image detection and classification and synthesises these approaches into four categories - feature extraction, machine learning, and deep learning. The findings underscore the importance of continued research efforts in this domain to combat the harmful effects of deepfake technology on society. This study provides recommendations for future research directions, emphasizing the significance of proactive measures in mitigating the spread of manipulated visual content.

## Keywords

Deepfake, Detection, Classification, Video, Image, Features, Techniques

## 1. Introduction

The proliferation of deepfake technology has ushered in an era fraught with concerns regarding the authenticity and reliability of visual content in videos. Leveraging advancements in artificial intelligence and deep learning algorithms, deepfake technology has enabled the creation of remarkably convincing manipulated videos and images, blurring the boundaries between reality and fiction [10]. Consequently, the imperative for robust techniques to detect and classify deepfake images within videos has become increasingly urgent. This paper addresses this concern by offering a comprehensive review of features, techniques, and challenges associated with deepfake image detection and classification in videos. Our

research delves deep into the challenges of combating the proliferation of manipulated visual content and explores the myriad methods employed in the field.

The primary research question driving this paper is: *How can the current state-of-the-art techniques for deepfake detection and classification of images within videos be evaluated and compared, and what are the key challenges and limitations faced by existing methodologies?* Addressing this question is pivotal in advancing the field of deepfake detection and classification. The aim of this paper is twofold: first, to systematically evaluate and compare the features, techniques, and challenges associated with deepfake detection and

\*Corresponding author: [laud.ochei@uniport.edu.ng](mailto:laud.ochei@uniport.edu.ng) (Laud Charles Ochei)

**Received:** 5 March 2024; **Accepted:** 18 March 2024; **Published:** 2 April 2024



Copyright: © The Author(s), 2023. Published by Science Publishing Group. This is an **Open Access** article, distributed under the terms of the Creative Commons Attribution 4.0 License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution and reproduction in any medium, provided the original work is properly cited.

classification of images within videos; and second, to propose a structured framework for comparing and evaluating deepfake detection and classification techniques, thereby facilitating a systematic assessment of performance metrics and criteria across different methodologies.

The main contributions of this research paper encompass several key aspects:

1. Conducts a comprehensive evaluation of state-of-the-art techniques for deepfake detection and classification, providing a panoramic overview of the diverse approaches employed in this domain.
2. Identifies the key challenges and limitations faced by existing methodologies, shedding light on areas that necessitate further investigation and improvement.
3. Synthesizes and compares different approaches for detecting and classifying deepfakes in video, categorizing them into four main categories: feature extraction, machine learning, and deep learning.
4. Presents a structured framework for comparing and evaluating deepfake detection and classification techniques, which facilitates a systematic assessment of performance metrics and criteria across different methodologies.
5. Provides recommendations for future research directions and potential areas of innovation in the field of deepfake detection and classification, aiming to propel advancements in this critical area of research [4].

This study is poised to make significant contributions by identifying the challenges and limitations of existing deepfake detection and classification methodologies, synthesizing and comparing different methodologies, and offering valuable insights into the effectiveness and applicability of various techniques for detecting and classifying deepfake images within video datasets. Moreover, it compares and evaluates deepfake detection and classification techniques, fostering a systematic assessment of performance metrics and criteria across different methodologies.

The subsequent sections of this paper are organized as follows: Section 2 provides an overview of related concepts, while Section 3 delves into a review of related work. Section 4 presents the findings from the review of related work, followed by a discussion of these findings in Section 5. Section 6 presents a structured approach for comparing and evaluating the different approaches for the detection and classification of deepfake images from videos. Finally, Section 7 encapsulates the conclusion drawn from this study, and Section 8 outlines potential avenues for future research.

## 2. Overview of Related Concepts

### 2.1. Deepfake Technology

Deepfake technology, propelled by advancements in deep learning algorithms, facilitates the creation of hyper-realistic but entirely synthetic visual content. This content can ma-

nipulate facial expressions, gestures, and even speech to deceive unsuspecting viewers, posing significant threats to various aspects of society [11]. Deepfake techniques leverage a combination of generative networks and encoder-decoder architectures to produce fake content in the form of images, videos, texts, or voices [10]. The popularity of deepfake technology has surged due to its accessibility and affordability, enabling creators to produce various forms of manipulated visual content.

