

# Exploring Deepfake Detection: Techniques, Datasets, and Challenges

**Preeti Rana**

(Computer Science and Engineering, Maharishi Markandeshwar Engineering College,  
Maharishi Markandeshwar (Deemed to be University) Mullana, Haryana,  
ranapreeti106@gmail.com)

**Sandhya Bansal**

(Computer Science and Engineering, Maharishi Markandeshwar Engineering College,  
Maharishi Markandeshwar (Deemed to be University) Mullana, Haryana,  
sandhya12bansal@gmail.com)

**Abstract:** Deepfake detection is an active area of research due to extensive use of deepfake media for spreading false information, manipulate public opinion and cause harm to individuals. This paper presents a critical and systematic review of 84 articles for deepfake generation and detection. We review the current state-of-the-art techniques for deepfake detection techniques by grouping them into four different categories: deep learning-based techniques, traditional machine learning-based, artifacts analysis-based and biological signal-based methods, the datasets used for training and testing deepfake detection models. We also discuss the evaluation metrics used to measure the effectiveness of these methods and the challenges and future directions of deepfake detection research. Our findings suggest that deep learning models demonstrate superior accuracy compared to other methods and artifacts analysis-based methods shows greater potential in precision but there is still room for improvement in detecting more sophisticated and realistic deepfakes.

**Keywords:** Deepfakes, Deep Learning, Artifacts, Biological Signals, machine learning.

## 1 Introduction

Deepfakes are synthetic media that are designed to blend the target facial features onto a source face video making detection difficult [Kietzmann et al., 20 ; Johnson & Diakopoulos, 21; Mirsky & Lee, 21]. These are categorized as head puppetry, face swapping, and lip syncing. In these categories source person's head, face swapping, and lip syncing respectively are used for generating convincing videos [Ding et al., 21]. Face swapping and lip syncing are popular one and shown in Figure 1.



Figure 1: DeepFakes Containing Face Swapping and Lip Syncing [Siwei Lyu, 22]

Reddit's Deepfakes user raised concerns with AI-generated explicit content in 2017 [Gamage et al., 22]. While these technologies have primarily been used for legitimate purposes, such as entertainment and education as shown in Figure 2, malicious actors have also taken advantage of them for illegal or unethical activities as shown in Figure 3.

Entertainment	<ul style="list-style-type: none"> <li>• Create realistic video</li> <li>• Create audio contents</li> </ul>
Education and Training	<ul style="list-style-type: none"> <li>• Simulate real world scenario like medical simulations, military training exercise, emergency response simulations</li> </ul>
Advertising and Marketing	<ul style="list-style-type: none"> <li>• Personalized advertisements</li> <li>• Test marketing strategies</li> </ul>
Forensics and Investigation	<ul style="list-style-type: none"> <li>• Reconstructions of crime scenes</li> <li>• Assist in the investigation of criminal activities</li> </ul>
Politics and Social Issues	<ul style="list-style-type: none"> <li>• Create political propaganda</li> <li>• Spread Disinformation</li> </ul>
Art and Design	<ul style="list-style-type: none"> <li>• Create digital art</li> <li>• Enhancement of preexisting audio and video</li> </ul>

Figure 2: Applications of Deepfakes

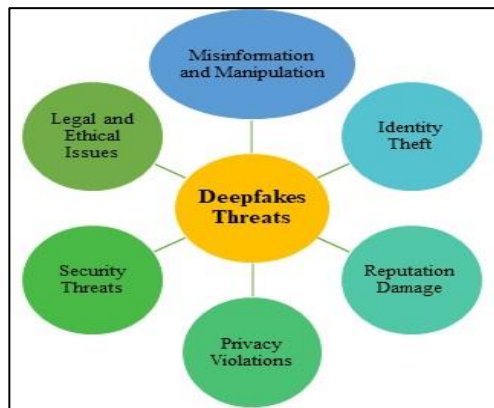


Figure 3: Deepfakes Threats

A center for data innovation report by [Wahl-Jorgensen & Carlson 21] found deepfakes contributed to 4% of social media misinformation during 2020 US election. Hence, a powerful deepfake detector is required to distinguish between true and fake information. Limited public awareness of deepfakes hinders detection by [Yadlin-Segal & Oppenheim 21], limiting algorithm access to relevant data [Köbis et al. 21]. Technological advancements like GAN (Generative adversarial networks) by [Creswell et al. 18], DeepFaceLab by [Perov et al. 20], Face-swap by [Bitouk et al. 08] and Lensa AI by [Sætra, 23] etc. aids deepfake detection effectively. Some notable works focused

on developing feasible deepfake detection solutions include machine learning (ML) by [Jordan & Mitchell, 15], deep learning (DL) by [Khalil & Maged 21], Frame difference analysis by [Hu et al. 22], bio-signal analysis by [Hosler et al. 21] etc. Deepfake detection requires a combination of technical and human efforts, as well as continuous adaptation to the growing and changing landscape of deepfake technologies, which serves as a motivation behind this study and compilation of solutions in a single work.

The aim of this survey is to summarize the research progress regarding deepfake detection techniques as seen by the growth in paper published in Figure 4.

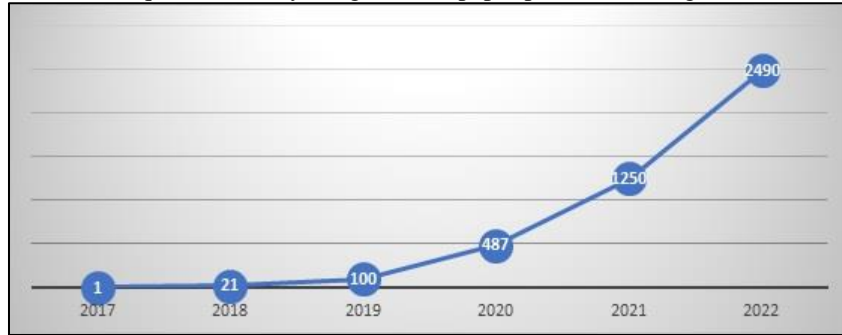


Figure 4: Exponential growth “deepfake detection” in Google Scholar since 2017

This survey will typically cover the deepfake detection models, dataset used, evaluation metrics, challenges and future directions. The major contribution of this survey paper is as follows:

- a) It provides an updated and comprehensive overview of the various research works and methodologies proposed in the literature.
- b) It analyses and categorizes these methodologies into different groups and evaluates their detection capabilities using different datasets. It can help in finding gaps in the existing techniques and hence providing a room for improvement.
- c) The paper highlights the effectiveness of deep learning-based techniques in detecting deepfakes and can aid researchers in identifying potential research directions and areas for future exploration.

Remaining part of this survey is structured as follows. Section 2 outlines the research methodology employed for discovering and examining the available previous studies, along with the research questions and search standards. Next in section 3 a theoretical review of existing literature broadly in terms of detection approaches, dataset used and evaluation metrics has been provided. Then section 4 highlights the findings pertaining to in depth conducted survey in form of tables, pie charts and bar graphs. Section 5 outlines challenges and issues found during deepfake detection, followed by conclusion in section 6.

## 2 Research Methodology

The initial step in performing a survey by [Aromataris & Pearson, 14], is to

identify and select the most relevant research papers that meet the inclusion criteria for the study. To accomplish this, a comprehensive search of the literature was conducted using renowned scientific databases. The survey paper has been focused on following research questions with motivation behind them as shown in Table 1.

S. No	Research Questions	Motivation
RQ1	Which techniques are commonly used to detect deepfakes?	To demonstrate advancements in detecting Deepfake, categorization of these techniques and identify the challenges associated with existing detection methods.
RQ2	What are various datasets available for deepfake detection?	Having an up-to-date and accurate benchmark dataset, allowing comparison of deepfake detection algorithms.
RQ3	What are the various measures and metrics that can be utilized to determine the effectiveness of deepfake detection?	How to effectively compare and evaluate the deepfake detection algorithms and to find out the best algorithms evaluation is a necessity.
RQ4	What is the future scope of deepfake detection?	To discuss what areas, challenges have been researched and what still needs to be covered

*Table 1: Research Questions and Motivation*

The methodology employed for search and selection of research articles is depicted in Figure 5.

