

A Survey of Papers Regarding Deep Fake Image and Video Detection

Avani PS
Department of Computer Science
College of Engineering
Karunagappally Kerala, India

Athul George
Department of Computer Science
College of Engineering
Karunagappally Kerala, India

Artsun T Kurian
Department of Computer Science
College of Engineering
Karunagappally Kerala, India

Anoop V
Department of Computer Science
College of Engineering
Karunagappally Kerala, India

Geethu Gopakumar
Department of Computer Science
College of Engineering
Karunagappally Kerala, India

Swathi S
Department of Computer Science
College of Engineering
Karunagappally Kerala, India

Abstract—This modern-day artificial intelligence technology has a worldwide threat that is the generation of hyperrealistic manipulated images and videos, often true to the last detail. The deepfake phenomenon has alarmed people from various walks of life, such as journalism, social networking, and cybersecurity, and has necessitated the provision of more sophisticated detection techniques. The review of approaches to finding the latest and modern deepfakes has been prepared based on the works from the top leading research studies. In the latest developments, various types of deep learning architecture and hybrid models are being used in their attempt to address image forgery and video forgery on various levels.

Improved Dense Convolutional Neural Networks (CNNs) have more compelling feature extraction capability in texture and illumination against manipulation-resisted images. Dual Attention Network is capable of fighting for both spatial and temporal attention mechanisms and detects forgery in the face using details at the frame level and using the whole connection between multiple frames. Video Vision Transformers (ViViTs) has an improved accuracy and efficiency in video processing when augmented with facial landmark detection and Depthwise Separable Convolutions techniques. Other approaches include the application of the Xception network with dual attention mechanisms and features fusion, which find more subtle manipulation artifacts by looking at the spatial and channel-wise attributes of fabricated images.

There are new and very different ways to address problems re-garding dataset diversity and model generalization. Graph-based deepfake detection frameworks, for instance, can harness the power of self-supervised learning to lessen dependency on labeled data, while on their end, convolutional architectures with Gabor filters could actually improve texture and frequency analysis because of the integrated features. These methods achieve a high detection accuracy on datasets like FaceForensics++, Celeb-DF, and DFDC, showing a considerable capacity for handling diverse deepfake scenarios. However, issues of computational complexity, real-time scalability, and the rapid pace of evolution of GANs continue to exist.

This survey thus points to the need for the development of rich and varied datasets that would help improve the generalizability of detection models against emerging deepfake such as those from StyleGAN2. Hybrid models that take on a combination of spatial and temporal plus trace-based feature extraction also present a very bright future for developing detection precision and efficiency. Future research is thus focusing within the domain of improving algorithms for real-time applications with less computational overhead, in addition to standardizing evaluation benchmark development so as to achieve comparable cross-method assessments. Advances in deepfake-detection work, therefore, become instrumental in fighting misinformation and enabling trust in digital media in an AI-powered age by con-fronting those issues.

INTRODUCTION

Dramatic rise of the deepfake technology and the method in which AI-based media threatens the integrity of digital content are the leading issues in today's world. The proposed paper therefore looks into almost all the most recent methods for detecting deep fakes, focusing more on the recent advancements in the field of image and video analytics as documented in 11 vital research papers. The Video Vision Transformer (ViViT) model was modified by Patel et al. using facial landmark detection, Depthwise Separable Convolutions, and self-attention mechanisms. The method successfully detects spatial-temporal inconsistencies across video frames, thus improving both accuracy and computational efficiency. Their method significantly improved performance on various datasets, including Celeb-DF and FaceForensics++.

Patel et al. have proposed an Enhanced Dense Convolutional Neural Network that connects neurons within a single layer to densify the connections in the network and improve the extraction of features and gradient flows. The architecture excellently performs as manipulation artifact identification based on texture and lighting inconsistencies in images. This method achieves detection accuracy ranging from 85 to 98

percent for datasets such as FaceForensics++. Similar to this, Luo and Chen then presented a dual attention network with integrated spatial and temporal attention for detecting facial forgery in videos. Their approach achieved great

detection accuracy in several datasets, including FaceForensics++ and DFDC, as it got to analyze both frame-level details and temporal relationships.

