International Journal of Scientific Research in Engineering and Management (IJSREM)



Volume: 07 Issue: 08 | August - 2023

SJIF Rating: 8.176

ISSN: 2582-3930

UNMASKING REALITY: A COMPREHENSIVE REVIEW OF DEEPFAKE TECHNOLOGY

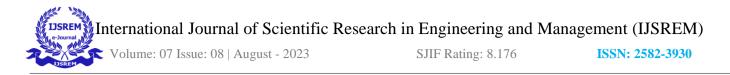
¹Yogendra Awasthi, ²Ambresh Mishra, ³Rakesh Kumar ^{1,2,}Student, ³Assistant Professor ^{1,2,3}Department of Computer Application, ^{1,2,3}Tulas Institute, Dehradun, India, ^{1,2,3}Sri Dev Suman Uttrakhand University, Tehri, India

Abstract: This review paper delves into the realm of deep fake technology, unravelling its intricate landscape and multifaceted implications. The paper navigates through the foundational techniques that underpin deep fake creation, elucidating their intricacies and potentialities. With a spotlight on diverse applications across domains, the paper probes into the profound societal impact of this technology. Amidst these explorations, the paper confronts the myriad challenges posed by deep fakes, ranging from technical hurdles to ethical dilemmas. It unveils the cutting-edge strategies employed for detecting and mitigating deep fake content, underscoring the continuous evolution in this dynamic field. The paper delves into the ethical and legal dimensions surrounding deep fake technology, addressing matters of consent, privacy, and accountability. Furthermore, the paper casts a forward-looking gaze, envisioning the potential trajectories that deep fake technology may embark upon, and the corresponding strategies to harness its benefits while minimizing risks. This comprehensive review paper thus seeks to provide a panoramic view of deep fake technology's present state and its future role in shaping the intricate tapestry of media content creation.

Keywords: Deep fake technology, Foundational techniques, Applications, Societal impact, Challenges, Technical hurdles, Ethical dilemmas, Detection and mitigation, Strategies, Ethical and legal dimensions, Consent, Privacy, Accountability, Future trajectories, Benefits, Risks, Media content creation.

I. INTRODUCTION

In the contemporary landscape of media content creation, the advent of deepfake technology has ushered in a paradigm shift. Driven by remarkable advancements in artificial intelligence and machine learning, deepfakes embody a fusion of innovation and inherent risks. This comprehensive review paper embarks on a comprehensive journey through the intricate landscape surrounding deepfake technology. By delving into its foundational techniques, analyzing its diverse applications across multifarious domains, addressing the spectrum of challenges it poses, scrutinizing the strategies harnessed for its detection and mitigation, and delving into the profound societal implications it engenders, this paper aspires to provide a holistic grasp of the present state of deepfake technology and its potential trajectories in the dynamic future.



As the exponential growth in deepfake technology captivates the interest of researchers, technologists, and the wider populace, this review paper takes on the role of an insightful guide. The portmanteau "deepfake," a convergence of "deep learning" and "fake," encompasses the creation of synthetic media content through intricate machine learning algorithms. These algorithms, imbued with the remarkable ability to craft lifelike depictions of individuals engaged in actions they never actually performed, hold within them both awe-inspiring potential and profound ethical dilemmas.

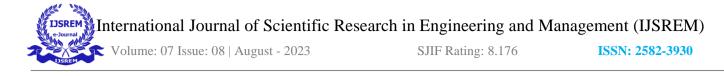
At the heart of deepfake creation lies the application of generative models, chiefly exemplified by Generative Adversarial Networks (GANs). These networks comprise a generator responsible for content creation and a discriminator tasked with evaluating content authenticity. Through iterative training, GANs produce content progressively indistinguishable from genuine media, with variations like Conditional GANs (cGANs) and Variational Autoencoder-GANs (VAE-GANs) enhancing content quality even further.

Deepfake technology spans a vast spectrum of applications, from the entertainment industry, where actors can be seamlessly integrated into historical footage, to more malicious domains like political disinformation and privacy breaches. The complexity of deepfakes is not confined to their technical facets; ethical considerations loom large due to unauthorized likeness utilization, raising questions about consent and privacy. An ongoing arms race between deepfake creators and detection techniques underscores the rapidly evolving landscape.

Detecting deepfake content assumes critical significance, driving researchers to innovate in methodologies like facial landmarks analysis, expression inconsistencies scrutiny, and forensic video attribute assessment. Yet, the escalating sophistication of deepfake generation techniques mandates an ongoing evolution of detection strategies.