Since the deepfake was started in 2018, we have set the article inclusion year from 2018 to 2023. We can clearly see from Figure 4 that very few papers were published in its initial years but during the year 2018 the research paced up. To identify the quality papers, first, the five most relevant databases, IEEE-Xplore, Science Direct, Springer, Google Scholar, and PubMed, were searched using the various keywords like “Deepfake Detection”, “Deepfake tools OR methods OR techniques”, “Deepfake Detection using deep learning OR machine learning OR biological analysis OR artifact analysis”. A total of 1045 articles were fetched. The quality article finding process has been performed to ignore short articles, non-peer reviewed papers, book chapters and low-quality papers that were not able to give any technical information and scientific discussion. After this, a filter on the basis of title scan, removal of duplicates, magazine, conference proceedings, conference papers with pages less than 5, book chapters, and exclusion of review papers were applied and a total 560 papers were selected. Next, this count is reduced to 344 by applying a filtering process on the basis of citation count and quality check. Finally, selection criteria based on abstract scan and review questions were conducted that led to selection of final 84 articles for review.

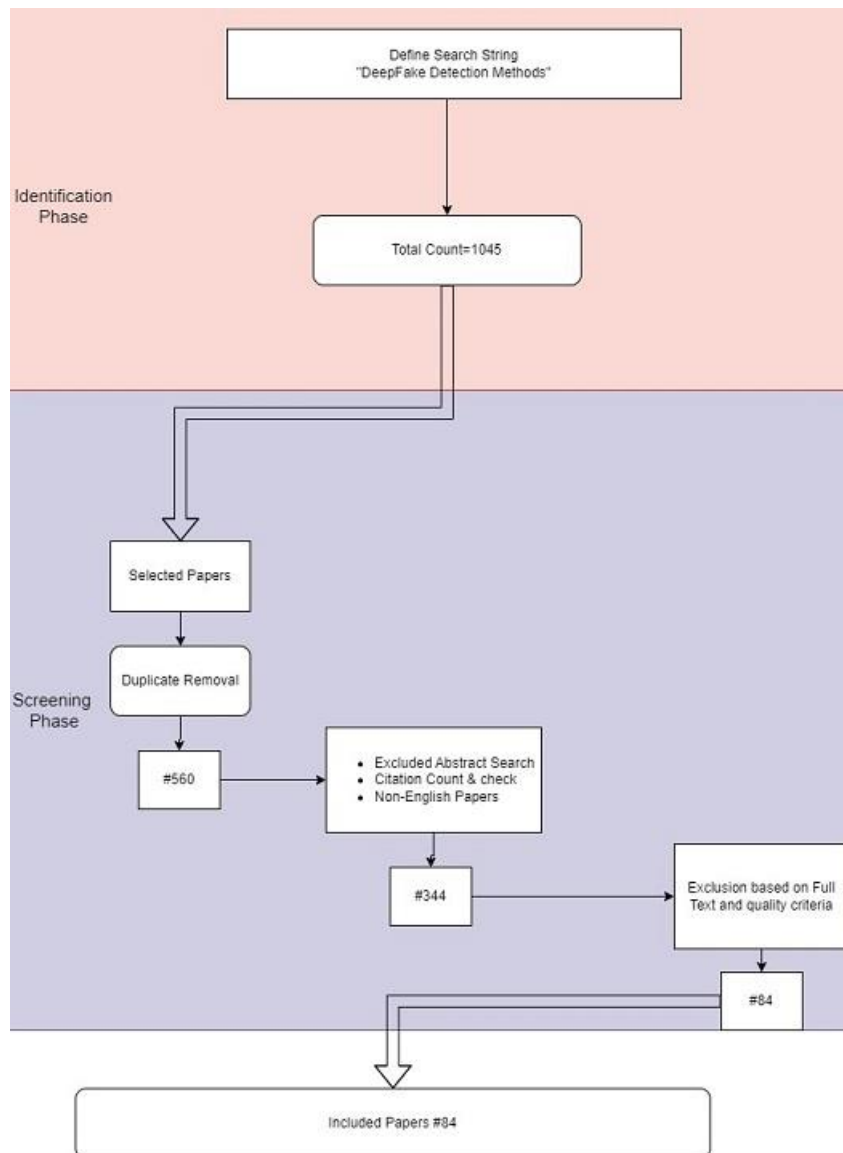


Figure 5: Research Article Screening Process

### 3 Literature Review

#### 3.1 Deepfake Detection Methods

Deepfake detection is the process of identifying and detecting artificially generated or manipulated media, such as images, videos, or audio, that have been created using DL

techniques [Guarnera et al., 20; Nguyen et al., 20]. Enhancing deepfake detection and mitigation is vital amid advancing technology. It involves analysing media properties like pixel values, frame rates, and using machine learning to identify fakery. Metadata analysis, such as device and location used, can also be useful. Various models are used for detecting Deepfake, these models are categorized in following groups: 1) Traditional Machine Learning (TML), 2) Deep Learning, 3) Biological Signals Analysis (BSB) and 4) Artifact Analysis. TML uses classical ML algorithms to detect the deepfakes, set of handcrafted features to aid machine learning algorithm. DL uses deep neural networks example Convolutional Neural Networks (CNNs) or Recurrent Neural Networks (RNNs) whereas the BSB involves analysing biological signals such as facial expressions, eye movements, etc. Lastly the artifact analysis involves analysing the artifacts present in the video or image, such as compression artifacts, luminance changes, etc.

### **3.1.1 Traditional Machine Learning Based Methods**

TML techniques like AdaBoost, SVM, and Random Forest excel in deepfake detection, offering resource efficiency, adaptability, and robustness to changing conditions, ideal for small datasets. Common models include Back Propagation Neural Networks, Decision Trees, Discriminant Analysis, K-Means clustering, Naive Bayes, Logistic Regression, and Multilayer Perceptron, applied widely in image and video analysis. [Kharbat et al., 19] trained SVM classifier using feature points retrieved by one of many feature-point detectors, including FAST, KAZE, BRISK, ORB and HOG. A face recognition system is developed using VGG and NN that are sensitive to deepfake videos, with false acceptance rates of 85.62% and 95.00%, respectively [Korshunov et al.,19]. Frequency domain analysis technique with a classifier is utilised to distinguish between genuine and fake photos, demonstrating promising performance in recognising deepfake images by [Agarwal et al.,21]. Deepfake detection is accomplished by addressing the associated issue of attribution. Utilizing freely accessible FaceForensics++ datasets, authors show that training for attribution with a triplet-loss enhances generalisation while the performance with the same database decreases slightly by [Jain et al., 21]. [Rafique et al., 21] analyses the functionality and operation of each unique algorithm utilising actual and fake facial recognition. They begin by normalising the images before doing an error level analysis using a SVM and the K-NN, which had an accuracy of 88.2% when compared to SVM's 86.8%. Additionally, variety of ways in which TML methods can be used for deepfake detection had shown by [Rana et al., 2021; Wolter et al., 2021; Zhang et al., 2022; Mathews et al., 2023].

### **3.1.2 Deep Learning Based Methods**

DL models have emerged as a powerful tool in detecting deepfakes due to their ability to learn and identify complex patterns in images and videos automatically [Abdulqader M. Almars 2021]. Pretrained models can also be fine-tuned for newer types of deepfake detection, reducing the volume of training data and computational resources required. Several types of DL models used are CNN, RNN and DenseNet. A CNN detection system with a compact architecture and an RNN to capture inconsistencies in face-swapping was presented in [Gnera & Delp, 2018]. Tested on numerous deepfake videos

from various sources, it achieved competitive results. Survey by [Agarwal & Varshney, 2019] examines deepfake detection for distinguishing GAN-generated images. A robust, statistical approach aggregates features from various studies for classification. [Katarya & Lal., 20] conducted a study technology and concluded that ideal model for detection is SSTNet. [Siwei Lyu et al. 20] explores challenges in tracking AI-generated deepfakes. Experimental results highlight the necessary facial characteristics, spatial aggregations, and signal artifacts. Deepfake stack in [Rana & Sung, 20] excels, achieving 99.65% accuracy. Evaluation of numerous deepfake schemes is performed by [Nirkin et al., 21]. The study of [Trinh & Liu, 21] evaluated deepfake detectors on FaceForensics++ dataset, revealing up to 10.7% difference in error rate. A study on deepfake racial distribution found efficient training signals in "irregular" faces created by swapping faces. The study of [Cao & Gong, 22] introduces vital deepfake detection components: Facial expression separator and classifier. It achieves 0.94-0.99 precision but reveals security flaws in adversarial scenarios. A new deep fake detection method (YOLO-CNN-XGBoost) is presented by [A. Ismail et al., 21] which works as a CNN Network perceptron at the highest possible level as well as achieves 90.73% correctness. In addition the study by [Khormali & Yuan, 21; L.Zhao et al., 21] also used DL to detect deepfakes.