Lin et al. enhanced the Xception network through dual attention mechanism and feature fusion for deepfake detection enhancement. This method emphasizes spatial and channel-wise features improving the ability of the model detection and subtle manipulation artifacts above others. Although computational complexity was increased, their model performed superiorly on standard deepfake datasets. Kim and Cho proposed another hybrid method combining content-based with trace-based feature extraction capturing both high-level visual content and subtle traces of manipulation. This method also achieved some state-of-the-art results in Celeb-DF and DFDC datasets.. Patel et al. point out the limitations of existing deepfake detection techniques and call for building stronger datasets for countering the advance challenges posed by newer deepfake creation techniques, such as StyleGAN2. The researchers made a comparative study of the recent techniques in deepfake detection focusing on performance evaluation on parameters, including accuracy, precision, and recall. They also strongly put forth the urgency for developing visualization tools that will help interpret a model's decisions and the necessity of more diverse datasets for a generalizable model. To detect local inconsistency in frequency with the dynamic model proposed by the researchers for consistency analysis of videos high potential detection even through superior-quality deepfakes was achieved only with consideration of local regions. Incorporating Gabor filters into CNNs has become a technique introduced by Khalifa et al. for achieving frequency texture analysis. This particular technique may be said to yield exemplary performance even on templates with a huge number of entries like that of Celeb-DF and FaceForensics++. Nonetheless, its computational cost is too much for real-time usage. In a systematic review on deepfake recognition, Rana et al. revealed the gaps that exist in the diversity of datasets and evaluation standards, reinforcing the urgent need for common benchmarks to better evaluate model capabilities. Follow the same tenor as with Khormali-Yuan, developing a self-supervised Graph Transformer architecture that examines deepfakes via the identification of spatial-temporal imbalances without needing market data. That said, it is far better but resource-hungry and doesn't allow real-time application.

LITERATURE REVIEW

I. IMPROVING VIDEO VISION TRANSFORMER FOR DEEPFAKE VIDEO DETECTION USING FACIAL LANDMARK, DEPTHWISE SEPARABLE CONVOLUTION AND SELF ATTENTION

Research has been completed under the study titled "Improving Video Vision Transformer for Detection of Deepfake Video with Facial Landmark, Depthwise Separable Convolution, and Self Attention" to make a mark in deepfake detection with the integration of video vision transformers along with facial landmark detection,

Depthwise Separable Convolution, and self-attention

strategies. It successfully incorporates prominent areas that help the model to temporally and spatially differentiate across video frames concentrating on facial landmarks. While DSC block efficiency increases the performance of ViViT, it reduces the computational power required across the blocks without harming the accuracy. Furthermore, the Convolution Block Attention Module (CBAM) enables improved visibility concerning minor differences across the frames, thus improving accuracy in deepfake video detection. Such a hybrid approach has shown magnificent performance on benchmark datasets of deepfake videos such as FaceForensics++ and Celeb-DF.

The spacetime characteristics of a video deepfake approach will involve further enhancements over the already brilliant existing methods in capturing small changes that other methods cannot detect. The face landmark-based representation of input videos highlights the most relevant parts of the video and thus makes the model more resilient in practice. Real-time implementation is now feasible in that the model achieves much higher accuracies while being less computationally expensive via advanced convolutional layers and attention mechanisms. This is an effective and important approach toward combating this rising threat of deepfake technology.

II. AN IMPROVED DENSE CNN

ARCHITECTURE FOR DEEPFAKE IMAGE DETECTION

The work is also an improvement over deepfake detection with regard to the CNN model since the architecture is designed to use dense connections, such that the feature reuse is enhanced, overfitting reduced, and gradient flow, which inevitably aids in the detection of amongst others subtle artifacts such as texture and lighting inconsistencies in manipulated images, improved. The architecture has manifested high accuracy and robustness, is hence effective on the assorted datasets such as FaceForensics++ and Celeb-DF.

The complex architecture also helps the model generalize better and ensure effective performance across categories deepfake manipulations. Moreover, by greatly enhancing feature extraction, this model overcame the performance of a conventional CNN model regarding detection capabilities. It indicates real potential for practical applications in the verification of deepfake images..