#### A. Autoencoders

Autoencoders, a subtype of feedforward neural networks, play a crucial role in the creation of deepfakes. Designed by Geoffrey Hinton in the 1980s, autoencoders replicate input data from the input layer to the output layer, aiming to reconstruct the original input as accurately as possible [3, 5]. These neural networks consist of an encoder, a code, and a decoder, facilitating the transformation of input data into a compact representation before decoding it to generate the reconstructed output. Deep autoencoders, with multiple layers representing encoding and decoding, enable the compression and dimensional reduction of images, a pivotal step in deepfake creation.

#### B. Generative Adversarial Network (GAN)

Generative Adversarial Networks (GANs) constitute a powerful class of deep neural networks. Generative Adversarial Networks (GANs) consist of two neural network models: a generator and a discriminator. GANs operate in an adversarial setting to generate synthetic data resembling real data distributions [31]. These models are trained concurrently using adversarial training, in which the generator attempts to produce realistic data samples to deceive the discriminator, while the discriminator attempts to distinguish between real and fake data samples. As training progresses, both models improve iteratively, with the generator producing more realistic samples and the discriminator becoming more adept at distinguishing between real and fake samples. This competitive process enables GANs to capture and replicate the variations found in the training dataset, producing high-quality synthetic data that closely resembles the original dataset.

### 2.2. Deepfake Detection

Deepfake detection encompasses the process of identifying manipulated visual content within videos, aiming to discern discrepancies indicative of synthetic alterations. Techniques for deepfake detection often analyze facial expressions, landmarks, lip synchronization, and artifacts within video frames to differentiate between authentic and manipulated content [32, 28]. Various approaches, including neural network-based methods and large-scale dataset construction, have been proposed to enhance deepfake detection capabilities [7]. Leveraging deep learning algorithms, particularly convolutional neural networks (CNNs), facilitates the accurate detection and classification of deepfake images extracted

from videos.

### 3. Review of Related Work on Detection and Classification of Deepfake Images from Videos

This section provides an overview of recent advancements in the field of deepfake detection and classification within videos. The literature is categorized into key techniques and approaches, spanning feature-based methods, machine learning algorithms, and deep learning techniques.

#### 3.1. Feature-based Approaches

Feature-based methods focus on extracting distinctive features from video frames to identify anomalies indicative of deepfake content. [15] proposed a deepfake predictor (DFP) approach leveraging a hybrid of VGG16 and convolutional neural network architecture. Their method, trained on a deepfake dataset comprising real and fake faces, achieved impressive precision and accuracy rates for deepfake detection. [2] developed practical facial finding and animation processes, utilizing techniques such as displaced dynamic expression (DDE) and real-time high-fidelity facial capture systems. These methods enable the reconstruction of detailed facial features in real time, contributing to the detection of manipulated visual content. [27] developed a customized CNN algorithm for deepfake image identification, demonstrating its efficacy through comparative analysis with alternative methods. Their approach, utilizing convolutional neural networks (CNNs), showcases the importance of robust classification techniques for detecting deepfake images within videos.

Other notable feature-based approaches include those by [29], who detected deepfake videos using an attribution-based confidence metric, and [32], who focused on detecting both machine and human-created fake face images in the wild. These studies highlight the importance of feature extraction and analysis in identifying manipulated visual content within video datasets.

#### 3.2. Machine Learning Algorithms

Various machine learning algorithms, including support vector machines (SVMs), decision trees, and neural networks, have been employed to distinguish between real and deepfake content. [13] proposed an automated method for deepfake image classification, integrating deep learning and machine learning methodologies. Their framework, combining error level analysis, convolutional neural networks (CNNs), and support vector machines (SVMs), demonstrated robustness and efficiency in detecting deepfake images. [26] introduced a hybrid deep learning approach, termed NOLANet, for deepfake video detection. By lever-

aging spatial, spectral, and temporal content consistently, their method effectively differentiated between real and fake videos, highlighting the importance of multimodal analysis in deepfake detection. [25] developed a novel deepfake detection method utilizing CNN, CNN-LSTM, and CNN-GRU models trained on diverse datasets. Their study emphasizes the significance of transfer learning and sequence detection for accurate deepfake classification. [18] conducted a comparative analysis of CNN models for deepfake detection, training various architectures on datasets sourced from Kaggle. Their findings underscore the effectiveness of convolutional neural networks in identifying manipulated visual content. Additionally, the work by [9] focused on fighting deepfakes using deep learning-based detection methods, while [19] discussed current challenges and next steps in deepfake detection.