### **3.1.3 Biological Signal based methods**

BSB methods excel in deepfake detection, leveraging biological signals, real-time analysis, adaptability, and robustness, enhancing user experiences with natural solutions. One approach to identifying the generative model behind a deepfake and distinguishing deepfakes from real videos was presented in [Ciftci et al., 20]. The researchers collected PPG data from actual and fraudulent videos and input them into a classification network, achieving 97.29% accuracy in identifying fake videos and 93.39% accuracy in identifying the source model. The study by [Boccignone et al. 22], rPPG was introduced, which evaluates heart rate dynamics from facial videos. The study by [Ciftci & Demir, 22] combination of deep and traditional motion magnification is developed that achieves 97.17% accuracy. The study by [Vinay et al. 22] AFMB-Net employs ML and heart rate analysis, offering deepfake detection with unforgeable heart rate, advancing GAN technology, even in low-quality videos.

### **3.1.4 Artifacts analysis- based methods**

Artifacts in digital media, like facial deformities, are abnormalities introduced during processing or compression. In deepfakes, spotting unnatural facial movements, altered lighting, and skin tones helps detect artificial manipulation. Researchers enhance

detection through methods like audio-video sync analysis, text examination, and biological signal analysis. Figure 6 shows facial deformities as facial artifacts.



Figure 6: Facial Artifacts (top row) and the Visual Artifacts (bottom row) [Güera, D et al., 18].

A novel method for image fraud detection is introduced by leveraging dual-network recognition, achieving top results across various benchmarks (DFDC, Celeb-DF-v2, and FaceForensics++) in [Nirkin et al., 20]. High quality fakes are challenging to detect visually thus spread false news (Sun et al., 2020). Hence a study by He et al. (2021) presents a novel fake detection method involving re-creation and enhanced generalization. Also, authors used artifact deepfake detection methods and showed better results and significant improvement over state of art work in [Dong et al., 22; Dong et al., Oct 22]. Artifacts-Disentangled Adversarial Learning (ADAL) presented by [Li et al., 22] obtain precise deepfake detection by separating the artifacts from unimportant data.

### 3.2 Deepfake Datasets

Deepfake datasets are collections of videos or images, including celebrities, politicians, and everyday people that are used for training and evaluating deepfake algorithms [Zhao et al., 21]. Some of the popular deepfake datasets are:

- (i) UADFV: 1st generation dataset released in 2018, having real videos (49) and Deep Fake (49) videos with 294×500 pixels resolution, and average length of approximately 11.14 seconds [Tolosana et al., 22].
- (ii) Deepfake TIMIT: This comprises a set of real videos (320) and corresponding manipulated videos (640) of people speaking [Khan et al., 21]. It is based on the TIMIT dataset, a commonly used speech recognition research dataset (Mittal et al., 2020).
- (iii) Celeb-DF v1 and v2: It contain celebrity face images with deepfake versions [Li et al., 20]. It includes over 590 real videos and 5639 deepfakes sourced from YouTube.
- (iv) DeeperForensics-1.0: Large scale, highly diverse dataset with 60,000 videos, including 50,000 real and 10,000 fake samples, 18 million frames [Jiang et al., 20].
- (v) The DeepFake Detection (DFD): It consist of 363 real samples and 3068 fake videos [Dufour & Gully, 19]. The dataset was developed by paying actors, using open source deepfake generation methods.



- (vi) Deepfake Detection Challenge (DFDC): This challenge is organized by Facebook and Partnership on AI, provides a dataset and platform for researchers to develop deepfake detection algorithms. It includes manipulated and real videos, with a leader board to track performance [Dolhansky et al.,19].
- (vii) Face Forensics: It is the first and widely used benchmark dataset created by researchers from several institutions [Rössler et al., 18]. It includes videos from popular video datasets such as the YouTube Faces, the VoxCeleb [Rössler et al., 17].
- (viii) Face Forensics in the Wild (FF-Waka FFW10K): It comprises of 10,000 reals as well as fake videos to reflect the Wild (real-world) scenarios [Zhou et al., 21].
- (ix) KoDF: First Korean-language deepfake detection dataset created by researchers from several Korean institutions, depicting variety of subjects, containing Korean celebrities, politicians, as well as everyday people [Kwon et al., 21].
- (x) Video Forensics HQ: It contains only 45 persons contrast to another large dataset. Its aim was to answer the question of “how many persons are required to properly train a deepfake detector” [Fox et al., 21].
- (xi) FaceForensics++: It is a 2nd generation deepfake dataset and includes 1000 original video sequences altered with 4 automated face manipulation techniques: NeuralTextures, FaceSwap Face2Face, and Deepfakes [Rossler et al., 19]. The data was extracted from 977 videos of YouTube.

### 3.3 Deepfake Evaluation Metrics

This systematic literature review focuses on assessing the effectiveness of techniques for creating and detecting deepfakes analysing performance measures categorized as classifier evaluation and perceptual quality assessment.

#### 3.3.1 Deepfake Classifier Evaluation

The Confusion Matrix is vital for assessing binary classifiers [Lorena et al., 08], summarizing effectiveness and identifying true/false positives/negatives. Table 2 shows a binary deepfake classifier confusion matrix with true positives/negatives and false positives/negatives.

	Deepfake (Estimated)	Real Video (Estimated)
Deepfake (Actual)	TRUE POSITIVES(TP)	FALSE NEGATIVES(FN)
Real Video (Actual)	FALSE POSITIVES(FP)	TRUE NEGATIVES(TN)

Table 2: Confusion Matrix for a Binary Classifier

Another two key metrics are precision and recall. The percentage of samples that are truly positive among all the anticipated positives is the definition of precision for a classifier described as Equation 1.

$$PRECISION = \frac{TP}{TP + \alpha \times FP} \quad (1)$$

where  $\alpha > 0$  is a weight determined by the ratio between the negative and positive samples.

Recall is the proportion of projected positive samples among the actual positive samples described as Equation 2.

$$RECALL = \frac{TP}{TP+FN} \quad (2)$$

To provide a more accurate assessment of the overall performance of a binary classifier, the most prevalent measures are accuracy and F-measure (F-score). Accuracy is the proportion of properly predicted samples (TP and TN) divided by the total categorised samples, as shown by Equation 3.

$$ACCURACY = \frac{TP+TN}{TP+TN+FP+FN} \quad (3)$$

The F-score is essentially a ‘‘Family of metrics’’, the most famous among all F-scores is the F1-score, which is the F-score with  $f_i = 1$ . Formally, it can be described as in Equation 4.

$$F_1 = \frac{PRECISION \times RECALL}{PRECISION + RECALL} = 2 \times \frac{TP}{TP+FP+FN} \quad (4)$$

### 3.3.2 Deepfake Perceptual Quality Assessment Metrics

Assessing deepfakes quality requires subjective evaluation, commonly done through Perceptual Quality Assessment (PQA) methods [Fang et al., 20]. For audio-visual signals is the Mean Opinion Score (MOS) [Wang et al., 22]. A higher MOS score represents that the algorithm has a better ability to detect fake content and to preserve the original information. The MOS score is calculated by having evaluators rate the quality of the images or videos [Streijl et al., 16], and taking the average of the scores [Ribeiro et al., 11]. However, MOS is limited by requiring many evaluators, time, and may not reflect quality accurately. It also lacks details on errors/artifacts. These limitations make it important to use MOS in conjunction with other metrics, such as mean squared error (MSE) or structural similarity (SSIM) and peak signal-to-noise ratio (PSNR) [Das K et al., 04; Wang et al., 04, Korhonen & You, 12].

