ISSN: 2347-1697

Impact Factor::8.057



An Open-Access, Double-Blind, Peer-Reviewed, Refereed International Journal of Multidisciplinary Research

Unique Research Paper Identification (URPI)::2025(6)-IJIFR/V12/E10/008

International Journal of Informative & Futuristic Research (IJIFR) Volume - 12, Issue -10, June 2025

# **Truth vs. Fabrication: Exploring AI's Role in the Detection of Deepfakes**

### Dipti A. Mirkute

Assistant Professor, Department of Computer Engineering, Jawaharlal Darda Institute of Engineering and Technology, Yavatmal, India

Abstract: The emergence of deepfakes - artificially generated media formats that manipulate or alter imagery, videos, and audio using AI - has created new risks with regard to verifying digital content authenticity and cybersecurity. They are a type of forgery made with AI using sophisticated techniques and with the help of deep learning frameworks (especially Generative Adversarial Networks (GANs) or autoencoders) to create media that can be indistinguishable from real sources of information. While there have been reasonable uses for this technology in entertainment and accessibility, the amount of misuse in the form of disinformation campaigns, identity theft, blackmail, and political manipulation is staggering. This paper provides a review of the various methods used for deepfake generation including methods such as face swapping, attribute manipulation, face reenactment, lip-syncing, and audio cloning. Then, we discuss ways to detect deepfakes using AI, including different types and frameworks such as a feature-based, deep-learning based, and hybrid methods and comparing the strengths, weaknesses and applicable performance in real-world examples. We will provide further discussion on the emergent multimodal methods, which focus on detecting both audio and visual or metadata aspects of deepfake content, as well as our ethical, societal, and regulatory considerations, and discuss the emergence of tools used in the research community and conclusions on the different challenges found in the study including lack of data, adversarial changes and concerns over performance and generalizability in social contexts.

Keywords: Deepfake Detection, Artificial Intelligence, Generative Adversarial Networks (GANs), Autoencoders, Multimedia Forensics, Synthetic Media, Facial Manipulation, Lip-Syncing, Audio Deepfakes, Image Processing, Deep Learning, Hybrid Detection Models, Multimodal Analysis, Digital Misinformation, Cybersecurity, Ethical AI, Blockchain for Media Authentication

#### 1. Introduction

Deepfake generally refers to images, videos, and audio that are originally created or manipulated entirely by generative models via machine learning [4]. Primarily we associate the term with the manipulated images and videos which is hardly a new practice as there has always been an effort to manipulate visual media, particularly digitally, since it was introduced to us. Because it has been applied exponentially using technology to historically forge image and video for deception and entertainment [2], the decision to manipulate faces and voices really muddles with the understanding of the internet's integrity and information in our digital and content securities. Of

This work is published under Attribution-NonCommercial-ShareAlike 4.0 International License Copyright©IJIFR 2025



course, we could argue the downsides of this technology include the invasion of privacy and the dissemination of misinformation [1], calculating the negatives which deepfake technology is renowned for the unethical and malicious appropriation for economic, political, and social reputation. We are noticing lost ground in recent years with regards to facial forgery, especially when these exploits require virtually no technical skill to execute [8].

Deepfakes are created using artificial intelligence to create fictitious media to alter images, audio, and videos to simulate actual media creating something that appears to be true but is actually false. With the invention of user-friendly tools to help generate deepfakes, making deepfakes more accessible is troublesome for researchers that study the evolving landscape of generative AI and the detection of these technologies [6]. The deepfake technology uses deep learning technology, such as GANs, to represent a person's likeness and convincingly place them on another person's and often made it nearly impossible to distinguish them from real media[7]. These technologies are being used to spread false information, encourage hate, manipulate public perception, and even commit criminal acts such as blackmail and identity theft which can inflict harm on both the individual and society [5].

# 2. Deepfake Manipulation Types

Deepfake manipulation techniques are changing rapidly along with improvements in generative artificial intelligence and deepfake manipulation capabilities have developed into five primary types of deepfake manipulations. These manipulations can generate synthetic media that is not only fake, but also deliberately misleading. Many methods blur the lines between the two. There are five main subcategories of deepfake manipulations and developments that have emerged, all relying on improvements in deep learning and video synthesis techniques.