#### 3.3. Deep Learning Techniques

Deep learning models have demonstrated remarkable success in deepfake detection and classification. Deep learning techniques offer a powerful framework for detecting and classifying deepfake images extracted from videos. These techniques leverage the temporal information present in video sequences to improve the accuracy of deepfake detection and classification.

##### (a) Recurrent Neural Networks (RNNs)

Recurrent Neural Networks (RNNs) are a type of supervised deep learning algorithm that has shown promise in sequential data analysis, including deepfake detection. [24] combined RNN with CNN models to develop a hybrid approach capable of capturing long-term dependencies and improving deepfake detection accuracy.

##### (b) Long Short-Term Memory Networks (LSTMs)

Long Short-Term Memory Networks (LSTMs), a variant of RNNs, excel in learning and memorizing long-term dependencies, making them suitable for time-series prediction tasks. By mitigating the vanishing gradient problem encountered in traditional RNNs, LSTMs contribute to enhanced deepfake detection capabilities [24].

##### (c) Convolutional Neural Networks (CNNs)

Convolutional Neural Networks (CNNs) are capable of extracting features from images or video frames using kernels and classifying them automatically. CNNs have been widely used in deepfake detection and classification due to their ability to efficiently learn spatial hierarchies of features from images or video frames [24, 14]. CNN-based techniques, such as those proposed by [19, 12], have showcased significant advancements in identifying and mitigating the spread of deepfake content.

Wang, S., Yoon, S., & Lee, J. [19] introduced CNN-based approaches that considerably improved the identification and mitigation of deepfakes. Their research addressed existing issues and potential future paths in deepfake detection. They emphasised the need of using advances in CNN architectures

to increase the accuracy and resilience of deepfake detection systems. [19] exhibited significant advances in identifying faked visual content within movies, hence improving the overall reliability of deepfake detection systems.

Similarly, Pavel, A. B. et al. [12] did substantial research on CNN-based techniques to deepfake detection and classification. Their work attempted to examine and analyse existing methodologies, stressing the benefits and limits of CNN models in dealing with the issues given by deepfakes. [12] offered novel approaches to improving the performance of CNN-based detection systems, including fine-tuning network designs and optimising training techniques. Through empirical assessments and comparative studies, they revealed vital insights into the efficiency of CNN-based strategies in fighting the spread of deepfake content in videos.

Yang et al. [21] provides a comprehensive examination of deepfake detection and attribution techniques within the realm of forensic applications. The research acknowledges the interconnectedness of these tasks and emphasises the importance of attribution alongside detection in fighting deepfakes. The study uses deep learning techniques like CNNs and RNNs and traditional forensic methods to create robust frameworks for identifying and attributing deepfake content.

Kim et al [23] propose a novel deep learning architecture for deepfake detection and classification in videos. Their approach combines convolutional neural networks (CNNs) for spatial feature extraction with recurrent neural networks (RNNs) for temporal analysis. By integrating both spatial and temporal information, their model achieves improved accuracy in identifying manipulated content. Gupta and colleagues present a deep learning-based method for detecting deepfakes that takes advantage of temporal inconsistencies in videos. Their method entails training recurrent neural networks (RNNs) to analyse the sequential nature of video frames and detect irregularities indicating deepfake manipulation. The experimental results show that their model can accurately classify deepfake content [33].

Singh et al. propose a CNN-RNN hybrid model for deepfake detection and classification, with the goal of addressing the challenges presented by evolving deepfake generation techniques. Their method employs convolutional neural networks (CNNs) for extracting spatial features from video frames and recurrent neural networks (RNNs) for capturing temporal relationships. Experimental evaluations yield promising results in accurately detecting and classifying deepfake images within videos [34]. Zhao and co-authors present a deep learning framework for deepfake detection that incorporates both CNN and RNN architecture. Their model uses convolutional neural networks (CNNs) to extract spatial features and recurrent neural networks (RNNs) to analyse temporal patterns. Their approach detects manipulated content in videos with high accuracy because it integrates spatial and temporal information [35].

Choi et al. propose a deep learning-based method for deepfake detection that combines CNN and RNN modules

into a single framework. Their method uses convolutional neural networks (CNNs) to capture spatial features and recurrent neural networks (RNNs) to model temporal dynamics. Experimental results show that their model can accurately distinguish between real and manipulated content in videos [36].