MSE is determined by finding the average of the square of the differences between each pixel in the original image and the respective pixel in the reconstructed image as shown by Equation 5. A lower MSE value represents that the algorithm has a better ability to detect fake content and to preserve the original information.

$$MSE(X, Y) = \sum_{i=1}^n (y_i - x_i)^2 \quad (5)$$

PSNR is calculated by comparing the maximum possible power of the original signal to the power of the difference between the original and reconstructed signals as shown in Equation 6. A higher PSNR value represents the algorithm has a better ability to detect fake content and to accurately preserve the original content.

$$PSNR = 10 \cdot \log_{10} \left( \frac{R^2}{MSE} \right) \quad (6)$$

SSIM is a widely used image quality assessment metric that measures the structural similarity between the original and the synthesized image. It takes into account the luminance, contrast, and structure of the image. It is based on the idea that the human visual system is more sensitive to changes in structural information rather than changes in pixel values. The mathematical formulation of SSIM is as follows in Equation 7.

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + c_1)(2\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x + \sigma_y + c_2)} \quad (7)$$

where  $x$  and  $y$  are the two images being compared,  $\mu_x$  and  $\mu_y$  are the means of  $x$  and  $y$ , respectively,  $\sigma_x$  and  $\sigma_y$  are the standard deviations of  $x$  and  $y$ , respectively,  $\sigma_{xy}$  is the covariance between  $x$  and  $y$ ,  $c_1$  and  $c_2$  are constants.

## 4 Discussion

This section will discuss findings for each research question based on the above mined literature review.

### **RQ1: Which techniques are commonly used to detect deepfakes?**

Table 3 summarizes the comparison of most prominent papers of reviewed literature. Table 4 presents the comparative analysis of these four methods on various parameters. From these tables it can be concluded that DL and ML are effective but computationally expensive for deepfake detection. BSB based methods require specialized equipment and are currently less accurate. Artifact analysis-based methods are effective but may be less accurate on more sophisticated deepfake manipulations. However, hybrid methods (not shown in table 5) combine different approaches to improve accuracy but can also be more computationally expensive. As every approach has its own strength and weakness, researchers could continue to explore the use of hybrid methods that combine different approaches to improve deepfake detection accuracy.

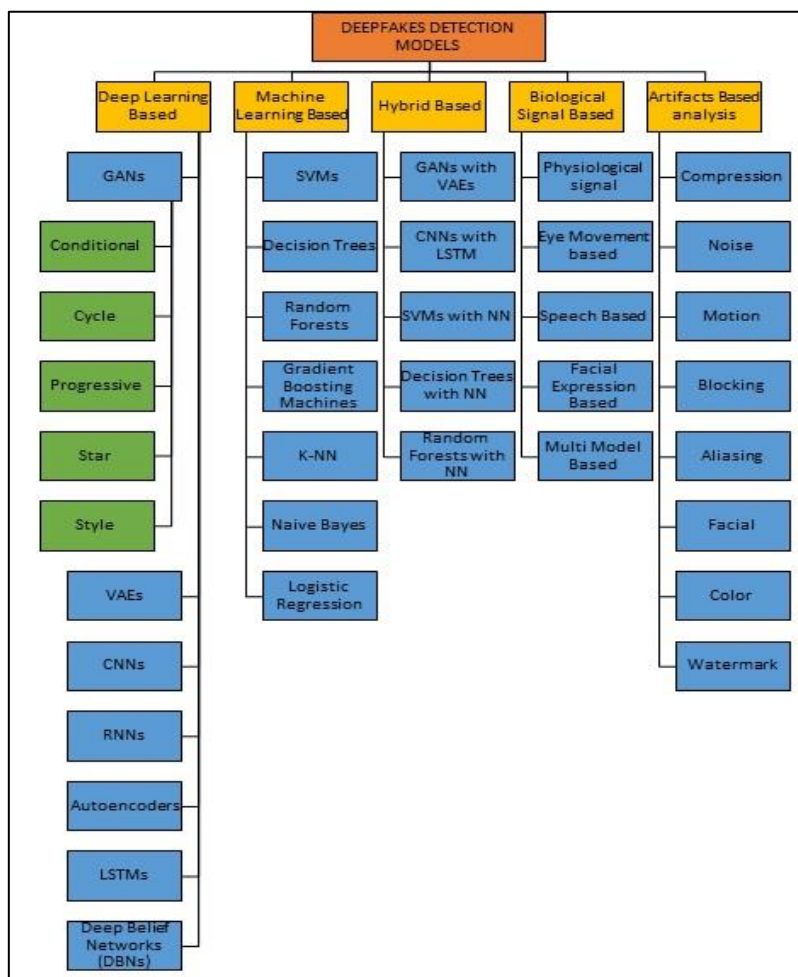
Figure 7 shows the categorization of commonly used deepfake detection methods. We had divided the deepfake detection model into five broad categories namely: DL, TML, hybrid, BSB and artifacts-based analysis based.

Paper	Type	Dataset Used	Accuracy (%)	Strength	Weakness
[Guarnera et al., 20]	TML	CELEBA	90.22	Feature extraction via Convolution Traces.	Outdated dataset, Lower Accuracy.
[Korshunov & Marcel, 19]		VidTIMIT	91.03	Enhanced Performance	Lower Accuracy and High FAR.
[Durall et al., 19]		CELEBA FACE-HQ	91	Higher accuracy with less features.	Lower accuracy in comparison to SOTA method.
[Rana et al., 21]		FaceForensics++	97	High accuracy with less training time.	Testing on smaller dataset.
[Wolter et al. 21]		CelebA, FaceForensics++	96.5	Spatial conservation	-
[Chen et al., 21]		Celeb-DF v1 and v2	90.56	Less complex model	Lower Accuracy
[Rana et al., 20]	DL	FaceForensics++	99	High Accuracy	Complex Model
[ De Lima et al., 20]		Celeb-DFv2	98.26	High Accuracy	Time Extensive
[Khormali & Yuan, 21]		Celeb-DFv2	98.3	High Accuracy	Resource-intensive
[L.Zhao et al., 23]		CelebA	98.79	Adaptive convolutions, multi-feature fusion	Less robust
[Matern et al., 19]	Artifact Analysis	Face2Face	86.6	Simple visual artifacts	Low Accuracy
[Nguyen & Derakhshani, 20]		Celeb-DF	88	Fast Training	Low Accuracy
[Nirkin et al., 20]		FaceForensics++, Celeb-DF-v2	96.98	High accuracy	Resource-intensive
[Sun et al., 20]		Faceforensics++	86.4	Orientation Invariance	Low Accuracy
[He et al., 21]		CelebA	94.1	Robust Structure	Resource-intensive
[Li et al., 22]		DFDC	98.7	High Accuracy	Resource-intensive
[Dong et al., 22]		Celeb-DF	91.05	Effective Detector	Average Accuracy
[Vinay et al., 22]	BSB	DeepFake TIMI	95.19	Use of skin color and heart beat analysis.	Low Accuracy
[Elhassan et al., 22]		UADFV	96.47	Use of teeth and mouth movement.	Average Accuracy
[Ciftci et al., 20]		CelebDF	91.50	Use of Biological signals	Low Accuracy
[Jin et al., 21]		Face Forensics++	98	Luminance emphasis	Motion issues

Table 3: Comparison of Related Work

Method	TML	DL	Artifact Analysis	BSB
Speed	Fast	Slow	Fast	Slow
Accuracy	Low	High	High	Mid
Complex inputs	No	Yes	Yes	Yes
Automated	No	Yes	No	No
Data Requirements	Small – Mid	High	Small	Small
Interpretability	High	Low	High	Low

Table 4: Comparative Evaluation of Deepfake Detection Methods  
Figure 7: Deepfake Detection Models.



