

# Developing Robust Detection Systems for Deepfake Media Using Advanced AI and Data Analytics Techniques to Enhance Cyber Security

Temitope Olubunmi Awodiji<sup>1</sup>, John Owoyemi<sup>2</sup>

Department of Information Technology, University of the Cumberlands, Kentucky, USA<sup>1</sup>

Department of Information Technology, University of the Cumberlands, Kentucky, USA<sup>2</sup>

**Abstract:** To improve cybersecurity, this study looks into the creation of a reliable deepfake media detection system employing cutting-edge AI and data analytics approaches. A hybrid detection approach is presented by combining Recurrent Neural Networks (RNNs) for temporal analysis, Generative Adversarial Networks (GANs) for adversarial training, and Convolutional Neural Networks (CNNs) for feature extraction. Twenty professionals in the fields of cybersecurity, data analytics, and artificial intelligence were surveyed; the results were analysed using the Relative Importance Index (RII). The results underscore the hybrid model's efficacy, expandability, versatility, and pragmatic nature, stressing its capacity for instantaneous processing and assimilation into current cybersecurity structures. By offering a thorough strategy to counter the growing danger of deepfake media, this study seeks to improve the security and dependability of digital environments.

**Keywords:** Deepfake Detection, Convolutional Neural Networks (CNNs), Generative Adversarial Networks (GANs), Recurrent Neural Networks (RNNs) and Cybersecurity

## I. INTRODUCTION

The proliferation of deepfake media, characterized by the replacement of a person's likeness with another's in synthetic media, has emerged as a significant challenge to individual privacy, public trust, and national security [1]. This research is centered on the development of robust detection systems for deepfake media through the application of advanced artificial intelligence and data analytics methodologies. Thus, the convergence of artificial intelligence and data analytics presents promising opportunities to combat the malevolent utilization of deepfake technology, which is progressively advancing in sophistication and accessibility.

The rapid advancement of deepfake technology has enabled the creation of highly realistic counterfeit videos and audio recordings that can deceive both individuals and current detection mechanisms [2]. These progressions introduce substantial risks across various sectors, including political manipulation, personal security, corporate espionage, and public trust [3].

Contextually, deepfakes have the potential to fabricate deceptive videos of political leaders, potentially influencing electoral outcomes and undermining democratic procedures [4]. As a result, individuals may become targets of deepfake media for purposes such as blackmail, defamation, and harassment. Organizations are susceptible to financial and reputational harm resulting from the dissemination of deepfake-generated misinformation [4] - the proliferation of deepfake media may diminish public trust in media, fostering widespread doubt and disorientation.

In light of these challenges, there exists an urgent necessity for efficient detection systems capable of accurately identifying deepfake media and mitigating its detrimental consequences [5]. The rationale behind this research is rooted in the escalating prevalence and complexity of deepfake technology. Conventional detection approaches, reliant on manual validation or basic algorithmic assessments, are inadequate to match the swift progress in deepfake generation techniques [6]. Thus, there is a compelling requirement to harness advanced artificial intelligence and data analytics for the development of more resilient and dependable detection systems.

Moreover, this research endeavors to address the following deficiencies in current research and application: refining the precision of deepfake detection to diminish false positives and negatives; devising systems capable of operating at scale, facilitating the analysis of extensive media content in real-time; and formulating adaptable detection frameworks that can evolve in response to emerging deepfake generation methods.

Thus, the principal objective of this research is to formulate and evaluate sophisticated artificial intelligence and data analytics methodologies to effectively identify deepfake media. Through this endeavor, the research aspires to enrich the realm of cybersecurity by furnishing tools and approaches that enhance the detection and mitigation of deepfake risks.