## 4. Findings from the Review of Related Work

Our in-depth exploration of the existing literature on deepfake detection and classification has unveiled crucial insights into the challenges encountered by current methodologies and the promising strategies employed to address them.

### 4.1. Challenges Faced by Current Techniques

There are several challenges faced by current techniques for deepfake detection and classification of images from video as discussed below.

#### a) Evolution of Deepfake Generation Methods

One of the foremost challenges in deepfake detection lies in the continuous evolution of deepfake generation methods. With the rapid advancements in generative models and algorithms, existing detection systems often struggle to keep pace, resulting in diminished performance when confronted with deepfakes created using novel techniques [8].

#### b) Scarcity of Diverse Datasets

A significant hurdle in the development of robust deepfake detection systems is the limited availability of diverse datasets for training and evaluation purposes. Existing datasets often lack representation from various generative models, making it challenging to generalize detection models across different deepfake variations [12, 16].

#### c) Ethical Considerations

Ethical considerations surrounding the development and deployment of deepfake detection systems are paramount. Upholding principles of privacy, transparency, and responsible use is essential to mitigate potential societal harms and maintain public trust in the technology [30].

### 4.2. Promising Approaches and Techniques

#### Identifying Inconsistencies

Feature-based approaches have emerged as promising techniques for identifying inconsistencies in deepfake content. By analyzing facial expressions, eye movements, and other distinctive features, these approaches offer valuable insights into potential indicators of deepfake manipulation (Hou et al., 2022; [19].

#### High Accuracy Rates

Machine learning algorithms and convolutional neural networks (CNNs) have demonstrated remarkable accuracy rates in deepfake detection tasks. However, challenges such as



susceptibility to adversarial attacks persist, underscoring the need for further research and refinement in this domain (Chen et al., 2020 [9]).

## 5. Discussion of Findings

This section analyses and interprets the findings from the review of related work in the context of the research objectives. We also compare different approaches and techniques and discuss the implications of the reviewed literature for future research and practical applications.

### 5.1. Analysis and Interpretation of Findings

The literature review revealed significant progress in deepfake detection and classification techniques, with researchers employing a variety of approaches ranging from feature-based methods to machine learning algorithms and deep learning techniques. Key findings include:

**Feature-based Approaches:** Techniques focusing on extracting distinctive features from video frames, such as facial landmarks and micro expressions, have shown promise in identifying deepfake content [15, 19].

**Machine Learning Algorithms:** Various machine learning algorithms, including support vector machines (SVMs) and decision trees, have been utilized to train models for distinguishing between real and deepfake content [13, 17].

**Deep Learning Techniques:** Deep learning models, particularly convolutional neural networks (CNNs) and recurrent neural networks (RNNs), have demonstrated remarkable success in deepfake detection tasks [9, 25].

### 5.2. Comparison of Approaches and Techniques

To provide a comprehensive overview of the different approaches and techniques employed in deepfake detection and classification, a comparison is presented in Table 1.

**Table 1.** Approaches and advantages of deepfake detection and classification techniques.

Approach/Technique	Advantages	Limitations
Feature-based	Effective in capturing subtle inconsistencies. Can analyse specific facial features	Limited to surface-level analysis. May not capture complex manipulations
Machine Learning	Versatile and interpretable. Can handle diverse data types	Requires labelled training data. May struggle with generalization to new datasets
Deep Learning	Capable of learning complex patterns. Suitable for large-scale data analysis	Prone to overfitting. High computational requirements

### 5.3. Implications for Future Research and Practical Applications

The findings from the review have several implications for future research and practical applications:

**Resilient Detection Models:** Future research should focus on developing resilient detection models capable of adapting to evolving deepfake generation techniques and effectively generalizing across diverse datasets [21].

**Creation of Diverse Datasets:** Collaborative efforts are needed to curate diverse and representative datasets encompassing a wide array of generative models and deepfake variations. Such datasets are crucial for training and evaluating robust deepfake detection systems [12, 16].

**Exploration of Novel Approaches:** Novel methodologies, including multimodal fusion and explainable AI, hold promise for enhancing detection accuracy and interpretability, thereby advancing the state-of-the-art in deepfake detection and classification [8].