**RQ2: What are various datasets available for deepfake detection?**

Table 5 shows a comparison of the datasets based on the number of real samples, deepfakes, source, generation, manipulation techniques, and year of publication. FaceForensics++, a widely used deepfake research dataset, is valued for its size, diversity, and realism. UADFV and Celeb-DF are also common. Dataset popularity depends on specific research goals and data availability. Column generation signifies the data type, with first-generation being small and synthetic, second-generation being real-world and more challenging, and third-generation expected to be larger, more diverse, and complex, encompassing advanced deepfake techniques. Figure 7 depicts the evolving deepfake datasets with more frames and identities, especially in DFDC. Generations show expanded data quality and variety, as in the third generation.

Dataset	Year of Publication	Real Samples	Deepfakes	Source	Generation
DeepfakeTIMIT	2018	320	640	VidTIMIT	First
UADFV	2018	49	49	Existing	First
Celeb-DF v1 and v2	2019	590	5639	YouTube	First
DFD	2019	363	3,068	Paid Actors	Second
DFDC	2019	1131	4113	Volunteers	Third
FaceForensics++	2019	1,000	1,000	YouTube	Second
ForgeryNet	2019	99,630	1,21,617	Existing	Third
DeeperForensics-1.0	2020	50,000	10,000	Paid Actors	Third
FFIW10K	2020	10,000	10,000	YouTube	Second
KoDF	2021	62,166	1,75,776	Volunteers	First
VideoForensicsHQ	2021	261	1737	YouTube	First
Wild Deepfake	2021	3,805	3,509	Internet	First

Table 5: *Deepfake Detection Dataset*

**RQ3: What are the various measures and metrics that can be utilized to determine the effectiveness of deepfake detection?**

Depending upon the metrics discussed in subsection 3.3.1 a comparative Table 6 has been made for four discussed approaches for deepfake detection.

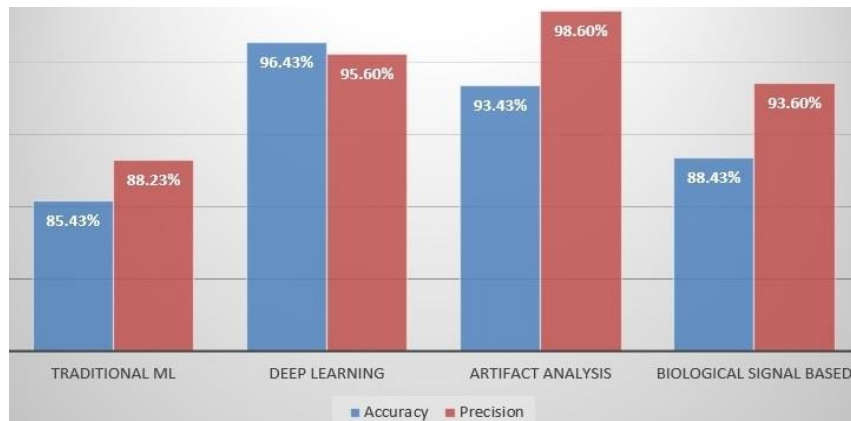
From the Table 6 it can be concluded that TML are simple and faster however very limited ability to handle complex images and videos with an average accuracy of 85%. DL methods on the other hand have this capability with average accuracy of 96%. Artifact-based methods outperforms others in terms of precision as shown in Figure 8. Accuracy of BSB method is higher than TML but lower than other two due to quality

of the physiological sensors used to collect the signals and presence of noise. The winner of all methods where speed meets performance is artifact-based methods. Their accuracy is closer to DL.

Category	Metric	Min	Max	Avg	Std
TML	Accuracy	52.31%	95.68%	85.43%	11.26
	Precision	56.11%	97.48%	88.23%	15.06
	Recall	51.01%	94.38%	83.63%	12.56
	F1-Score	0.498	0.998	0.9268	13.04
DL	Accuracy	70.31%	100%	96.43%	7.26
	Precision	71.11%	100%	95.60%	13.06
	Recall	69.01%	100%	93.63%	7.56
	F1-Score	0.648	1.00	0.887	9.04
Artifact Analysis	Accuracy	81.31%	99.80%	93.43%	5.26
	Precision	82.11%	100.00%	98.60%	10.06
	Recall	82.01%	97.98%	97.63%	3.56
	F1-Score	0.7876	0.989	90.70%	5.04
BSB	Accuracy	69.31%	94.80%	88.43%	11.26
	Precision	66.11%	95.00%	93.60%	16.06
	Recall	67.01%	91.98%	91.63%	8.56
	F1-Score	63.76%	93.90%	84.70%	11.04

Table 6: Metric Evaluation of Deepfake Detection Approaches

Figure 8: Average Accuracy and Precision for Deepfake Detection Methods



#### **RQ4: What is the future scope of deepfake detection?**

Deepfake detection is a rapidly evolving field, but there are still a number of challenges that need to be addressed in order to improve the accuracy and reliability of these approaches. Here are some of the key challenges:

- (i) **Limited Availability of High-Quality and Diverse Datasets:** Deepfake detection hindered by scarce high-quality, diverse datasets. Artifacts like splicing borders, low-quality faces impede algorithm effectiveness. Vital dataset enhancement needed [Siwei Lyu et al., 20].
- (ii) **Scalability:** Researchers face a challenge with limited high-quality datasets. Scaling issues in current DL approaches hinder effective detection of Deep Fake techniques [Katarya et al., 20].
- (iii) **Poor Quality Datasets and Limited Real-World Relevance:** Existing datasets for training deepfake detection lack real-world relevance due to poor visual quality. Performance on these may not translate to practical success. Limited diversity hampers algorithm efficiency [Celebi et al., 22].
- (iv) **Computational Optimization:** DL struggles to keep up with escalating Deepfake quality. Algorithmic upgrades needed for effective recognition. Optimal layers and architecture uncertain [Trinh & Liu, 21].
- (v) **Accuracy Optimization:** The Deep Fake detection algorithm had a 65% accuracy and only identified 1/3rd of the Deep Fakes. 50% misclassification of real videos, 50% undetected, 35% false positives. Prioritize accuracy and efficiency.
- (vi) **Multiclass, cross-label, and localized recognition:** Detecting Deep Fakes limited by binary categorization. Multiclass, cross-label, and localized recognition essential for detailed identification in complex scenarios.

The future of deepfake detection involves synergizing artifact-based and DL methods to overcome individual limitations. While artifact-based methods offer advantages, their dependency on specific creation processes necessitates integration with more versatile deep learning techniques. Addressing dataset bias and enhancing robustness are crucial for improving overall effectiveness in countering advanced deepfake techniques.

## **5 Conclusion**

This survey examines the current state-of-the-art methods for detecting deepfake, various datasets available for validation of methods, metrics for measuring performance of approaches. It also categorizes the existing approaches into 5 major groups. Additionally, it compares these categorized approaches in terms of various performance measuring parameters. The following summarization is provided:

- FaceForensics++ is one of the most popular datasets used for deep fake research.
- Detection accuracy is the most widely used performance metric.
- Experimental results show that DL techniques are effective in detecting Deepfake. DL models outperform non-DL models in terms of accuracy (96.34 %) whereas artifact analysis-based methods in terms of precision (98.60%).