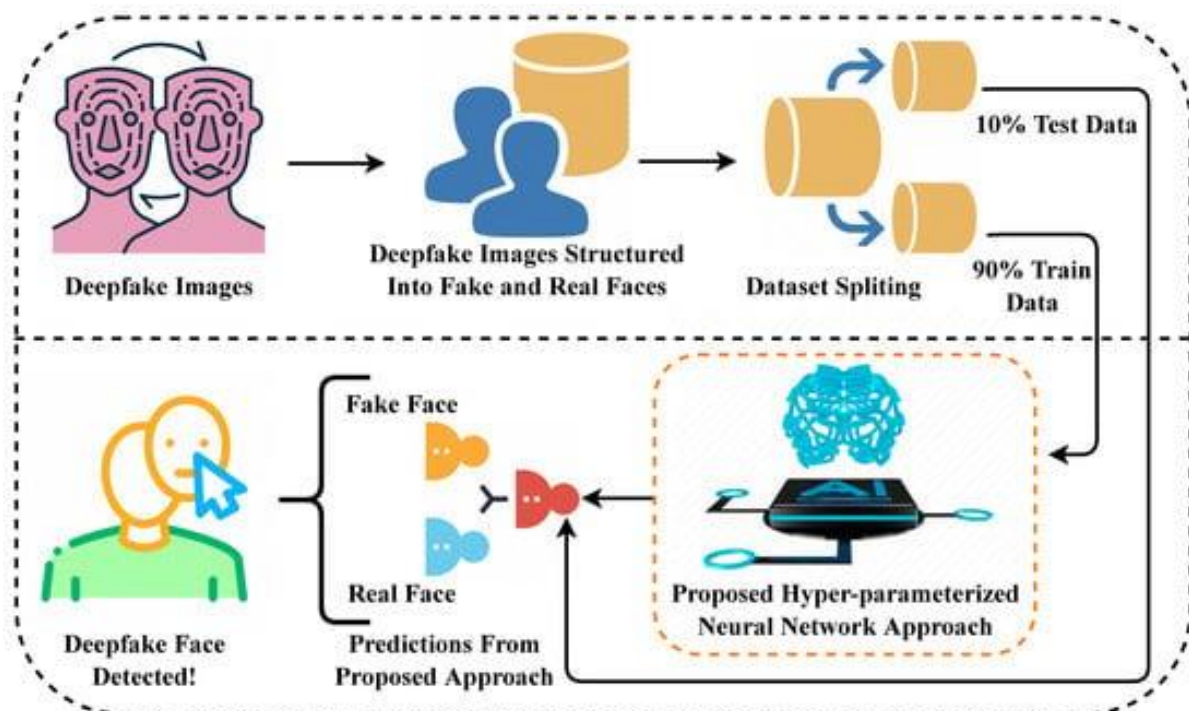
Specifically, the research aims to propose a model capable of accurately discriminating between authentic and deepfake media, employ data analytics to unveil patterns and irregularities suggestive of deepfake content, and establish a comprehensive framework for integrating these technologies into existing cybersecurity infrastructures.

## II. RELATED STUDIES

### AI and Deepfake Detection

MesoNet is a lightweight deep learning model that is specifically intended for deepfake detection [7]. This model is effective for real-time applications since it concentrates on the mesoscopic characteristics of video frames [8]. This work is noteworthy because it shows how to use compact neural networks to detect deepfakes in a way that balances computational efficiency and accuracy. This is in line with the objective of the current work, which is to create scalable and reliable detection systems.