- A. *Face synthesis:* Face synthesis is the creation of completely synthetic human faces that do not exist in reality. The synthetic images are made by models trained on huge datasets of facial characteristics, which allows models with this training to create realistic, human-like faces. As mentioned, the images often appear indistinguishable from real images, and they can be used for fake accounts online or for creating online anonymity [6].
- *B. Attribute manipulation:* Attribute manipulation involves adjusting many aspects or regions of a face, including addition or subtraction of eyeglasses, skin tone, style, and amount of hair, and representation of age and/or gender. Only the portions of the face that relate to the altered identity of the original face must be edited, allowing the overall identity to remain recognizable while still representing a very altered version of themselves. Attribute manipulations produce very subtle artifacts in video content that can be powerful for establishing fake content that appears credible [6].
- C. Face swapping: Most people are already familiar with a process that is sometimes referred to as face swapping, or identity swapping, where the face of one person is swapped with the face of another person's face in a video or image. In this process, the expressive movement of the original person is preserved, but the visual identity is swapped for the target identity. The process of face swapping takes it a step further, as it integrates seamlessly and believably and allows it to look as though the target person is completing the actions captured in the original footage [6].
- D. Face reenactment: Face reenactment is the act of modifying the facial expressions of a person as presented in a video. Face reenactment allows changes in the facial behavior of a subject either in real time or post-processing, which ultimately includes making them smile when they didn't, make a frown when they didn't, etc. In the end, face reenactment will give a user a fully believable video of a subject behaving in a way that they completely did not [6].
- *E. Lip-syncing:* Lip-syncing is where you take the visual of a person's lips moving in a video and sync the combination of that physics with some outside audio typically speech that did not



happen. The resulting outcomes looks like that person's lips in the video are speaking aloud the dubbed audio content. Lip-syncing involves very precise mapping of facial motion, typically with complicated processes and and post-processing that make the resulting visuals look natural and believable [6].

# 3. Techniques Used To Generate Deepfakes

Deepfakes, which consist of written information, images, audio, and videos are likely the most prevalent type of fake media. The first "deepfake" video was produced in 2017 where one actor's face was replaced with the face of another actor. Shortly after, deepfakes garnered attention and began to go viral when a Reddit user called "Deepfake" demonstrated how the face of a celebrity could be morphed to provide them a starring role in an illicit video clip [2].

While deepfake video generation algorithms have recently seen a lot of development with advances in deep learning, computer vision and probabilistic modeling, altogether [1]. Deepfakes utilize deep learning methods, including CNNs, autoencoders, and GANs, to generate fake yet realistic images, audio or video that look like real people [5]. GANs have become the most popular method to produce deepfakes, which consist of a generator that makes synthetic media and a discriminator that distinguishes between real data and synthetic data, where each component improves realism through adversarial training [6].

Mostly, deepfakes use deep neural networks to manipulate image/video content. The most important 2 deep learning models, viz., autoencoders and GANs are the crux for the genesis of deepfakes [5].

- A. Autoencoders: Autoencoders are referred to as the first model can be utilized to create deepfake content. It was introduced in 2017 in a script in 2017, now it is known as FakeApp. Autoencoders have typically been used for dimensionality reduction, image compression, or generative learning model in machine learning tasks. In fact, autoencoders produced the best compressible images with the least loss function objective than other compressive image techniques. Generally, the autoencoder trains three main components: encoder, latent space, and decoder. The encoder is responsible for compressing the input image. The encoder is going to compress the input image or the image features while converting the input image features into the important features such as skin tone, skin texture, facial expression, structure of the face, state of eyes, and any other necessary features. This compressed data is given to the latent space, which is able to find the mapping of patterns and structural similarities among one or more data points. Finally, the decoder is left with to reconstruct an image, the image it has learned from, based on the latent space data. The decoder is tasked with making its image as realistic as the original input image based on the naturalistic image created in the latent space. These different learning processes can also be adapted for deepfakes [5]. Autoencoders can be employed as feature extractors to encode and decode facial features. When training, the preceding autoencoder learns to compress an input facial image onto a representation in a lower dimension that preserves the defining facial attributes. From this latent space representation, the original image can then be reconstructed [6].
- B. Generative Adversarial Network: GAN's contain pairs of neural networks, a generator and a discriminator; together they comprise a competitive process. The generator network generates synthetic images, which are then passed to the discriminator network together with real images. The generator network is learning to produce images that trick the discriminator, while at the same time, the discriminator network is being trained to determine whether images are real or synthetic. It is through this interactive training process that GANs are able to get better at generating images that