III. DUAL ATTENTION NETWORK APPROACHES TO FACE FORGERY VIDEO DETECTION

This paper proposes a dual attention network, which has proposed the integration of spatial and temporal attention mechanisms for detecting face forgery in videos. The proposed model has shown considerable advancements in accuracy, emphasizing spatial features, such as still individual frames, and temporal features, such as the interrelationship among frames, thus recognizing subtle differences in deepfake videos. The approach has been validated on popular deepfake video datasets, such as FaceForensics++ and DFDC, and has shown

that the dual attention approach enhances detection effectiveness. The accuracy is significantly improved over previous methods, thus indicating that prior work has been better able to recognize a forgery. Rates of accuracy for artifacts of different manipulation types approach nearly 95 percent. This method has a few benefits, such as being exceedingly robust and accurate, particularly against subtle deepfake artifacts over time since it handles spatial as well as temporal features at the same time. Therefore, this model surpasses most of the conventional ones, which get biased on spatial aspects or temporal aspects. The dual attention network also shows a remarkable generalization capability over various datasets and face forgery manipulations types, making the network even more flexible. Yet, the intricacy embedded in this model might also serve as a hindrance for the fine dual attention mechanism, which can prove computation-intensive thus hinder real-time deployment, especially for high-fidelity videos. The model performs efficiently on common benchmark datasets but is still deemed susceptible to emerging high-quality deepfake techniques, even those built on advanced GANs, such as StyleGAN2; forgeries of this type often appear to escape current methods of detection.

IV. IMPROVED XCEPTION WITH DUAL ATTENTION MECHANISM AND FEATURE FUSION FOR FACE FORGERY DETECTION

This competitive manuscript presents a version of the Xception network that is engineering for detecting faces counterfeit. The protection includes a dual attention mechanism and fusion of features to improve accuracy and robustness. Dual attention focuses spatially and from the channel side. While the fusion detects subtle inconsistencies in forged faces, the input also contains the information collected through different internal layers to enhance performance. The technique tested with standard datasets of face forgery is expected to deliver high accuracy compared to several modern ones. But then again, the same approach includes computational complexity-introducing aspects of dual attention and feature fusion, which might prove detrimental in speed or memory usage with regard to inference, especially when applied to real-time applications. Moreover, it can have difficulties when it comes to external novel face forgery techniques, even though it is excellent on benchmark datasets and its performance is nebulous under different forgery types in resource-based environments.

V. EXPOSING FAKE FACES THROUGH DEEP NEURAL NETWORKS COMBINING CONTENT AND TRACE FEATURE EXTRACTORS

The approach presented in this paper is a new way of detecting deepfakes, using models trained for hybrid content-and trace-based feature extraction with deep neural networks. The dual feature approach thus allows improving detection precision by having both high visual factors and other subtle indicators. Among the manipulations are artefacts most frequently seen in an image. Assessment of the method uses datasets such as the DeepFake Detection Challenge (DFDC)

and Celeb-DF. It proves to have state-of-the-art performance with regard to distinguishing real faces from shams. However, the method involves several problems, especially in the computational cost and the generalization of the model on newer or more diverse forgery techniques. It does well in existing datasets; however, it shows scalability and real-world applicability, especially in resource-constrained environments, remaining areas for further improvement.

VI. DEEPPAKE GENERATION AND DETECTION: CASE STUDY AND CHALLENGES

The current situation in deepfake production and hacking has been well detailed in the paper. Deepfake technology-in-use case studies as well as practical challenges faced are dissected minutely in the paper. Several different approaches for fake image generation and for detection are elaborated in this paper. Deepfakes: Advancing the limits and advantages of current detection methods. The authors tested several datasets like FaceForensics++ and DeepFakeTIMIT. They also have a discussion on the accuracy of detection models across these datasets and against the argument that achieving a score high in truth would often be challenging since evolving deepfake techniques. Perhaps one of the most pertinent merits of this research is its enormous extent in terms of the technical dimension and societal concerns, which-the research is not limited to, but also encompasses applications of deepfake technology-which gives considerable views into the real-life challenges that detection could face. Despite the above, the narrative also reveals some critical limitations: the absence of robustness with regard to new and novel unseen deepfake generations, and even in highly advanced detection methods, they have high computational cost. Moreover, The study shows that present datasets do not reflect the depth of deepfake diversity. manipulations limiting the ways how the models of detection could be applied toward real-world scenarios.