## 6. Structured Approach for Comparing and Evaluating Techniques

To systematically assess the performance of various techniques for detecting and classifying deepfakes from videos, it is essential to establish a structured approach. This framework should encompass key performance metrics and criteria that enable a comprehensive comparison of different methodologies. Below, a structured approach for evaluating deepfake detection and classification techniques is proposed, along with an illustrative example of different techniques presented in a tabular form.

### 6.1. Steps in the Structured Approach

**Step 1: Identify Performance Metrics and Criteria** - Define key performance metrics and criteria that are essential for evaluating the effectiveness of deepfake detection and classification techniques. These metrics may include accuracy,

precision, recall, false positive rate, false negative rate, computational efficiency, robustness, and generalization [6, 30].

**Step 2: Select Relevant Techniques** - Identify a comprehensive set of deepfake detection and classification techniques from the existing literature. Ensure that the selected techniques represent a diverse range of approaches, including traditional machine learning methods, deep learning models, and hybrid approaches [20, 1].

**Step 3: Gather Evaluation Data** - Collect a dataset comprising a diverse set of videos containing both genuine and deepfake content. Ensure that the dataset covers various aspects such as different video resolutions, lighting conditions, facial expressions, and manipulation techniques [20, 1].

**Step 4: Implement Techniques** - Implement each selected technique using appropriate frameworks and libraries. Fine-tune the parameters of the techniques based on the characteristics of the evaluation dataset [25, 14, 3].

**Step 5: Evaluate Performance** - Evaluate the performance of each technique based on the defined performance metrics and criteria. Conduct rigorous testing and validation to ensure the reliability and reproducibility of the results [6, 30].

**Step 6: Compare Results** - Compare the performance of different techniques using a structured framework. Present the results in a tabular form, highlighting the strengths and weaknesses of each technique across various metrics [24, 13, 3].

**Step 7: Draw Conclusions** - Analyze the findings of the evaluation process and draw conclusions regarding the effectiveness of different techniques for deepfake detection and classification. Identify trends, patterns, and areas for future research [6, 20, 21].

## 6.2. Performance Metrics and Criteria

The following performance metrics and criteria can be utilized for evaluating deepfake detection and classification techniques:

1. **Accuracy:** The overall correctness of the detection and classification results.
2. **Precision and Recall:** Precision measures the proportion of correctly identified deepfakes among all detected instances, while recall measures the proportion of correctly identified deepfakes among all actual deepfakes.
3. **False Positive Rate (FPR):** The rate of falsely identifying genuine videos as deepfakes.
4. **False Negative Rate (FNR):** The rate of failing to detect actual deepfakes.
5. **Computational Efficiency:** The computational resources required for detection and classification.
6. **Robustness:** The ability of the technique to perform effectively under various conditions, such as different video resolutions, lighting conditions, and facial expressions.
7. **Generalization:** The extent to which the technique can

detect and classify deepfakes created using different methods and tools.

## 6.3. Illustrative Example: Comparison of Techniques Using a Structured Approach

This section demonstrates how the novel structural method, as well as the performance metric and criteria discussed above, may be applied to a typical CNN-based technique and an RNN-based technique for detecting and classifying deepfakes from video. We also include a table with an example of comparing several deepfake detection and classification algorithms (CNN vs RNN) using the proposed structured approach. The structured approach for comparing CNN-based and RNN-based Techniques for Deepfake Detection and Classification translates to the following steps:

**Step 1: Problem Definition and Scope**

**CNN-based Technique:** Utilizes convolutional neural networks to extract spatial features from video frames for deepfake detection and classification. **RNN-based Technique:** Employs recurrent neural networks to capture temporal dependencies in sequential video data for deepfake detection and classification.

**Step 2: Data Collection and Preprocessing**

**CNN-based Technique:** Requires a large dataset of labelled videos containing both real and deepfake content for training. **RNN-based Technique:** Needs sequential video data with appropriate preprocessing to handle temporal dependencies and ensure data continuity.

**Step 3: Feature Extraction**

**CNN-based Technique:** Employs convolutional layers to automatically extract spatial features from individual frames of the video. **RNN-based Technique:** Utilizes recurrent layers to capture temporal patterns and dependencies across frames in the video sequence.

**Step 4: Model Architecture**

**CNN-based Technique:** Typically consists of convolutional layers followed by pooling layers for spatial feature extraction, and fully connected layers for classification. **RNN-based Technique:** Comprises recurrent layers such as LSTM or GRU to capture sequential information, followed by fully connected layers for classification.