## References

- [Ismail, 21] Ismail, M. Elpeltagy, M.S. Zaki & K. Eldahshan. A New Deep Learning-Based Methodology for Video Deepfake Detection Using XG Boost. *Sensors* 21,5413, 2021. <https://doi.org/10.3390/s21165413>
- [Agarwal, 21] Agarwal, H., Singh, A., & Rajeswari, D. Deepfake Detection Using SVM. In 2021 Second International Conference on Electronics and Sustainable Communication Systems (ICESC) (pp. 1245-1249), August 2021. IEEE. <https://doi.org/10.1109/ICESC51422.2021.9532627>
- [Agarwal, 19] Agarwal, S., & Varshney, L. R. Limits of deepfake detection: A robust estimation viewpoint. arXiv preprint arXiv:1905.03493. 2019. <https://doi.org/10.48550/arXiv.1905.03493>
- [Almars, 21] Almars, A. M. Deepfakes detection techniques using deep learning: a survey. *Journal of Computer and Communications*, 9(5), 20-35. 2021. <https://doi.org/10.4236/jcc.2021.95003>
- [Aromataris, 14] Aromataris, E., & Pearson, A. The systematic review: an overview. *AJN The American Journal of Nursing*, 114(3), 53-58. 2014. <https://doi.org/10.1097/01.NAJ.0000444496.24228.2c>
- [Bitouk, 08] Bitouk, D., Kumar, N., Dhillon, S., Belhumeur, P., & Nayar, S. K. Face swapping: automatically replacing faces in photographs. In *ACM SIGGRAPH 2008 papers* (pp. 1-8). <https://doi.org/10.1145/1399504.1360638>
- [Boccignone, 22] Boccignone, G., Bursic, S., Cuculo, V., D'Amelio, A., Grossi, G., Lanzarotti, R., & Patania, S. DeepFakes Have No Heart: A Simple rPPG-Based Method to Reveal Fake Videos. In *Image Analysis and Processing - ICIAP 2022: 21st International Conference, Lecce, Italy, May 23-27, 2022, Proceedings, Part II* (pp. 186-195). Cham: Springer International Publishing. [https://doi.org/10.1007/978-3-031-06430-2\\_16](https://doi.org/10.1007/978-3-031-06430-2_16)
- [Cao, 22] Cao, X., & Gong, N. Z. Understanding the Security of Deepfake Detection. In *International Conference on Digital Forensics and Cyber Crime* (pp. 360-378) 2022. Springer, Cham. [https://doi.org/10.1007/978-3-031-06365-7\\_22](https://doi.org/10.1007/978-3-031-06365-7_22)
- [Celebi, 22] Celebi, N., Liu, Q., & Karatoprak, M. A Survey of Deep Fake Detection for Trial Courts 2022. arXiv preprint arXiv:2205.15792. <https://doi.org/10.48550/arXiv.2205.15792>
- [Chen, 21] Chen, H. S., Rouhsedaghat, M., Ghani, H., Hu, S., You, S., & Kuo, C. C. J. Defakehop: A light-weight high-performance deepfake detector. In 2021 IEEE International Conference on Multimedia and Expo (ICME) (pp. 1-6) July 2021. IEEE. <https://doi.org/10.1109/ICME51207.2021.9428361>
- [Ciftci, 22] Ciftci, U. A., & Demir, I. How Do Deepfakes Move? Motion Magnification for Deepfake Source Detection. 2022. <https://doi.org/10.48550/arXiv.2212.14033>
- [Ciftci, 20] Ciftci, U. A., Demir, I., & Yin, L. Fakecatcher: Detection of synthetic portrait videos using biological signals. *IEEE transactions on pattern analysis and machine intelligence*. 2020. <https://doi.org/10.1109/TPAMI.2020.3009287>
- [Ciftci, 20] Ciftci, U. A., Demir, I., & Yin, L. How do the hearts of deep fakes beat? deep fake source detection via interpreting residuals with biological signals. In 2020 IEEE international joint conference on biometrics (IJCB) (pp. 1-10) September 2020. IEEE. <https://doi.org/10.1109/IJCB48548.2020.9304909>

- [Creswell, 18] Creswell, A., White, T., Dumoulin, V., Arulkumaran, K., Sengupta, B., & Bharath, A. A. Generative adversarial networks: An overview. *IEEE signal processing magazine*, 35(1), 53-65 2018. <https://doi.org/10.1109/MSP.2017.2765202>
- [Das, 04] Das K., Jiang j., Rao J.N.K. Mean squared error of empirical predictor." *Ann. Statist.* 32 (2) 818 - 840, 2004. <https://doi.org/10.1214/009053604000000201>.
- [De Lima, 20] De Lima, O., Franklin, S., Basu, S., Karwoski, B., & George, A. Deepfake detection using spatiotemporal convolutional networks. 2020. arXiv preprint arXiv:2006.14749. <https://doi.org/10.48550/arXiv.2006.14749>
- [Ding, 21] Ding, F., Zhu, G., Li, Y., Zhang, X., Atrey, P. K., & Lyu, S. Anti-forensics for face swapping videos via adversarial training. *IEEE Transactions on Multimedia*, 24, 3429-3441 2021. <https://doi.org/10.1109/TMM.2021.3098422>
- [Dolhansky, 19] Dolhansky, B., Howes, R., Pflaum, B., Baram, N., & Ferrer, C. C. The deepfake detection challenge (dfdc) preview dataset. 2019. arXiv preprint arXiv:1910.08854. <https://doi.org/10.48550/arXiv.1910.08854>
- [Dong, 22] Dong, S., Wang, J., Ji, R., Liang, J., Fan, H., & Ge, Z. Towards A Robust Deepfake Detector: Common Artifact Deepfake Detection Model. 2022. arXiv preprint arXiv:2210.14457. <https://doi.org/10.48550/arXiv.2210.14457>
- [Dong, 22 October] Dong, S., Wang, J., Liang, J., Fan, H., & Ji, R. Explaining Deepfake Detection by Analysing Image Matching. In *Computer Vision – ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XIV* (pp. 18-35). Cham: Springer Nature Switzerland. 2022, October. [https://doi.org/10.1007/978-3-031-19781-9\\_2](https://doi.org/10.1007/978-3-031-19781-9_2)
- [Dufour,19] Dufour, N., & Gully, A. Contributing data to deepfake detection research. *Google AI Blog*, 1(3) 2019.
- [Durall, 19] Durall, R., Keuper, M., Pfreundt, F. J., & Keuper, J. Unmasking deepfakes with simple features. 2019 arXiv preprint arXiv:1911.00686. <https://doi.org/10.48550/arXiv.1911.00686>
- [Elhassan, 22] Elhassan, A., Al-Fawa'reh, M., Jafar, M. T., Ababneh, M., & Jafar, S. T. DFT-MF: Enhanced deepfake detection using mouth movement and transfer learning. *SoftwareX*, 19, 101115. 2022. <https://doi.org/10.1016/j.softx.2022.101115>
- [Fang, 20] Fang, Y., Zhu, H., Zeng, Y., Ma, K., & Wang, Z. Perceptual quality assessment of smartphone photography. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 3677-3686) 2020. <https://doi.org/10.1109/CVPR42600.2020.00373>
- [Fox, 21] Fox, G., Liu, W., Kim, H., Seidel, H. P., Elgharib, M., & Theobalt, C. VideoforensicsHQ: Detecting high-quality manipulated face videos. In *2021 IEEE International Conference on Multimedia and Expo (ICME)* (pp. 1-6) July 2021. IEEE. <https://doi.org/10.1109/ICME51207.2021.9428101>
- [Gamage, 22] Gamage, D., Ghasiya, P., Bonagiri, V., Whiting, M. E., & Sasahara, K.. Are deepfakes concerning? analyzing conversations of deepfakes on reddit and exploring societal implications. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (pp. 1-19) April 2022. <https://doi.org/10.1145/3491102.3517446>
- [Gnera, 18] Gnera, D., & Delp, E. J. Deepfake video detection using recurrent neural networks. In *2018 15th IEEE international conference on advanced video and signal-based surveillance (AVSS)* (pp. 1-6) November 2018. IEEE. <https://doi.org/10.1109/AVSS.2018.8639163>