look progressively more realistic deepfakes [6]. A numbers of advanced generative adversarial network (GAN) models have propelled the emergence and sophistication of deepfakes. These approaches facilitate the manipulation of facial content and attributes in a way that is nearly indistinguishable from real content. One such model is AttGAN, which utilizes GANs to specifically manipulate facial attributes, and exemplifies attribute-awareness for high-quality editing of content including identity-preserving age progressive/regressive edits and face swaps. Using AttGAN, changes can be made to subtle and major facial attributes that maintain identity consistency. StyleGAN is one of the most powerful and detailed generative adversarial networks because it can produce remarkably detailed and realistic facial images, and manipulate any number of aspects of facial features for true-to-content deepfakes. Furthermore, STGAN can effectively manipulate facial attributes through a GAN architecture which supports even unlabeled data. STGAN is also favorable for deepfakes because it provides fine-grained editing while maintaining identity consistency [6]. The original architecture of StarGAN was enhanced in StarGANv2 to allow for a multi-domain version of image-to-image translation. While the original StarGAN only allowed one-to-one translations, StarGANv2 allows for one-to-many translations in a single model, such as translated images of faces in multiple styles or looks. The other model, CycleGAN, has a similar trade where style is transferred in some form from one domain to transfer some visuals from one domain. This style transfer is ideal for real-life face swapping by transferring some visuals from one person to another. For example, an interesting model is RSGAN, which encodes faces and hair visuals into the latent space such that latent manipulations can manipulate the internal encodings and the original images. From an audio-video synchronization perspective, LipGAN takes only audio input to create real-looking lip movements and synchronized speech for the purpose of producing realistic dubbed and/or synthetic video [6]. Both the architecture of the neural networks used and the dataset relied upon, help dictate the type of GAN models used. The most common software tools for video deepfake generation are Faceswap, DeepFaceLab, and DFaker; for audio generation WaveNet, MelNet, Char2Wave, and WaveGlow. Modern deepfake tools are relatively user friendly and this contributes to the prominence of deepfakes across various social media sites [5]. Mobile applications like ZAO, Auto FaceSwap, and FaceApp have simplified the process of creating and sharing deepfake images and video, accelerating the approximate spread of false content [9]. A very widely used deepfake approach is face swapping using GANs or Generative Adversarial Networks, which allow the user to easily transport an entire face from one image directly onto another face image and create outputs that can be indistinguishable from real world photography using the latest GAN models such as StyleGAN, StyleGAN2, and StyleGAN2-Ada [9].

C. Audio Deepfakes: Imitation-based deepfakes are a type of deepfake audio. Imitation-based deepfakes are created when the speaker's voice is transformed to sound like another speaker and is done for sometimes secondary privacy protection methods by using methods such as the Efficient Wavelet Mask (EWM) and through human voice imitation [3]. In synthetic deepfake audio, generated using a TTS or text-based speech system, the only available training data to train such models, such as Tacotron 2, Deep Voice 3, and FastSpeech 2. Important to note that, like many other machine learning techniques, training a TTS system requires structured audio and transcripts, and the model produces natural-sounding speech by using spectrograms and vocoders [3].

There are other advanced tools that also provide you with a way to create lip-sync aligned videos in which audio perfectly matches visual face gestures in the footage for aligned facial movement and content, such as Obamanet, which provides photorealistic video content with the output of a fake video being created from an image and then trained through GANs to leverage negative or face swapping resources or possibly a phone number, email address, social media handles, GPS



coordinates, and other personalized digital resources [5]. The misuse aspect of these tools has been significant when made available for consumers and further, many individuals and companies have suffered financially and reputational losses when utilized for fake news, derogatory propaganda, and defending against defamation. There may be a way for these tools to spur improvement or expand the possibilities of post-product Dubbing in realistic realistic dubbing during film production [5].

Besides the techniques listed above, there are also a number of open-source tools that are available to use digitally, making it easy for individuals to produce deep fakes, such as FaceApp, Reface, DeepBrain, DeepFaceLab, and Deepfakes Web.When these tools happened to be introduced they certainly prompted public fascination because of the ease in producing these seemly viable deep fakes [6].