**Step 5: Training and Optimization**

**CNN-based Technique:** Trained using backpropagation with optimization techniques such as gradient descent or Adam to minimize classification loss. **RNN-based Technique:** Trained using backpropagation through time (BPTT) or truncated backpropagation with optimization techniques to update weights and biases iteratively.

**Step 6: Performance Evaluation**

**CNN-based Technique:** Evaluated based on accuracy, precision, recall, false positive rate (FPR), false negative rate (FNR), computational efficiency, robustness, and generalization. **RNN-based Technique:** Assessed using similar perfor-

mance metrics including accuracy, precision, recall, FPR, FNR, computational efficiency, robustness, and generalization.

#### Step 7: Validation and Testing

**CNN-based Technique:** Validated and tested on separate datasets to assess generalization and robustness. **RNN-based Technique:** Similarly validated and tested on independent datasets to ensure reliability and generalization of the model.

**Table 2.** Comparison of CNN and RNN Techniques for Deepfake Detection and Classification.<sup>1</sup>

Criteria	CNN	RNN
Accuracy	High	Moderate
Precision	Moderate	High
Recall	High	Moderate
FPR	Low	Low
FNR	Low	Low
Computational Efficiency	High	Moderate
Robustness	Moderate	High
Generalization	High	Moderate

In this example, the CNN-based technique demonstrates higher accuracy and recall compared to the RNN-based technique. However, the RNN-based technique exhibits higher precision and computational efficiency. Both techniques show robustness and moderate generalization capabilities, with slight variations in performance across different metrics. The choice between CNN and RNN will then depend on the specific requirements and constraints of the application scenario, considering factors such as accuracy, computational efficiency, and robustness.

## 7. Future Work

As the field of deepfake detection and classification rapidly evolves, several avenues for future research and development emerge. This section delineates key directions for advancing the state of the art in deepfake detection and classification, highlight emerging challenges, and suggest practical recommendations.

### 7.1. Potential Areas for Future Research and Development

**Enhancing Adversarial Robustness:** Future studies should

concentrate on bolstering deepfake detection models against adversarial attacks. Techniques for detecting and mitigating sophisticated manipulation methods aimed at circumventing detection algorithms warrant exploration [6].

**Multimodal Fusion Approaches:** Investigating multimodal fusion methods, such as integrating visual and audio cues, offers promise for enhancing the accuracy and reliability of deepfake detection systems [30].

**Advancing Explainable AI:** The development of explainable AI techniques for deepfake detection is imperative to enhance transparency and interpretability. Research efforts should focus on elucidating the decisions of detection models, enabling users to comprehend classification outcomes [4].

**Real-time Detection Systems:** Efforts should be directed towards devising real-time deepfake detection systems capable of processing video streams instantaneously. Optimizing detection algorithms for efficiency and scalability is pivotal in combating the rapid proliferation of manipulated content online [20].

### 7.2. Recommendations for Improving Existing Techniques

**Diverse Dataset Construction:** Priority should be given to constructing diverse and representative datasets encompassing various deepfake generation techniques, actors, and scenarios. This would facilitate the development of more robust detection models capable of identifying a broader spectrum of deepfake content (Rössler et al., 2019).

**Interdisciplinary Collaboration:** Collaborative endeavours between researchers from diverse disciplines, including computer vision, machine learning, psychology, and sociology, are essential. Interdisciplinary approaches can yield more effective detection strategies that account for broader societal implications [1].

### 7.3. Emerging Challenges and Trends

**Evolution of Deepfake Generation Techniques:** The continuous evolution of deepfake generation techniques poses a significant challenge for detection algorithms. Researchers must adapt detection techniques to keep pace with advancements in generative models and manipulation methods (Nguyen et al., 2020).

**Deepfake Attribution and Forensics:** With the proliferation of deepfake content, there is a pressing need for tools and techniques for attribution and forensics. Research endeavors should focus on developing methods to trace the origin of deepfake content and identify responsible parties [21].

**Regulatory and Policy Considerations:** Policymakers face challenges in formulating regulations and policies to address the ethical, legal, and societal ramifications of deepfake technology. Future research should inform policy discussions and provide evidence-based recommendations for mitigating the adverse effects of manipulated content [22].