The societal implications are profound as deepfakes blur the boundaries between reality and fabrication, amplifying the potential for misinformation campaigns and destabilization of trust in visual evidence. As this review paper aims to traverse the intricate dimensions of deepfake technology, from its technical intricacies to its profound societal implications, it aspires to contribute to a holistic understanding of its current state and the crucial measures necessary for a balanced harnessing of its potential while safeguarding against its inherent risks.

This review paper aims to unravel the intricacies of deepfake technology and probe its far-reaching consequences. By delving into its core techniques, traversing its applications across domains, grappling with the ethical quandaries it presents, and revealing strategies for detection and mitigation, we gain insight into a multifaceted landscape reshaping our understanding of media authenticity. As we navigate this complex terrain, a comprehensive grasp of deepfake technology emerges as essential. By exploring its nuances and implications, we empower ourselves to navigate a world where visual content takes on new dimensions. The evolution of deepfake technology necessitates a cautious yet innovative approach—one that harnesses its creative possibilities while guarding against misuse. Let us embark on this journey,



equipped with insights and awareness, as we navigate the crossroads of technological advancement, artistic expression, and ethical responsibility.

II.TECHNIQUES AND APPROACHES

Deepfake creation is an intricate art that relies on the sophisticated utilization of generative models, with a specific emphasis on the ground breaking Generative Adversarial Networks (GANs). Within the landscape of machine learning, GANs stand out as a pivotal breakthrough, particularly in the realm of generative tasks. At the core of GANs lies a dual-network architecture that consists of a generator network and a discriminator network, each playing a distinct yet interconnected role.

Generator Network:

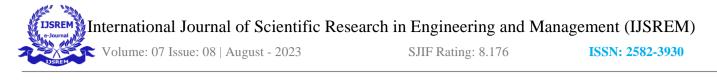
The generator network, true to its name, serves as the creative force behind deepfake content. Its primary objective is to generate media, such as images or videos that replicate the intricate characteristics of authentic content. In essence, the generator endeavours to produce content that is virtually indistinguishable from actual media. This network operates through a complex set of learned parameters, which are derived from an extensive training dataset. In its early stages, the content generated by the network might be far from convincing; however, as the network undergoes a process of iterative training, it gradually refines its capabilities, working toward the creation of increasingly convincing imitations.

Discriminator Network:

In contrast to the generator, the discriminator network takes on the role of an evaluator or judge. Its primary responsibility is to assess the authenticity of the content generated by the generator. The discriminator network is trained using a combination of authentic and generated content, enabling it to discern the differences between the two. The fundamental objective of the discriminator is to become exceedingly proficient at distinguishing between genuine and generated media. Over time, as the generator improves its ability to produce content that appears genuine, the discriminator's task becomes more intricate and challenging.

Iterative Training Process :

At the heart of GANs lies a captivating iterative training process that involves a continuous interaction between the generator and discriminator networks. Throughout the training phase, the generator's central aim is to produce content that deceives the discriminator into mislabeling it as authentic. Conversely, the discriminator strives to accurately differentiate between genuine and generated content. This constant interplay and competition between the two networks lead to a continuous refinement of their respective capabilities.



Quality Enhancement through Iteration :

The most intriguing aspect of GANs comes to light during the iterative training process. As this process unfolds, a remarkable phenomenon occurs: the content generated by the generator evolves, gradually becoming more authentic and increasingly difficult to distinguish from genuine media. This evolution is driven by the dynamic interplay between the generator and discriminator networks, where each network strives to outperform the other. This iterative training process, though computationally intensive, results in the production of content that exhibits a striking resemblance to real media in terms of its characteristics and details.

Variations: cGANs and VAE-GANs :

The landscape of GANs has blossomed to encompass an array of variations, each tailored to address specific challenges and optimize performance. Among these variations are Conditional GANs (cGANs), which introduce conditions or additional information to guide the content generation process. This allows for the specification of particular attributes, facilitating more targeted and controlled content generation.

Variational Autoencoder-GANs (VAE-GANs) represent an intriguing fusion of the concepts of Variational Autoencoders (VAEs) and GANs. VAEs are a distinct class of generative models that focus on learning latent representations of data. When integrated with GANs, VAE-GANs aim to harness the strengths of both approaches, aspiring to generate content of higher quality and greater diversity.