# 4. AI Based Deepfake Detection Approaches

Deepfake video detection methods are typically classified into one of the three main approaches namely: Feature-Based Methods, Deep Learning-Based Methods, and Hybrid Methods. Feature-Based methods (hand-crafted) seek to find specific visual or audio anomalies in synthetic media that a human can perceive through abnormal blinking of subjects, light inconsistencies, or artifacts introduced into the digital synthetic medium. Deep Learning-based methods utilize model-driven approaches to automatically learn meaningful differences from large datasets with a mix of real and fake videos, and make distinctions with little human intervention. As a hybrid of both, Hybrid methods use both hand-crafted features and deep learning models together to use the advantage of hand-crafted features and deep learning together to create more accurate, effective, and robust systems [1].

Many detection models are built on controlled test conditions; however, compelling deepfake videos on the internet are often provided in low-quality visual or audio quality, which makes visual anomalies less perceptible and prevents reliable detection. This low-resolution audio-visual content, as a result, many deepfakes use multiple lower-resolution sources resulting in further difficulty for traditional tools to assimilate the nuanced real portions of the media separate from deepfake content. Recent studies have elevated the urgency of detecting deepfakes by showing clear limitations of detection systems to generalize across various types of deepfakes while also clarifying the performance level of detection data on degraded, compressed video that is currently disbursed on social media platforms [1].

### 5. AI Based Deepfake Detection Method

Deepfake Image vs. Deepfake Video DetectionMethods: Deepfake image detection is the process of determining whether or not a still image has been altered for the purpose of deceiving the viewer. Alterations include changing a person's appearance, adding or deleting objects, and changes in lighting or backgrounds. Detection methods include examining the metadata of an image, analyzing it for pixel-level inconsistencies, and comparing it to databases of known real and fake images. Since images are static, detection is focused solely on spatial features.Deepfake video detection presents a greater challenge because it requires processing more information, is temporal in nature, and takes into consideration a much larger volume of data. Changes may have been made to facial expressions, human movement, audio, or even deleting and/or adding objects. Detecting these changes requires analysis of not just a single frame but also the other frames in order to assess the temporal coherence of the frames. To affect temporal analysis, video detection techniques may utilize a machine learning algorithm to detect inconsistencies that are essentially cross-sectional in



nature in space and time. The time component associated with video makes video detection a greater challenge and far more computationally expensive than image deepfake detection. In short, while deepfake image detection primarily incorporates spatial analysis of a single image, deepfake video detection must contain spatial and temporal analyses in order to identify presence or absence of manipulation, and demands greater sophistication and resources [2].

## 6. AI Based Deepfake Detetection Techniques

Deepfake detection techniques primarily focus on identifying inconsistencies or artifacts introduced during the manipulation of images, videos and audios.