VII. PERFORMANCE COMPARISON AND VISUALIZATION OF AI-GENERATED- IMAGE DETECTION METHODS

This paper presents a vast juxtaposition between various techniques of AI-generated image detection. Methods emphasizing the importance of evaluating performance across different detection scenarios. The above methods are assessed by the authors using several benchmark datasets, including ImageNet and CelebA, while performance metrics such as accuracy, precision, and recall are detailed. to test their performance in the domain of detection of AI-generated images. One significant advantage A precise account of many latest advanced methods will be found, offering vital Insights into which models excel in which scenarios and why. And while the authors Models enable one to visualize techniques that allow interpretation of decisions made by models. transparency and an improved understanding. However, in this narrow scope, there are serious limits in that the datasets used for testing can best represent the richness of diversity in-built in AI generation. Some techniques used today may restrict the extension of research results. The paper Some challenges

are also mentioned in creating uniform detection accuracy across the two forms of images on which AI has been created, requiring further works towards developing more robust and adaptable techniques for detection.

VIII. LOCAL REGION FREQUENCY GUIDED DYNAMIC INCONSISTENCY NETWORK FOR DEEPFAKE VIDEO DETECTION

The paper presents a new mechanism for detecting deepfake videos with dynamic inconsistencies across a guest network controlled by frequency characteristics of local domains. This approach considers identifying inconsistencies with frequency patterns such as local regional differences in a video: an approach that helps strengthen the capacity of the model toward identifying subtle changes in a video. The proposed methodology is evaluated using benchmark dataset comparisons with FaceForensics++ and DeepFakeTIMIT. Their results show a high accuracy percentage in discriminating real videos from deepfakes. One of the most significant advantages of this method is its ability to clearly capture and harness local frequency differences, improving effectiveness even for highly manipulated videos. During this time, at the same time, the text refers to challenges in computation-related complexity, which requires a very huge amount of annotated training data. Future implementable methods such as these will carry out further investigation on their applicability in an event such as future deepfake attacks or less controlled scenarios (for example, low-resolution videos).

IX. CONVOLUTIONAL NEURAL NETWORK BASED ON DIVERSE GABOR FILTERS FOR DEEPFAKE RECOGNITION

The work proposes a deepfake detection technique based on the combined application of several Gabor filters within a convolutional neural network framework. It is worth pointing out that Gabor filters are known to be very efficient in extracting fine texture and frequency information, thereby improving the capability of the CNN to detect subtle artifacts present in deepfake manipulations. Such techniques were evaluated in accordance to the data present on standard deepfake detection datasets faceforensics++ and celeb-df and exhibit very high accuracy-rates in distinguishing between real and fake media. Some of the most important advantages that the expected method could tentatively yield are The ability to extract subtle frequency information from very different. Orientations and scales thus making it more robust to attacks from all deepfake generation methods. However, one limitation associated with this paper is the heavy computations incurred with the use of different Gabor filters that could very substantially slow down the continued processing and increase the consumption of resources, especially in real-time applications. In addition, although the method works well on established datasets, the efficacy against new or more sophisticated deepfake generating techniques. The methods of the information are still unknown and thus need more testing on a more diverse and demanding data.

X. DEEPFAKE DETECTION: A SYSTEMATIC LITERATURE REVIEW

Here extensive reviews of out of numerous methodologies available for identifying deep The imitations along with evaluation on different approaches with respect to accuracy, datasets adopted, advantages and disadvantages are presented here. Challenges. The analysis finds these algorithms for deep fake detection usually have quite high accuracy as well. Rates, in particular Convolution Neural Networks that use state-of-the-art deep learning techniques. works (CNNs) and Recurrent Neural Networks (RNNs), where performance usually relies on the quality and diversity of datasets used. FaceForensics++ , DeepFakeTIMIT some of the most widely quoted datasets, as often appearing for training and testing, although limitations pertaining to the size of dataset, representational bias, and the generalization of models across different domains are acknowledged. However, an enormous drawback is the lack of a unified, a uniform method of evaluation will bring in-consistencies within reported results and making comparisons difficult across studies.