<sup>1</sup> Note that the values in the table are hypothetical and serve as an example to demonstrate the structured comparison approach.

## 8. Conclusion

In this paper, a comprehensive review of features, techniques, and challenges related to deepfake detection and classification in video content was conducted. Through the analysis, the aim was to contribute to the understanding of the evolving landscape of deepfake technology and its implications for society.

Our review highlighted various approaches and methodologies employed in the detection and classification of deepfake images from videos. We discussed feature-based methods, machine learning algorithms, and deep learning techniques, showcasing their strengths and limitations. Additionally, key findings were identified from recent research studies, including advancements in detection accuracy and the emergence of novel approaches such as multimodal fusion and explainable AI.

The proliferation of deepfake technology poses significant threats to the integrity of visual content in videos, with potential implications for misinformation, privacy violations, and societal unrest. Effective detection and classification of deepfake images are essential for preserving the authenticity and trustworthiness of digital media, safeguarding individuals' reputations, and maintaining public trust in visual information sources.

Despite the progress made in deepfake detection and classification, several challenges and limitations persist. These include the rapid evolution of deepfake generation techniques, the scarcity of diverse and representative datasets, and the vulnerability of detection models to adversarial attacks. Addressing these limitations requires concerted research efforts and interdisciplinary collaboration to develop robust and resilient detection systems.

The findings from the review has highlighted several practical implications for researchers, practitioners, and policymakers. Firstly, there is a need for continued investment in research and development to enhance the effectiveness and reliability of deepfake detection technologies. Secondly, stakeholders must prioritize the creation of diverse and inclusive datasets to ensure the generalizability and scalability of detection models across different contexts and scenarios. Finally, policymakers should consider implementing regulations and guidelines to govern the responsible use of deepfake technology and mitigate its potential negative consequences on society.

## Conflicts of Interest

The authors declare no conflicts of interest.

## References

- [1] Avcı, E., et al. (2021). Interdisciplinary Approaches to Deepfake Detection: A Perspective. *Journal of Interdisciplinary Studies in Artificial Intelligence*.
- [2] Chen, C., Hou, Q., & Zhou, K. (2014). Displaced Dynamic Expression Regression for Real-Time Facial Tracking and Animation. *ACM Transactions on Graphics (TOG)*.
- [3] Cheng, Z., et al. (2019). Deep Autoencoders for Image Compression and Dimensionality Reduction. *Neural Computing and Applications*.
- [4] Cheng, Z., et al. (2020). Explainable Deepfake Detection: A Survey. *ACM Computing Surveys (CSUR)*.
- [5] Chorowski, J., et al. (2019). Learning Representations by Maximizing Compression. *arXiv preprint arXiv: 1804.07723*.
- [6] Dang-Nguyen, D. T., & Boato, G. (2021). Advances in Deepfake Detection: A Review. In *Proceedings of the European Conference on Computer Vision Workshops (ECCVW)*.
- [7] Dolhansky, B., et al. (2020). DeepFake Detection Challenge (DFDC) Dataset. *arXiv preprint arXiv: 2006.07397*.
- [8] Hsu, W. T., Wu, J. S., & Lyu, S. (2021). Deepfake Detection: A Survey. *arXiv preprint arXiv: 2101.07761*.
- [9] Huh, M., Liu, A., & Owens, A. (2020). Fighting Deepfakes: Deep Learning-based Detection Methods. *arXiv preprint arXiv: 2006.07835*.
- [10] Mahmud, B. U., & Sharmin, A. (2021). Deep insights of deepfake technology: A review. *arXiv preprint arXiv: 2105.00192*.
- [11] Mirsky, Y., & Lee, M. (2020). Deepfake Detection: A Review. *IEEE Signal Processing Magazine*.
- [12] Pavel, A. B., Kiran, R. U., & Chen, H. (2020). Deepfake Detection Techniques: A Review. *IEEE Access*.
- [13] Rafique, M., et al. (2023). Deepfake Detection using Error Level Analysis and Convolutional Neural Networks. *IEEE Transactions on Information Forensics and Security*.
- [14] Rafique, R., Gantassi, R., Amin, R., Frnda, J., Mustapha, A., & Alshehri, A. H. (2023). Deep fake detection and classification using error-level analysis and deep learning. *Scientific Reports*, 13(1), 7422.
- [15] Raza, M., Khan, M., & Ahmed, S. (2022). A Hybrid Deep Learning Approach for Deepfake Detection. *IEEE Access*.
- [16] Rossler, A., et al. (2019). FaceForensics++: Learning to Detect Manipulated Facial Images. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
- [17] Shad, H. S., Rizvee, Md. M., Roza N. T., Ahsanul Hoq, S. M., Khan, M. M., Singh, A., Zaguia, A., & Bourouis, S. (2021). Comparative Analysis of Deepfake Image Detection Method Using Convolutional Neural Network. *Hindawi Computational Intelligence and Neuroscience Volume 2021, Article ID 3111676*.
- [18] Shad, R., Malik, S., & Khan, M. (2021). Comparative Analysis of CNN Models for Deepfake Detection. *IEEE Transactions on Information Forensics and Security*.
- [19] Wang, S., Yoon, S., & Lee, J. (2020). Deepfake Detection: Current Challenges and Next Steps. *IEEE Transactions on Multimedia*.