**A**. Deepfake Image/Video Detection Techniques: Filters and sensors used earlier detection methods were largely classic, that is, human-specified either handcrafted characteristics, or employing handcrafted rules. Now with many systems employing a deep neural network (DNN), and the use of other artificial intelligence (AI)-based multimedia generators the detection domain has matured and existing collection actions have been artificially strategically distributed which allows many DNN encoded multimedia detectors whenever sophisticated deepfake threat come up into conversations [5].DeepFake face images and video detection largely led research efforts on monitoring multimedia information and was generally intended to improve the specific confidentiality and integrity of multimedia content. Also, it is not an easy task to detect such recently altered multimedia content. The detection task has also become challenging due to the conditions of generative models. In a simplified view, the detection of forgery in multimedia content basically tests the original multimedia content to see if the generated multimedia has been tampered with or not. Forgery detection research was considered traditional before the emergence of DNN (AI-based)-based generated mediums. But lately DNN-generated multimedia detection is becoming more commonplace [9]. Deepfake detection methods can be divided into three categories: physical/physiological attribute-based, signal-level feature-based and data-driven deep learning models. Physical attributes identify visible inconsistencies such as unnatural blinking, mismatched lighting, or head pose inconsistencies. While physical attributes are the least computationally demanding methods, they are related to the evolution of deepfakes since they involve identifying features that are becoming increasingly incoherent. Signal-level features analyze pixel-level features and patterns resulting from the synthesis process via noise and non-noise [e.g., Convolutional Neural Networks (CNN), steganalysis, and Photo Response Non-Uniformity (PRNU)] and can enhance physical methods, thereby providing more robust modeling of deepfakes [5].Commonly used traditional CNN-based deepfake detection methods are largely dependent upon frequency-domain and statistical features. Networks like XceptionNet and ResNet are often used as backbones, where the models are applied to extract features and detect fake features from compression artifacts and device fingerprints. Capable of detecting intermediate fake features, however, modern approaches to deepfake creation have incorporated sophisticated post-processing techniques to make the fake product almost indistinguishable from the authentic. As a consequence, the feature distance between fake samples and genuine samples decrease, making it hard for binary classifiers like Support Vector Machines to do well as they concentrate the features onto specific areas of hyperplane. A downfall of using traditional CNNs for detection is their ability to detect only real frame-based samples, while causing temporal inconsistency, resulting in failures to generalize on unseen information, or samples of quality quality due to addition of new manipulation and risk overfitting in training. Alternatively, Capsule Networks were designed to address CNN restrictions in accurately represented spatial hierarchies and relationships of object division. Originally developed for inverse graphics problems, Capsule Networks attempt the best preservation of features orientation and spatial relationships, with the ability to correctly model



representations of 3D information, and using a smaller quantity of parameters and training data while performing comparably to CNNs makes them a feasible approach for deepfake detection tasks. [8]The most sophisticated detection techniques are data-driven deep learning models such as XceptionNet, MesoNet, and Capsule Networks, which learn highly complex patterns by using large datasets of both real and fake content. These advanced models performed particularly well on social media content (e.g., TikTok, Facebook) despite requiring a significant amount of computational power and the need for large datasets. Common feature extraction techniques employed in cases of poor quality include Scale Invariant Feature Transform (SIFT), Histograms of Oriented Gradients (HOG), Oriented FAST and Rotated BRIEF (ORB), and Image Quality Metrics (IQM) paired with classifiers such as Support Vector Machines (SVM). Taken together, combining methods, particularly deep learning, provide the scale and sophistication needed in the battle against deepfakes [5].

**B**. Deepfake Audio Detection Techniques: With regard to Deepfakes, many detection methods have been introduced to discern fake audio files from real speech. A number of ML and DL models have been developed that use different strategies [3]. With the tools that can generate fake audio becoming easier to access, audio deepfake (AD) detection continues to be an active area of research. Typically, AD detection can be categorized into two main classes of techniques: Machine Learning (ML) and Deep Learning (DL) [5]. Machine learning Techniques : Mostly MLbased AD detection includes building datasets using imitation techniques, then extracting entropy or statistical information from real and fake audio to create a feature-set that can be fed into traditional models, like Logistic Regression, that have been able to obtain success rates of up to 98% for detection [3].Deep learning Techniques : In contrast, DL methods especially using Convolutional Neural Networks (CNNs) take advantage of sophisticated methods. For instance, CNN-based models such as EfficientCNN and RES-EfficientCNN have been successfully employed for detecting synthetic audio, successfully classifying known types of artificial audio, achieving F1-scores of over 97.61 on benchmark datasets. Other DL methods use an audio signal distributed into a 2D representation, such as spectrograms and histograms, allowing the CNN architectures to analyze and classify for live detection. Both DL and ML models have been shown to achieve detection accuracies of up to 98.5%, although their performance is often limited by data transformation and scalability [3]. Overall comparisons on DL vs ML methods demonstrate that ML methods like Support Vector Machine (SVM) can attain very high accuracy (as high as 99%) whereas CNNs can specifically learn to extract and generalize small deviations in audio like no other approach. The major downside to CNN based models is that they cannot directly input raw audio because audio data must be converted into a visual representation. Even so, CNNs are valuable because they automatically extract features and create spuriously correlated synthetic audio patterns [3].