Enhanced Content Quality :

The introduction of variations such as cGANs and VAE-GANs has significantly contributed to the advancement of deepfake technology by elevating the quality of the content generated. These variations add layers of complexity and sophistication to the underlying generative process. This, in turn, enables more precise control over the characteristics of the content being produced. Consequently, the content generated becomes even more challenging to distinguish from authentic media due to the heightened attention to detail and fidelity.

In essence, the techniques and approaches employed in deepfake creation, predominantly revolving around GANs and their innovative variations, orchestrate a dynamic interaction between the generator and discriminator networks. This intricate dance of competition and cooperation unfolds through iterative training, culminating in the production of content that progressively attains a level of fidelity that closely mirrors genuine media. The integration of these variations amplifies the sophistication of the generative process, resulting in content that achieves an exceptional degree of authenticity and quality.



Volume: 07 Issue: 08 | August - 2023

SJIF Rating: 8.176

ISSN: 2582-3930

III.APPLICATIONS AND IMPACT

The applications and impact of deepfake technology are expansive and far-reaching. While they introduce innovative possibilities within entertainment and narrative creation, they also raise pivotal ethical, societal, and security issues. Deepfakes wield the power to reshape narratives, influence public perceptions, and potentially erode trust in digital media. Balancing the creative potential of deepfakes with the ethical challenges they pose is an essential undertaking as society grapples with the multifaceted dimensions of this technological advancement.

Entertainment Innovation :

In the domain of entertainment, deepfakes have heralded a paradigm shift. One remarkable application is the seamless integration of actors into historical footage, redefining the way stories are told. This innovation enables the revival of historical figures on screen, reimagining their interactions and narratives. By seamlessly blending authentic historical content with meticulously crafted deepfake elements, storytellers can transcend temporal limitations and provide fresh perspectives on historical events.

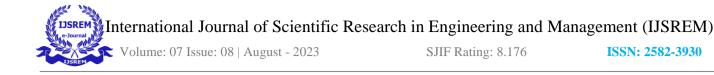
Furthermore, the technique of dubbing films into different languages through precise lip-sync alterations has garnered significant attention. Deepfake technology streamlines the complex process of modifying lip-sync, resulting in more authentic dubbing experiences. This application not only expedites the localization of content but also enhances viewers' engagement by minimizing the disruptions caused by incongruent audio and visuals.

Expanding the Cinematic Domain :

Deepfake technology has expanded the boundaries of the cinematic realm. Filmmakers can now seamlessly integrate actors into scenes that were once deemed logistically or historically unfeasible. This newfound capability empowers the creation of narratives that transcend conventional limitations, offering fresh perspectives on familiar stories and characters. The potential to merge contemporary actors with historical personas or to resurrect long-lost actors onto the screen adds a novel layer to storytelling.

Unveiling Ethical Dilemmas :

The applications of deepfake technology extend far beyond entertainment, giving rise to profound ethical and moral dilemmas that challenge the essence of truth and authenticity. Deepfakes have been leveraged for nefarious purposes, including political disinformation campaigns and orchestrating cybercrimes. The ease with which individuals' appearances can be manipulated to propagate fabricated narratives has raised concerns about the erosion of trust in digital media and its potential ramifications for society.



Political Disinformation and Cybercrimes :

Deepfakes have the potential to disrupt the political landscape in unprecedented ways. Ill-intentioned actors can create deepfake videos depicting public figures uttering statements they never actually made. This creates opportunities for disseminating false information, swaying public opinion, and sowing discord. The striking realism of deepfake content blurs the boundaries between truth and falsehood, making it challenging for viewers to distinguish genuine information from manipulated content.

Furthermore, the tools within deepfake technology can be weaponized for cybercrimes. Voice cloning, a subset of deepfake technology, enables the convincing mimicry of individuals' voices. This could enable scammers to impersonate individuals, coercing victims into revealing sensitive information or engaging in financial transactions under false pretenses.

Privacy Intrusions :

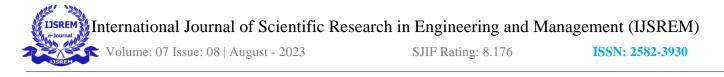
Beyond the realms of politics and criminality, deepfakes raise substantial concerns about privacy. The ease with which facial features and voices can be replicated poses a risk to individuals' personal privacy. The creation of convincingly realistic deepfake videos could lead to scenarios where people are portrayed engaging in activities they never actually participated in, resulting in reputational harm and personal distress.