- [Guarnera, 20] Guarnera, L., Giudice, O., & Battiato, S. Deepfake detection by analyzing convolutional traces. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops (pp. 666-667) 2020. [https://doi.org/ 10.1109/CVPRW50498.2020.00341](https://doi.org/10.1109/CVPRW50498.2020.00341)
- [He, 21] He, Y., Yu, N., Keuper, M., & Fritz, M. Beyond the spectrum: Detecting deepfakes via re-synthesis 2021. arXiv preprint arXiv:2105.14376. [https://doi.org/ 10.48550/arXiv.2105.14376](https://doi.org/10.48550/arXiv.2105.14376)
- [Hosler, 21] Hosler, B., Salvi, D., Murray, A., Antonacci, F., Bestagini, P., Tubaro, S., & Stamm, M. C. Do deepfakes feel emotions? A semantic approach to detecting deepfakes via emotional inconsistencies. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 1013-1022) 2021. [https://doi.org/ 10.1109/CVPRW53098.2021.00112](https://doi.org/10.1109/CVPRW53098.2021.00112)
- [Hu, 22] Hu, J., Liao, X., Liang, J., Zhou, W., & Qin, Z. Finfer: Frame inference-based deepfake detection for high-visual-quality videos. In Proceedings of the AAAI conference on artificial intelligence (Vol. 36, No. 1, pp. 951-959) June 2022. <https://doi.org/10.1609/aaai.v36i1.19978>
- [Jain, 21] Jain, A., Korshunov, P., & Marcel, S. Improving generalization of deepfake detection by training for attribution. In 2021 IEEE 23rd International Workshop on Multimedia Signal Processing (MMSP) (pp. 1-6) October 2021. IEEE. [https://doi.org/ 10.1109/MMSP53017.2021.9733468](https://doi.org/10.1109/MMSP53017.2021.9733468)
- [Jiang, 20] Jiang, L., Li, R., Wu, W., Qian, C., & Loy, C. C. Deepforensics-1.0: A large-scale dataset for real-world face forgery detection. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 2889-2898) 2020. [https://doi.org/ 10.48550/arXiv.2001.03024](https://doi.org/10.48550/arXiv.2001.03024)
- [Jin, 21] Jin, X., Ye, D., & Chen, C. Countering spoof: towards detecting deepfake with multidimensional biological signals. Security and Communication Networks, 2021, 1-8. [https://doi.org/ 10.1155/2021/6626974](https://doi.org/10.1155/2021/6626974)
- [Johnson, 21] Johnson, D. G., & Diakopoulos, N. What to do about deepfakes. Communications of the ACM, 64(3), 33-35 2021. <https://doi.org/10.1145/3447255>
- [Jordan, 15] Jordan, M. I., & Mitchell, T. M. Machine learning: Trends, perspectives, and prospects. Science, 349(6245), 255-260 2015. [https://doi.org/ 10.1126/science.aaa8415](https://doi.org/10.1126/science.aaa8415)
- [Katarya, 20] Katarya, R., & Lal, A. A study on combating the emerging threat of deepfake weaponization. In 2020 Fourth International Conference on I-SMAC (IoT in Social, Mobile, Analytics, and Cloud) (I-SMAC) (pp. 485-490) October 2020. IEEE. [https://doi.org/ 10.1109/I-SMAC49090.2020.9243588](https://doi.org/10.1109/I-SMAC49090.2020.9243588)
- [Khalil, 21] Khalil, H. A., & Maged, S. A. Deepfakes creation and detection using deep learning. In 2021 International Mobile, Intelligent, and Ubiquitous Computing Conference (MIUCC) (pp. 1-4) May 2021. IEEE. [https://doi.org/ 10.1109/MIUCC52538.2021.9447642](https://doi.org/10.1109/MIUCC52538.2021.9447642)
- [Khan, 21] Khan, S. A., Artusi, A., & Dai, H. Adversarially robust deepfake media detection using fused convolutional neural network predictions. 2021 arXiv preprint arXiv:2102.05950. [https://doi.org/ 10.48550/arXiv.2102.05950](https://doi.org/10.48550/arXiv.2102.05950)
- [Kharbat, 19] Kharbat, F. F., Elamsy, T., Mahmoud, A., & Abdullah, R. Image feature detectors for deepfake video detection. In 2019 IEEE/ACS 16th International Conference on Computer Systems and Applications (AICCSA) (pp. 1-4) November 2019. IEEE. [https://doi.org/ 10.1109/AICCSA47632.2019.9035360](https://doi.org/10.1109/AICCSA47632.2019.9035360)
- [Khomali, 21] Khomali, A. and Yuan, J-S. ADD: Attention-Based DeepFake Detection Approach, Comput.5,49 2021. [https://doi.org/ 10.3390/bdcc5040049](https://doi.org/10.3390/bdcc5040049)

- [Kietzmann, 20] Kietzmann, J., Lee, L. W., McCarthy, I. P., & Kietzmann, T. C. Deepfakes: Trick or treat?. *Business Horizons*, 63(2), 135-146 2020. <https://doi.org/10.1016/j.bushor.2019.11.006>
- [Köbis ,21] Köbis, N. C., Doležalová, B., & Soraperra, I. Fooled twice: People cannot detect deepfakes but think they can. *Iscience*, 24(11), 103364 2021. <https://doi.org/10.1016/j.isci.2021.103364>
- [Korhonen, 12 July] Korhonen, J., & You, J. Peak signal-to-noise ratio revisited: Is simple beautiful?. In 2012 Fourth International Workshop on Quality of Multimedia Experience (pp. 37-38) July 2012. IEEE. <https://doi.org/10.1109/QoMEX.2012.6263880>
- [Korshunov, 19] Korshunov, P., & Marcel, S. Vulnerability assessment and detection of deepfake videos. In 2019 International Conference on Biometrics (ICB) (pp. 1-6) June 2019. IEEE. <https://doi.org/10.1109/ICB45273.2019.8987375>
- [Kwon, 21] Kwon, P., You, J., Nam, G., Park, S., & Chae, G. Kodf: A large-scale korean deepfake detection dataset. In Proceedings of the IEEE/CVF International Conference on Computer Vision (pp. 10744-10753) 2021. <https://doi.org/10.1109/ICCV48922.2021.01057>
- [L.Zhao, 21] L.Zhao, M. Zhang & H. Ding & X. Cui . MFF-Net: Deepfake Detection Network Based on Multi-Feature Fusion, *Entropy* 23,1692 2021. <https://doi.org/10.3390/e23121692>
- [Li, 22] Li, X., Ni, R., Yang, P., Fu, Z., & Zhao, Y. Artifacts-Disentangled Adversarial Learning for Deepfake Detection. 2022. *IEEE Transactions on Circuits and Systems for Video Technology*. <https://doi.org/10.1109/TCSVT.2022.3217950>
- [Li, 20] Li, Y., Yang, X., Sun, P., Qi, H., & Lyu, S. Celeb-df: A large-scale challenging dataset for deepfake forensics. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 3207-3216) 2020. <https://doi.org/10.1109/CVPR42600.2020.00327>
- [Lorena, 08] Lorena, A. C., De Carvalho, A. C., & Gama, J. M. A review on the combination of binary classifiers in multiclass problems. *Artificial Intelligence Review*, 30, 19-37 2008. <https://doi.org/10.1007/s10462-009-9114-9>
- [Lyu, 20] Lyu, S. Deepfake detection: Current challenges and next steps. In 2020 IEEE international conference on multimedia & expo workshops (ICMEW) (pp. 1-6) July 2020. IEEE. <https://doi.org/10.1109/ICMEW46912.2020.9105991>
- [Lyu, 22] Lyu, S. Deepfake detection. In *Multimedia Forensics* (pp. 313-331). Singapore: Springer Singapore. H. T. Sencar et al. (eds.), *Multimedia Forensics, Advances in Computer Vision and Pattern Recognition 2022*. [https://doi.org/10.1007/978-981-16-7621-5\\_12](https://doi.org/10.1007/978-981-16-7621-5_12)
- [Matern, 19] Matern, F., Riess, C., & Stamminger, M. Exploiting visual artifacts to expose deepfakes and face manipulations. In 2019 IEEE Winter Applications of Computer Vision Workshops (WACVW) (pp. 83-92) January 2022. IEEE. <https://doi.org/10.1109/WACVW.2019.00020>
- [Mathews, 23] Mathews, S., Trivedi, S., House, A., Povolny, S., & Fralick, C. An explainable deepfake detection framework on a novel unconstrained dataset. *Complex & Intelligent Systems*, 1-13 2023. <https://doi.org/10.1007/s40747-022-00956-7>
- [Mirsky, 21] Mirsky, Y., & Lee, W. The creation and detection of deepfakes: A survey. *ACM Computing Surveys (CSUR)*, 54(1), 1-41 2021. <https://doi.org/10.1145/3425780>
- [Mittal, 20] Mittal, T., Bhattacharya, U., Chandra, R., Bera, A., & Manocha, D. Emotions don't lie: An audio-visual deepfake detection method using affective cues. In Proceedings of the 28th ACM international conference on multimedia (pp. 2823-2832) October 2020. <https://doi.org/10.1145/3394171.3413570>