# 7. Multimedia-Enabled Deepfake Detection

Multimedia-enabled deepfake detection uses audio, video, and image data to identify synthetic and manipulated media faster and more accurately. Traditional detection approaches investigate a single modality either facial images or audio separately. In contrast, multimedia-enabled deepfake detection methods leverage multimodal deepfake detection, which enables the detection of inconsistencies present within various modalities of responses and cues, with the assumption that all of them would be in alignment. In other words, compared to a traditional approach that would factor each video component and audio component individually, beyond multimodal synchrony, a researcher or innocent bystanders can consider dissociation methodologies across integrated video content—like examining the specific lip movements that match the audio track and facially displayed emotional expression that accompanies the speech sample. Dispositions that creep into



the observed finding could present lip-syncing issues, unnatural facial motions or expressions, or inconsistent emotions aligned with speech, all of which could highlight signs of deepfake manipulations without the researcher being aware when examining each discrete assessment on its own [7].In terms of detection, audio spoofing also looks at the authenticity of words to identify voice impersonation and synthesis through text-to-speech. VMD addresses inconsistencies in visual experiences in both audio and video, or still pictures, like changes in texture, motion, lighting, and facial asymmetry. When both visual and sound are measured together, the detection system is stronger and more credible. In this sense detection approaches are useful for purposes of detecting deepfake materials as a means of misinformation and impersonation, with realistic generated impersonation through the use of Generative Adversarial Networks (GANs) and neural networks that continue to widen the positive gap for legitimate uses of AI through machine learning. On the backend of the detection of all of this unprecedented experimental use of media, there are systems based on blockchain-based protocols and cryptography as a means of enhancing traceability and verifying media authenticity and integrity as both the private and public digital exchanges of media expand. For instance, Blockchain Distributed Ledger Technology (BDLT) would be an ideal storage mechanism for hashes of original media data that could be used to identify and make analyses of original media as it became compromised, by comparing it with existing hashes that could be collected from newly made original media embedded in time-stamped media secured with hash identification, behaviours during creation and longitudinal metadata records from the recording devices. Watermarking media adds the additional opportunity to add invisible identifiers to verify media, along with metadata such as timestamps, device identifiers, licencing and registration numbers and statistics to identify background information and concensus during forensic investigations and consistency. Forensic investigation largely uses deep learning in the form of either Convolutional Neural Networks or Transformer-based models systems trained to capture details in vast datasets that are invisible to humans, while recognising subtle and almost invisible signals of visual and sound manipulation to aid in the responsible and credible display of media in these unprecedented times [2],[7]. The combination of these different ways not only improves deepfake detection, but also enables systems to operate in real-time situations, such as during live streaming or using a social media platform. New lightweight models and methods, such as federated learning, explainable AI and self-supervised learning, are emerging to combat challenges like adversarial attacks and generalization to new datasets. More broadly, multimediaenabled detection represents an inclusive and future-oriented approach to protecting the authenticity of digital content in a world where synthetic media is complex and widespread [7].

### 8. Tools For Deepfake Detection

There are many real-time tools and platforms for deepfake detection with multimedia-enabled environments. Sensity AI is an example of a sophisticated detection platform that scans image and video metadata to discover whether a synthetic content exists. It employs machine learning models to deepfakes, but also to prevent deepfakes in the first place. Lip Sync AI is another example, which can swap speech from one video synchronizing it to mouth movements. Although, it exists for creating fake videos, it can clearly take advantage of the technology it employs to identify and understand patterns of manipulation. There is also a variety of other open-sourced tools such as Faceswap, DeepFaceLab and others, that allow for utilizing many neural networks and configurations to train and popularize deepfake models; all of which can be used for evidence-based forensic analysis and detecting manipulated media. [7].



# 9. Ethical And Societal Implications

- A. Disinformation & Propaganda: Deepfakes are also commonly used to spread fake news, alter political narratives, and defame individuals, all of which are catalysts to civil strife and sometimes violence [5].
- *B. Psychological & Emotional Harm:* Individuals who are subject to deepfake forgeries, or character assassinations, could potentially sustain humiliation, mental anguish, and negative social consequences [5].
- *C. Erosion of Trust*: The emergence of hyper-realistic deepfakes have rendered it difficult to trust visual evidence, which breaks down trust by the public in media, or the systems of justice, or digital content authenticity [11].
- D. Privacy intrusions & Identity Theft: Risks of privacy infringements, blackmail, and assumed identities could result from facial switching and voice duplicating [5].
- *E. Legal &Regulation:* Deepfakes are advancing faster than laws can regulate, especially around consent, liability, and admissibility of evidence [11].