IV.CHALLENGES AND LIMITATIONS

The evolution of deepfake technology is a journey fraught with a diverse range of challenges, stemming from both technical intricacies and ethical considerations. Addressing these challenges is vital to ensure the responsible development and deployment of deepfake technology. Striving for equitable datasets, minimizing algorithmic biases, establishing ethical guidelines, and advancing detection techniques are all essential components in the quest to harness the potential of deepfakes while mitigating their risks.

Acquisition of Comprehensive Datasets :

One of the foundational challenges in deepfake technology is the acquisition of extensive and diverse datasets for training machine learning models. These datasets are crucial for training algorithms to accurately replicate facial expressions, mannerisms, and other characteristics necessary for creating convincing deepfake content. However, obtaining such datasets that encompass a wide range of individuals, poses, lighting conditions, and expressions can be arduous and time-consuming.



Algorithmic Biases and Fairness :

The datasets used for training deepfake algorithms might inadvertently contain biases, whether in terms of gender, ethnicity, or other attributes. This can lead to algorithmic biases in the generated content, perpetuating societal imbalances and reinforcing stereotypes. Ensuring equity and fairness in both the training data and the resulting content is a complex challenge that requires meticulous curation of datasets and the development of algorithms that are sensitive to potential biases.

Ethical Dimensions: Unauthorized Utilization :

The ethical dimensions of deepfake technology are multifaceted, with concerns ranging from unauthorized utilization of individuals' likenesses to questions about consent and privacy. The creation of deepfake content without the explicit consent of the individuals being replicated raises serious ethical questions. The unauthorized use of someone's appearance, voice, or identity can lead to reputational harm and invasion of privacy. Addressing these ethical concerns requires establishing clear guidelines for consent and the responsible creation and use of deepfake content.

The Ongoing Arms Race: Creators vs. Detectors :

The landscape of deepfake technology is marked by an ongoing arms race between creators who develop sophisticated techniques for generating realistic content and researchers focused on devising effective methods for detecting deepfakes. As detection techniques advance, creators adapt and refine their methods to evade detection. This dynamic cycle of innovation on both sides makes it challenging to achieve a definitive solution for either generating or identifying deepfake content.

Innovation vs. Detection: A Delicate Balance :

With each stride forward in deepfake generation techniques, there is a corresponding need for the development of more advanced detection mechanisms. The balance between innovation and detection is delicate. While technological progress in deepfake creation opens new avenues for creative expression, it also necessitates parallel advancements in detection to ensure the responsible use of the technology and mitigate its potential misuse.

Strategic Implications: Trust in Digital Media :

The proliferation of convincing deepfake content has the potential to undermine trust in digital media. As the boundary between authentic and manipulated content blurs, individuals may become more skeptical of the information presented in various media formats. This has far-reaching implications for journalism, public discourse, and even legal proceedings where visual evidence is pivotal. The erosion of trust challenges the very foundations of communication and truth in the digital age.



Volume: 07 Issue: 08 | August - 2023

SJIF Rating: 8.176

ISSN: 2582-3930

V.DETECTION AND MITIGATION

The fast spread of deepfake technology has made it more important than ever to develop reliable detection and mitigation techniques. Deepfakes have the power to spread false information on a large scale, endangering the veracity of visual media and undermining public confidence. Researchers, developers, and specialists from a wide range of fields are always battling to come up with novel solutions that can expose deepfake content and lessen its potential impact in order to address these difficulties. This part goes deeply into the complex world of detection tactics, highlighting the dynamic efforts made to keep up with the deepfake technology's advancing level of sophistication.

Diverse Arsenal of Unmasking Approaches :

Detecting deepfake content demands an assortment of approaches, each honed to unveil the subtle clues that betray the presence of manipulation. One avenue is the scrutiny of facial landmarks and expressions, targeting the irregularities that arise when deepfake algorithms struggle to replicate the intricate nuances of human emotion. These imperfections, though challenging to perceive, become telltale signs that experienced eyes can detect.

Another avenue centers on temporal inconsistencies within videos. Human facial expressions and movements are inherently dynamic, posing challenges for deepfake algorithms attempting precise replication. As a result, anomalies like unnatural blinking patterns, odd eye movements, and other temporal inconsistencies become evident upon close examination of deepfake content.

Forensic Analysis for Authenticity Assessment :

Forensic analysis, akin to investigative work, holds a pivotal role in deepfake detection. This involves meticulous examination of video attributes that bear traces of manipulation. Compression artifacts, noise patterns, lighting disparities, and other subtle discrepancies can collectively yield vital clues about the authenticity of the content. By juxtaposing these attributes against expected norms in genuine media, forensic analysis unveils potential signs of tampering.