- [Nagrani, 17] Nagrani, A., Chung, J. S., & Zisserman, A. Voxceleb: a large-scale speaker identification dataset. 2017. arXiv preprint arXiv:1706.08612. <https://doi.org/10.48550/arXiv.1706.08612>
- [Nguyen, 20] Nguyen, H. M., & Derakhshani, R. Eyebrow recognition for identifying deepfake videos. In 2020 international conference of the biometrics special interest group (BIOSIG) (pp. 1-5) September, 2020. IEEE.
- [Nguyen, 22] Nguyen, T. T., Nguyen, Q. V. H., Nguyen, D. T., Nguyen, D. T., Huynh-The, T., Nahavandi, S., Nguyen, T.T., Pham, Q.V. and Nguyen, C.M. Deep learning for deepfakes creation and detection: A survey. *Computer Vision and Image Understanding*, 223, 103525 2022. <https://doi.org/10.1016/j.cviu.2022.103525>
- [Nirkin, 20] Nirkin, Y., Wolf, L., Keller, Y., & Hassner, T. Deepfake detection based on the discrepancy between the face and its context. 2020. arXiv preprint arXiv:2008.12262. <https://doi.org/10.48550/arXiv.2008.12262>
- [Nirkin, 21] Nirkin, Y., Wolf, L., Keller, Y., & Hassner, T. DeepFake detection is based on discrepancies between faces and their context. *IEEE Transactions on Pattern Analysis and Machine Intelligence* .2021. <https://doi.org/10.1109/TPAMI.2021.3093446>
- [Perov, 20] Perov, I., Gao, D., Chervoniy, N., Liu, K., Marangonda, S., Umé, C., ... & Zhang, W. DeepFaceLab: Integrated, flexible and extensible face-swapping framework. 2020. arXiv preprint arXiv:2005.05535. <https://doi.org/10.48550/arXiv.2005.05535>.
- [Rafique, 21] Rafique, R., Nawaz, M., Kibriya, H., & Masood, M. Deepfake detection using error level analysis and deep learning. In 2021 4th International Conference on Computing & Information Sciences (ICCIS) (pp. 1-4) November 2021. IEEE. <https://doi.org/10.1109/ICCIS54243.2021.9676375>
- [Rana, 20] Rana, M. S., & Sung, A. H. Deepfakestack: A deep ensemble-based learning technique for deepfake detection. In 2020 7th IEEE international conference on cyber security and cloud computing (CSCloud)/2020 6th IEEE international conference on edge computing and scalable cloud (EdgeCom) (pp. 70-75) August 2020. IEEE. <https://doi.org/10.1109/CSCloud-EdgeCom49738.2020.00021>
- [Rana, 21] Rana, M. S., Murali, B., & Sung, A. H. Deepfake Detection Using Machine Learning Algorithms. In 2021 10th International Congress on Advanced Applied Informatics (IIAI-AAI) (pp. 458-463) July 2021. IEEE. <https://doi.org/10.1109/IIAI-AAI53430.2021.00079>
- [Ribeiro, 11] Ribeiro, F., Florêncio, D., Zhang, C., & Seltzer, M. Crowdmos: An approach for crowdsourcing mean opinion score studies. In 2011 IEEE international conference on acoustics, speech and signal processing (ICASSP) (pp. 2416-2419) May 2011. IEEE. <https://doi.org/10.1109/ICASSP.2011.5946971>
- [Rossler, 19] Rossler, A., Cozzolino, D., Verdoliva, L., Riess, C., Thies, J., & Nießner, M. Faceforensics++: Learning to detect manipulated facial images. In Proceedings of the IEEE/CVF international conference on computer vision (pp. 1-11) 2019. <https://doi.org/10.1109/ICCV.2019.00009>
- [Rossler, 18] Rössler, A., Cozzolino, D., Verdoliva, L., Riess, C., Thies, J., & Nießner, M. Faceforensics: A large-scale video dataset for forgery detection in human faces. 2018. arXiv preprint arXiv:1803.09179. <https://doi.org/10.48550/arXiv.1803.09179>
- [Sættra, 23] Sættra, H. S. Generative AI: Here to stay, but for good?. *Technology in Society*, 75, 102372.2023. <https://doi.org/10.1016/j.techsoc.2023.102372>

- [Streijl, 16] Streijl, R. C., Winkler, S., & Hands, D. S. Mean opinion score (MOS) revisited: methods and applications, limitations and alternatives. *Multimedia Systems*, 22(2), 213-227 2016. <https://doi.org/10.1007/s00530-014-0446-1>
- [Sun, 20] Sun, X., Wu, B., & Chen, W. Identifying invariant texture violation for robust deepfake detection. 2020. arXiv preprint arXiv:2012.10580. <https://doi.org/10.48550/arXiv.2012.10580>
- [Tolosana, 22] Tolosana, R., Romero-Tapiador, S., Vera-Rodriguez, R., Gonzalez-Sosa, E., & Fierrez, J. DeepFakes detection across generations: Analysis of facial regions, fusion, and performance evaluation. *Engineering Applications of Artificial Intelligence*, 110, 104673 2022. <https://doi.org/10.1016/j.engappai.2022.104673>
- [Trinh, 21] Trinh, L., & Liu, Y. An examination of fairness of AI models for deepfake detection. 2021. arXiv preprint arXiv:2105.00558. <https://doi.org/10.48550/arXiv.2105.00558>
- [Vinay, 22] Vinay, A., Bhat, N., Khurana, P. S., Lakshminarayanan, V., Nagesh, V., Natarajan, S., & Sudarshan, T. B. AFMB-Net: DeepFake Detection Network Using Heart Rate Analysis. *Tehnički glasnik*, 16(4), 503-508 2022. <https://doi.org/10.31803/tg-20220403080215>
- [Wahl-Jorgensen, 21] Wahl-Jorgensen, K., & Carlson, M. Conjecturing fearful futures: Journalistic discourses on deepfakes. *Journalism practice*, 15(6), 803-820 2021. <https://doi.org/10.1080/17512786.2021.1908838>
- [Wang, 22] Wang, J., Wu, Z., Ouyang, W., Han, X., Chen, J., Jiang, Y. G., & Li, S. N. M2tr: Multi-modal multi-scale transformers for deepfake detection. In *Proceedings of the 2022 International Conference on Multimedia Retrieval* (pp. 615-623) June 2022. <https://doi.org/10.1145/3512527.3531415>
- [Wang, 04] Wang, Z., Bovik, A. C., Sheikh, H. R., & Simoncelli, E. P. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4), 600-612, 2004. <https://doi.org/10.1109/TIP.2003.819861>
- [Wolter, 21] Wolter, M., Blanke, F., Hoyt, C. T., & Garcke, J. Wavelet-packet powered deepfake image detection. 2021. arXiv preprint arXiv:2106.09369.
- [Yadlin-Segal, 21] Yadlin-Segal, A., & Oppenheim, Y. Whose dystopia is it anyway? Deepfakes and social media regulation. *Convergence*, 27(1), 36-51. 2021. <https://doi.org/10.1177/13548565209239>
- [Zhang, 22] Zhang, J., Cheng, K., Sovrnigo, G., & Lin, X. A Heterogeneous Feature Ensemble Learning based Deepfake Detection Method. In *ICC 2022-IEEE International Conference on Communications* (pp. 2084-2089) May 2022. IEEE. <https://doi.org/10.1109/ICC45855.2022.9838630>
- [Zhao, 21] Zhao, H., Zhou, W., Chen, D., Wei, T., Zhang, W., & Yu, N. Multi-attentional deepfake detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 2185-2194) 2021. <https://doi.org/10.1109/CVPR46437.2021.00222>
- [Zhao, 21] Zhao, T., Xu, X., Xu, M., Ding, H., Xiong, Y., & Xia, W. Learning self-consistency for deepfake detection. In *Proceedings of the IEEE/CVF international conference on computer vision* (pp. 15023-15033) 2021. <https://doi.org/10.1109/ICCV48922.2021.01475>
- [Zhou, 21] Zhou, T., Wang, W., Liang, Z., & Shen, J. (2021). Face forensics in the wild. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 5778-5788) 2021. <https://doi.org/10.1109/CVPR46437.2021.00572>