## **10. Challenges And Limitations**

- *A.* Generalization Issues. Many detection models falter in real-world situations, in particular when tested on unseen datasets, or compressed files (e.g., WhatsApp, Facebook) [4].
- *B. Adversarial Arms Race :* As detection models develop, so too do deepfake generation models to circumvent detection, establishing a never-ending arms race between defender and attacker [2]. Attackers can leverage types of perturbations and laundering approaches to trick detection systems, such as compression, noise, or resized [10].
- *C. Data Quality &Availability :* There is a shortage of reliable, balanced, and diverse datasets for which to extract training data that ideally improves accuracy and robustness [2].
- D. Computational Complexity :Deep learning-based detection methods have high computational requirements that preclude feasible deployment in real-time scenarios and low-resource environments [2].
- *E. Overconfidence in Detection Accuracy* :Some projects include reported accuracy rates with overconfidence and do not translate readily to practical applications due to dataset bias and lack of variation between environments [2].
- *F. Multimodal Detection Complexity :* Established detection across image, video, and audio modalities continues to be complex, and detection across modalities remains under researched [11].
- *G. Unreliable Real-World Performance:* Detection models tend to fail on compressed or edited content viewed on social media, in which the forensic signals are destroyed [10].
- *H. Data Drift and Catastrophic Forgetting:* Models can learn more types of deepfakes and variants, but they can only do this at the cost of forgetting what they learned previously, necessitating in-depth continual learning [10].
- *I. Lack of Interpretability:* Most deep learning detectors are black boxes, making them not only hard to trust, but hard to use in critical areas such as the law and forensics [10].
- J. Quickly Evolving Threat Landscape: Due to novel generative processes and changing social media (e.g., processing) requirements, the threat could change rapidly. Ongoing continual learning or few-shot learning is needed to update models continuously [10].

### **11. Future Directions**

Future opportunities in deepfake detection point to a clear need for reliable and generalizable models capable of providing accurate and reliable distinctions outside the constraints of an experimental setting, especially when considering common post-processing distortions such as JPEG or re-encoded versions of the original on social media. Presently, deepfake detection models



often maintain task-specific robustness to some primary form of distortion, but poorly generalize to any usages beyond the specific dataset used in their training, a barrier which has resulted in a shift toward creating diagnostic detection algorithms capable of "in the wild detection" [4]. As recognition that deepfake techniques and manipulation abilities also evolve, an industry-wide interest in multimodal strategies has arisen, where visual cues, audio cues and textual cues are merged into the detection algorithm to provide a more robust and well-balanced detection, adapting across multiple forms and types of manipulative techniques [5], [11]. The expansion of preferred and formed datasets is, therefore, a very needed improvement; along with the advancement of standardized, large-scale datasets, which will drive what datasets drive learning and performance assessment; and providing models the same context during both training and testing, decreasing the likelihood of overfitting to limited content types [2], [11]. Another avenue that is promising can include the use of Explainable AI (XAI) techniques and methods, to promote users trust with existing seeding systems - meaning the user understands the rationale behind detection results [11]. Last and now part of more discussion is emerging development of blockchain or watermarking possibilities for future options for proactive loss prevention and benefits that accompany an integrity and reliability of authenticity at the document's beginning or at dissemination [2], [5]. From an ethical, policy, and regulatory perspective, the future work will also need to focus on establishing comprehensive and strong regulatory frameworks, interdisciplinary collaboration and public education to enable an effective counter to the misuse of deepfake technologies, while enabling their responsible use in valuable areas such as education, health care, and entertainment [2], [11].