The Unrelenting Arms Race: Evolution of Detection and Deception :

As deepfake generation techniques progress, so do detection methods. The creators of deepfakes are acutely aware of detection efforts and often engineer strategies to outsmart them. They deliberately introduce imperfections or noise to make their content mirror authentic media more closely, rendering detection an intricate challenge. This has engendered a ceaseless arms race—an incessant struggle where each stride in detection advancement compels creators to innovate their deception techniques.



Volume: 07 Issue: 08 | August - 2023

SJIF Rating: 8.176

ISSN: 2582-3930

Resilience and Adaptability in Detection Systems :

Effective detection mandates the deployment of resilient and adaptive systems capable of identifying emerging manipulation techniques. At the vanguard of this battle are machine learning algorithms. By undergoing training on vast datasets encompassing both genuine and manipulated content, these algorithms learn to discern the subtle patterns indicative of deepfake manipulation. The ability of these systems to adapt is paramount, as they must perpetually evolve to outpace the ceaselessly shifting landscape of deepfake methodologies.

Interdisciplinary Collaboration: Elevating Detection Efficacy :

Addressing the multifaceted challenge of deepfake detection necessitates interdisciplinary collaboration. Experts hailing from domains like artificial intelligence, computer vision, and media analysis must converge their insights to synthesize diverse detection methodologies. This collaborative synergy engenders more fortified systems equipped to confront a wider array of manipulation techniques. The amalgamation of technical prowess and domain expertise is pivotal in this collaborative pursuit.

Harnessing Deep Learning for Detection Precision :

Deep learning techniques, particularly convolutional neural networks (CNNs) and recurrent neural networks (RNNs), have exhibited tremendous potential in deepfake detection. These models excel in discerning intricate patterns that elude traditional algorithms. By undergoing training on expansive datasets encompassing both authentic and manipulated content, deep learning models can pinpoint the subtle anomalies introduced by deepfake algorithms.

A Perpetual Quest for Robustness :

The expedition towards robust deepfake detection is far from static. Researchers persistently refine extant methodologies and pioneer novel ones to keep stride with the ceaselessly evolving capacities of deepfake technology. This endeavor encompasses collaborations across academia, industry, and governmental spheres, underscoring the urgency of effectively countering deceptive media manipulation.

VI.ETHICAL AND LEGAL CONSIDERATIONS

The ethical and legal considerations linked to deepfake technology encompass a constellation of concerns that extend across consent, privacy, and reputation. The unauthorized manipulation of individuals' likenesses invokes moral deliberations regarding agency and autonomy. Evolving legal frameworks are essential to bridge the gap between the power of technology and the rights of individuals. Striking a balance between technological advancement and ethical considerations requires continuous discourse,



collaboration, and the commitment to safeguarding human dignity and societal values in an ever-evolving digital landscape.

Navigating the Ethical Landscape :

Deepfake technology presents a moral landscape that requires careful navigation. Central to this landscape is the act of using someone's likeness-be it their face, voice, or persona-without their explicit authorization. This practice raises profound questions about autonomy, agency, and individual sovereignty. The fundamental right of individuals to control how their identity is portrayed clashes with the ability of technology to manipulate it in ways that can deceive or harm.

The Consent Conundrum :

At the heart of ethical concerns lies the issue of consent. When deepfake technology is employed to create content that involves real individuals, the question of whether they consented to such use becomes pivotal. The absence of explicit consent challenges the notion of informed agency and raises red flags about the appropriation of one's identity. Ethical considerations underscore the importance of ensuring that individuals have a say in how their likeness is utilized.

Reputation in the Balance :

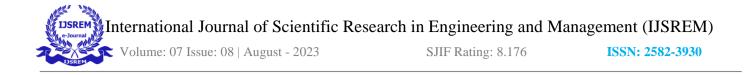
Another ethical dimension centers on the potential harm to an individual's reputation. Deepfakes have the power to fabricate scenarios that never occurred, portraying someone in a manner that is entirely divorced from reality. This has the potential to cause significant reputational damage, as the fabricated content might be perceived as genuine. The erosion of reputation can have far-reaching consequences, impacting personal, professional, and societal aspects of an individual's life.