- [20] Yang, H., et al. (2020). Real-time Deepfake Detection: Challenges and Opportunities. In Proceedings of the IEEE International Conference on Multimedia and Expo (ICME).
  - [21] Yang, Y., et al. (2021). Deepfake Forensics: From Detection to Attribution. In Proceedings of the ACM Workshop on Deep Learning for Forensic Applications.
  - [22] Yan, W., et al. (2021). Policy Implications of Deepfake Technology: A Review. *Journal of Policy and Technology*.
  - [23] Kim, S., Park, J., & Lee, S. (2023). A Novel Deep Learning Architecture for Deepfake Detection and Classification in Videos. *Journal of Artificial Intelligence Research*, 45(2), 123-137.
  - [24] Zhou, J., et al. (2022). Towards Comprehensive Deepfake Detection: Challenges and Opportunities. *IEEE Transactions on Multimedia*.
  - [25] Hande, N., Patil, A., & Jain, A. (2022). Deepfake Detection using Hybrid CNN-RNN Models. *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*.
  - [26] Hande, R., Goon, S., Gondhali, A., & Singhaniya, N. (2022). A Novel Method of Deepfake Detection. *ITM Web of Conferences* 44, 03064 ICACC-2022  
<https://doi.org/10.1051/itmconf/20224403064>
  - [27] Lewis, J. K., Imad, E. T., Chen, H., Sandesera, V., Prasad, C., & Palaniappan, K. (2020). Deepfake Video Detection Based on Spatial, Spectral, and Temporal Inconsistencies Using Multi-modal Deep Learning. *IEEE Applied Imagery Pattern Recognition Workshop (AIPR)*.
  - [28] Li, Z., et al. (2018). Statistical Analysis of GAN-Generated Images for Fake Face Detection. *Journal of Visual Communication and Image Representation*.
  - [29] Fernandes, S., Raj, S., Ewetz, R., Pannu, J. S., Jha, S. K., Ortiz, E.,... & Salter, M. (2020). Detecting deepfake videos using attribution-based confidence metric. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (pp. 308-309).
  - [30] Kamiran, F., & Calders, T. (2021). Data preprocessing techniques for classification without discrimination. *Knowledge and Information Systems*, 64(2), 493-512.
  - [31] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative Adversarial Networks. *arXiv preprint arXiv: 1406.2661*.
  - [32] Tariq, U., et al. (2018). Neural Network-Based Methods for Detecting Fake GAN Videos. *IEEE Transactions on Information Forensics and Security*.
  - [33] Gupta, A., Sharma, V., & Patel, D. (2023). Exploiting Temporal Inconsistencies for Deepfake Detection: A Deep Learning Approach. *IEEE Transactions on Image Processing*, 32(4), 567-582.
  - [34] Singh, R., Kumar, A., & Jain, P. (2022). CNN-RNN Hybrid Model for Deepfake Detection and Classification. *International Journal of Computer Vision*, 78(3), 215-230.
  - [35] Zhao, H., Li, X., & Wang, Y. (2022). Integrating CNN and RNN Architectures for Deepfake Detection in Videos. *Pattern Recognition Letters*, 41(5), 701-716.
- Choi, H., Lee, D., & Kim, Y. (2021). Unified CNN-RNN Framework for Deepfake Detection in Videos. *IEEE Transactions on Multimedia*, 29(1), 45-59.