## 12. Conclusion

The mass production of deepfake technologies is putting more and more strain on the validity of everything in digital media, cybersecurity, and trust in the public sphere. By leveraging constructs like Generative Adversarial Networks (GANs), and Linear and Non-Linear autoencoders for generative models, deepfake technologies can produce synthetic content that betrays authentic content in highly-credible ways. Further, it's making it near impossible for traditional forms of detection to ultimately become effective. While there have been detection frameworks developed for deepfake detection—featuring both hand-crafted features (modeling) and new deep learning detection models (the XceptionNet, and Capsule Networks, 2019)—these frameworks struggle to achieve generalized success across the range of datasets involved in detecting deepfake content, also suffering from post-processing and compression artifacts. The paper discusses the continued promise of multimodal detection that exploits audio, visual, and metadata together, as well as the potential benefits of blockchain and digital watermarking leveling methods to support content verification and content tracing. However, with deepfake generation methods continuing to evolve until an understanding of misrepresentation or falsification is well-defined, we need adaptive, explainable AI working in real-time and in real-life environments. The study concluded with the need for international regulations, ethical practices and practices for the public to understand the dangers reducing harm to society's misuses. It called for inter-discipline collaboration, ongoing models to learn from experience, and new laws and regulations to protect digital content. All of these components are necessary to develop a clear defendable posture against the increasing threat of deepfakes while being able to apply their advantages.

### 13. References

- [1] M. Alrashoud, "Deepfake video detection methods, approaches, and challenges," *Alexandria Engineering Journal*, vol. 125, pp. 265–277, 2025
- [2] A. Kaur, A. N. Hoshyar, V. Saikrishna, S. Firmin, and F. Xia, "Deepfake video detection: challenges and opportunities," *Artificial Intelligence Review*, vol. 57, no. 6, p. 159, 2024.



- [3] Z. Almutairi and H. Elgibreen, "A review of modern audio deepfake detection methods: challenges and future directions," *Algorithms*, vol. 15, no. 5, p. 155, 2022.
- [4] L. Guarnera, O. Giudice, F. Guarnera, A. Ortis, G. Puglisi, A. Paratore, L. M. Q. Bui *et al.*, "The face deepfake detection challenge," *Journal of Imaging*, vol. 8, no. 10, p. 263, 2022.
- [5] Y. Patel, S. Tanwar, R. Gupta, P. Bhattacharya, I. E. Davidson, R. Nyameko, S. Aluvala, and V. Vimal, "Deepfake generation and detection: Case study and challenges," *IEEE Access*, vol. 11, pp. 143296–143323, 2023.
- [6] A. Naitali, M. Ridouani, F. Salahdine, and N. Kaabouch, "Deepfake attacks: Generation, detection, datasets, challenges, and research directions," *Computers*, vol. 12, no. 10, p. 216, 2023.
- [7] A. A. Khan, A. A. Laghari, S. A. Inam, S. Ullah, M. Shahzad, and D. Syed, "A survey on multimedia-enabled deepfake detection: state-of-the-art tools and techniques, emerging trends, current challenges & limitations, and future directions," *Discover Computing*, vol. 28, no. 1, p. 48, 2025.
- [8] L. Y. Gong and X. J. Li, "A contemporary survey on deepfake detection: Datasets, algorithms, and challenges," *Electronics*, vol. 13, no. 3, p. 585, 2024.
- [9] A. Malik, M. Kuribayashi, S. M. Abdullahi, and A. N. Khan, "DeepFake detection for human face images and videos: A survey," *IEEE Access*, vol. 10, pp. 18757–18775, 2022.
- [10] I. Amerini, M. Barni, S. Battiato, P. Bestagini, G. Boato, V. Bruni, R. Caldelli*et al.*, "Deepfake media forensics: Status and future challenges," *Journal of Imaging*, vol. 11, no. 3, p. 73, 2025.
- [11] N. U. R. Ahmed, A. Badshah, H. Adeel, A. Tajammul, A. Duad, and T. Alsahfi, "Visual Deepfake Detection: Review of Techniques, Tools, Limitations, and Future Prospects," *IEEE Access*, 2024.



Dipti A. Mirkute is presently working as Assistant Professor in Department of Computer Engineering, Jawaharlal Darda Institute of Engineering and Technology, Yavatmal, India. Her area of inerests includes Artificial Intelligence & Deepfake techniques. She has attended & published various papers in International Journals & conferences throughout her carrier.

PAPER CITATION: Mirkute, A.D. (June 24, 2025) :: Truth vs. Fabrication: Exploring AI's Role in the Detection of Deepfakes. International Journal of Informative & Futuristic Research (IJIFR), Volume - 12, Issue -10, June 2025, Pg. No. 40-50, URPI-2025(6) IJIFR/V12/E10/008 Available Online through:: https://ijifr.org/pdfsave/24-06-2025278IJIFR-V12-E10-008.pdf