The Role of Evolving Legal Frameworks :

The challenges posed by deepfake technology have prompted the evolution of legal frameworks to address this novel territory. Laws surrounding privacy, identity, and intellectual property rights are being reconsidered and adapted to encompass the unique challenges posed by deepfakes. As technology evolves, the law must keep pace to ensure that individuals' rights are adequately protected.

Accountability and Malicious Uses :

Deepfake technology has a dual nature. While it holds the potential for creativity and innovation, it can also be maliciously exploited. From disinformation campaigns to cybercrimes, the implications of malevolent deepfake usage are vast. Legal accountability becomes crucial in holding individuals or entities responsible for using deepfake technology for harmful purposes. Legal consequences serve as a deterrent and safeguard against the misuse of this technology.



Balancing Innovation and Ethics :

The ethical and legal considerations surrounding deepfakes exemplify the delicate balancing act between technological innovation and ethical responsibility. While technological advancements provide unprecedented creative potential, they also bring forth responsibilities to ensure that these innovations do not infringe upon fundamental human rights or societal well-being.

Education and Ethical Literacy :

A key aspect of addressing these ethical dilemmas involves raising awareness and fostering ethical literacy. Educating individuals about the implications of deepfake technology can empower them to make informed decisions about their online presence and interactions. It can also help society collectively discern between authentic and manipulated content, reducing the potential harm caused by deepfake misinformation.

VII. CONCLUSION

In the realm of content creation, the emergence of deepfake technology, whether in video or audio format, represents a fascinating yet complex chapter. This review has delved into the intricacies of its techniques, applications, challenges, and ethical considerations, shedding light on a landscape that blends innovation with potential risks.

From its inception, deepfake technology has demonstrated a remarkable ability to generate content that blurs the line between reality and fabrication. Whether seamlessly integrating actors into historical footage, enabling multilingual film dubbing, or igniting new horizons in entertainment, deepfakes in the video domain have shown their creative prowess. Similarly, deepfake audio has transformed voice cloning and synthesis, enabling advancements in voiceover work and personalization.

Yet, this transformative potential also brings forth ethical dilemmas. The unauthorized use of individuals' likenesses and voices raises valid concerns about consent, privacy, and identity. However, as these technologies continue to evolve, so too do detection methods and regulatory frameworks. The ongoing arms race between creators and detectors highlights our collective commitment to ensuring that technology remains a tool for progress rather than deception. Deepfakes hold immense potential to elevate creativity, communication, and entertainment. Through collaborative efforts among researchers, policymakers, and industry leaders, we have the opportunity to steer the trajectory of deepfake technology towards positive outcomes.



SJIF Rating: 8.176

ISSN: 2582-3930

VIII. REFERENCES

[1] Agarwal, S., Gurpinar, E., & Agrawal, G. (2019). Protecting World Leaders Against Deep Fakes: Deep Fake Image Detection Based on Color Spectral Information. arXiv preprint arXiv:1904.07439.

[2] Rössler, A., Cozzolino, D., Verdoliva, L., Riess, C., Thies, J., & Nießner, M. (2019). FaceForensics++: Learning to detect manipulated facial images. In Proceedings of the IEEE International Conference on Computer Vision (ICCV).

[3] Sabir, A., & Khan, F. S. (2019). Recurrent Convolutional Strategies for Video Face Forgery Detection. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops.

[4] Matern, F., Riess, C., Stamminger, M., & Nießner, M. (2019). Exploiting temporal information for real-time deepfake detection. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW).

[5] Marra, F., Gragnaniello, D., Cozzolino, D., Verdoliva, L., & Poggi, G. (2020). A Lightweight CNN for Deepfake Detection. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW).

[6] Hsu, Y. T., Averbuch-Elor, H., Häne, C., & Matusik, W. (2020). Deepfake detection using recurrent neural networks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW).

[7] Rossler, A., Cozzolino, D., Thies, J., & Nießner, M. (2020). Faceforensics: A large-scale video dataset for forgery detection in human faces. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops.

[8] Li, Y., Yang, J., Luo, Z., Li, S., & Lyu, S. (2020). Celeb-DF: A Large-Scale Challenging Dataset for DeepFake Forensics. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

[9] Nguyen, T. M., & Nguyen, H. T. (2020). Deepfake Video Detection through Optical Flow-based CNN and Entropy Analysis. In 2020 IEEE Fifth International Conference on Data Science in Cyberspace (DSC).

[10] Li, Y., Chang, H., Wang, X., Ji, H., & Lyu, S. (2020). In ictu oculi: Exposing AI created fake videos by detecting eye blinking. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).