XI. SELF-SUPERVISED DEEPFAKE DETECTION METHOD USING A GRAPH TRANSFORMER FRAMEWORK

In a research study of 2024, Khormali and Yuan advance a self-supervised detection mechanism of deepfake. An innovative method is presented through the Graph Transformer framework to source interdummy-relationship of features: it deals with the interlinks in relation to graph-based relationships relating to visual features of a digital medium. It takes advantage of the Graph Transformer model, which captures spatial-temporal interdependencies in enhancing the detection. Such an inconsistency is typical of the deepfake. This framework is self-supervised. It is to allow the model to learn features on unsupervised data so that it maximally improves scalability and adaptability in divergent datasets. The first was a graph-fa construction Social landmarks and transformer to analyze the relationships make it possible to detect subtle artifacts in detail. This method has very high accuracy in detection. showing the ability to withstand widespread manipulations and distortions. Nonetheless, Though the Graph Transformer can show a good percentage of detection rates, it is reduced dependence on Very costly in computational requirements and may limit its use. The main point here is that, in a real-time application.

DISCUSSION

It is seen that the study highlighted above has brought to the fore the innovative development of various deep learning methodologies for detecting this emerging technology's threat. Deepfake technology has become a challenge to journalism, social media, and security by creating realistic, fake media. The current state-of-the-art techniques, though exhibiting great

TABLE I
COMPARISON OF METHODS FOR DEEPPFAKE DETECTION.

Title	Author	Pros	Cons
Method	Overview	Strengths	Limitations
Deepfake Detection Techniques	Early methods focused on pixel-level analysis.	Simple to implement	Struggle with diverse deepfake types
Convolutional Neural Networks	Widely used for image classification.	Effective for images	Misses temporal features in videos
Recurrent Neural Networks	Designed for sequential data.	Good for capturing sequences	Limited by vanishing gradients
3D Convolutional Networks	Extend 2D convolutions into the time dimension.	Capture spatial and temporal features	Computationally intensive
Vision Transformers	Emerging in image analysis; applied to video is still developing.	Effective global context capture	Challenges with temporal coherence
Hybrid Approaches	Combine strengths of CNNs and RNNs or ViTs.	Improved detection accuracy	Increased complexity and resource demands

advances, have also shown gaps and need more exploration and innovation.

The reviewed techniques incorporated advanced state-of-the-art technologies that boost the detection effectiveness. Improved Dense Convolutional Neural Networks (CNNs) have their superior feature extraction capabilities allowing them to excel in the detection of fine artifacts like inconsistencies in texture and lighting anomalies. Video Vision Transformers (ViViT), blended with Depthwise Separable Convolutions and self-attention mechanisms, are shown to improve both spatial and temporal analysis relative to video-based detection. The focus on key features makes accuracy better while improving computational efficiency through facial landmark detection. Dual Attention Networks work by using mechanisms of spatial and temporal attention to analyze details contained in frames and relationships that exist between frames. These thus achieved high detection rates for video deepfakes.

Hybrid approaches incorporate spatial, temporal, and frequency analysis when analyzing an individual. Some popular models include applications of diverse Gabor filters and Xception networks with dual attention mechanisms that can recognize different forms of manipulations for use in samples like FaceForensics++, Celeb-DF, and DFDC. Moreover, the potential role of self-supervised methodologies through Graph Transformers in reducing reliance on labeled data resulting in high scalability and adaptability is notable.

Noteworthy by definition, methods still encounter hurdles, especially on computational intensity, thus making real-time applications impossible. Modalities such as dual attention mechanisms and graph-based frameworks would need heavy computational requirements and so would be found unsuitable in resource-constrained solutions. Generalizability to new deepfake generation techniques, for example, those generated using StyleGAN2, is also a huge concern. Although current data sets are comprehensive, there is no capturing of reality, leading to possible biases. Models fail most especially with low-resolution or poorly compressed media, and they may further perform very badly in uncontrolled environments.

The un-uniformity of benchmarking and common evaluation criteria therefore aggravate the comparison between approaches and make it difficult to find out their best works. And while most of those issues involve interpretability by detection models, it does limit trust and further adoption in sensitive fields, like forensic analysis and legal proceedings. And when it comes to these limitations, future advances in research should mainly focus on light architectures that are quite energy efficient yet on real-time level detection. Another essential aspect for advancing the generalizability of models is the collection of datasets that use other cultural contexts as well as different resolutions and compressions. Unified standard description benchmarks with standardized evaluation frameworks are then critical in creating meaningful comparisons so that they can help drive innovation. The other techniques like explainability based on artificial

intelligence, say using visualization tools, can complement the possibility that the model can be more transparent and increase user trust and adoption.

This warrants a multi-disciplinary approach, with researchers coming together with policymakers and the industry to avert the detection-tailored technology misuse effects. The serious restriction on ethical grounds must also be enforcing this technology for degrading the possible privacy-related contents in intelligence detection analysis measures. Various efforts are taken jointly to formulate effective impacting practices to deal with the misuse of deepfake technology while propagating public awareness and caution about the possible hazards and effects that it may bring.

CONCLUSION

The papers collated in this article highlight the most recent advancements in deepfake investigation through image and video analytics. Such methods involve using techniques like dense connections in CNNs, self-attention mechanisms, facial landmark detection, Depthwise Separable Convolutions, and temporal and spatial attention mechanisms to improve precision and efficiency in detecting manipulation artifacts at fine detail levels in images or videos. Even though these approaches are promising, there are still challenges, such as computational complexity, real-time detection, and generalization to new deepfake generation techniques. Future work can balance an optimal Model, which performs fast and accurate detection, hybrid method exploration, and increased data diversity and evaluation benchmarks. Such moves become very necessary in combating misinformation and ensuring the integrity of digital media as the world gets more AI driven.

REFERENCES

- [1] K. N. Ramadhani, R. Munir and N. P. Utama, "Improving Video Vision Transformer for Deepfake Video Detection Using Facial Landmark, Depthwise Separable Convolution and Self Attention," IEEE Access vol. 12, pp. 8932-8939, 2024, doi: 10.1109/ACCESS.2024.3352890.
- [2] Y. Patel et al., "An Improved Dense CNN Architecture for Deepfake Image Detection," IEEE Access, vol. 11, pp. 22081-22095, 2023, doi: 10.1109/ACCESS.2023.3251417
- [3] Y. -X. Luo and J. -L. Chen, "Dual Attention Network Approaches to Face Forgery Video Detection," IEEE Access, vol. 10, pp. 110754-110760, 2022, doi: 10.1109/ACCESS.2022.3215963.
- [4] A. H. Khalifa, N. A. Zaher, A. S. Abdallah and M. W. Fakhri, "Convolutional Neural Network Based on Diverse Gabor Filters for Deepfake Recognition," IEEE Access, vol. 10, pp. 22678-22686, 2022, doi: 10.1109/ACCESS.2022.3152029.
- [5] M. S. Rana, M. N. Nobi, B. Murali and A. H. Sung, "Deepfake Detection: A Systematic Literature Review," in IEEE Access, vol. 10, pp. 25494-25513, 2022, doi: 10.1109/ACCESS.2022.3154404..
- [6] H. Lin, W. Luo, K. Wei and M. Liu, "IMPROVED XCEPTION WITH DUAL ATTENTION MECHANISM AND FEATURE FUSION FOR FACE FORGERY DETECTION," 2022 4th International Conference on Data Intelligence and Security (ICDIS), Shenzhen, China, 2022, pp. 208-212, doi: 10.1109/ICDIS55630.2022.00039.
- [7] E. Kim and S. Cho, "Exposing Fake Faces Through Deep Neural Networks Combining Content and Trace Feature Extractors," IEEE Access, vol. 9, pp. 123493-123503, 2021, doi: 10.1109/ACCESS.2021.3110859.
- [8] Y. Patel et al., "Deepfake Generation and Detection: Case Study and Challenges," IEEE Access, vol. 11, pp. 143296-143323, 2023, doi: 10.1109/ACCESS.2023.3342107.
- [9] D. Park, H. Na and D. Choi, "Performance Comparison and Visualization of AI-Generated-Image Detection Methods," IEEE Access, vol. 12, pp. 62609-62627, 2024, doi: 10.1109/ACCESS.2024.3394250.
- [10] BIG DATA MINING AND ANALYTICS, ISSN: 2096-0654 21/25 pp889904, DOI: 10.26599/BDMA.2024.9020030 Volume 7, Number 3, September 2024.