Figure 1: AI and Deepfake Face Detection Framework

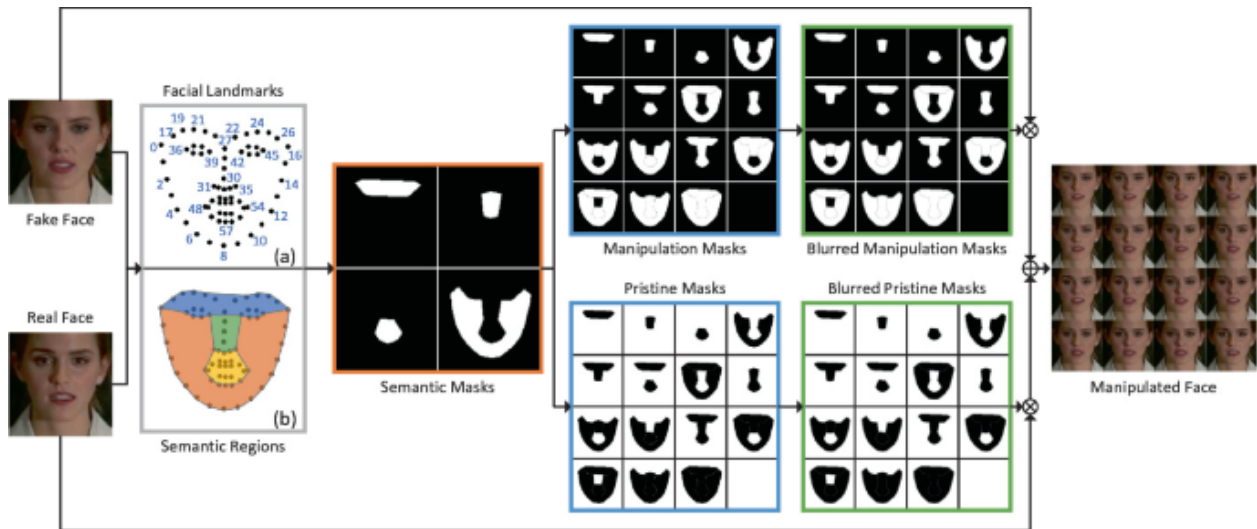


Source: Raza, Munir and Almutairi [9]

The multitask learning approach proposed by Weerawardana and Fernando in their relevant article "Multitask Learning for Detecting and Segmenting Manipulated Facial Images and Videos" allows for the simultaneous detection and segmentation of altered regions in facial photos and videos [10].

The detection performance of this technique is enhanced by utilising shared representations. This work highlights the role that multitask learning plays in improving deepfake detection capabilities and offers insights into how grouping related activities together might increase system resilience as a whole.

Figure 2: Image for Multitask Learning for Detection and Segmenting Manipulated Facial Images and Videos

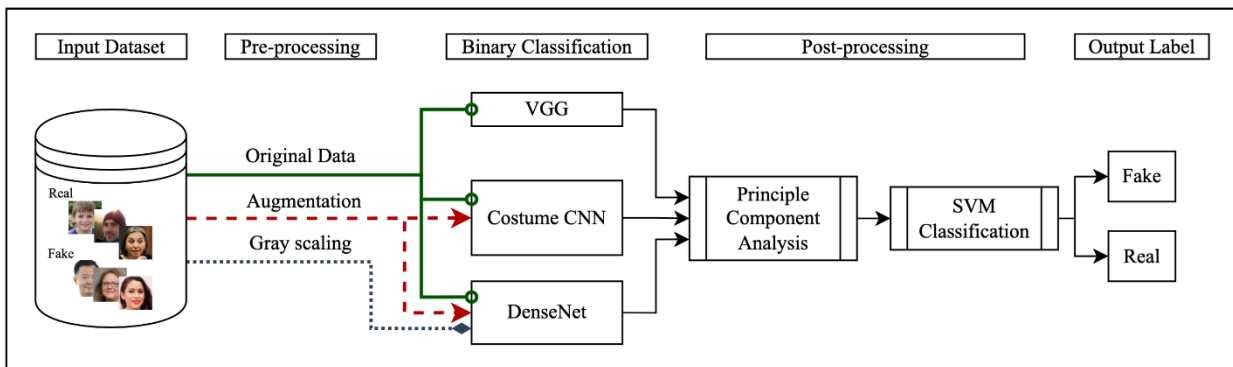


Source: Wang et al. [11]

**Data Analytics and Deepfake Detection**

In "FaceForensics++: Learning to Detect Manipulated Facial Images," the framework used a large-scale dataset and benchmarking framework to provide a thorough analysis of several detection algorithms [12]. FaceForensics++ assesses the potency of various methods for identifying alterations to facial images [12]. Deductively, the present research underscores the vital function that vast datasets perform in the training and assessment of deepfake detection models, hence directly advancing the creation of dependable detection systems.

Figure 3. Image for Data Analytics and Deepfake Detection



Source: Taeb and Chi [13]

The use of CNNs for identifying deepfake films by examining temporal and spatial discrepancies is investigated in the paper "Detecting Deepfakes Using Convolutional Neural Networks" [14]. The study shows that CNN-based models are good at spotting altered content. This study's emphasis on CNNs for deepfake detection is especially pertinent because it makes use of CNNs' capacity to analyse image and video data to precisely identify deepfakes.

**Ethical and Legal Perspectives**

The wider effects of deepfake technology on privacy, democracy, and national security are covered in the groundbreaking study "Deepfakes: A Looming Challenge for Privacy, Democracy, and National Security," [15]. The authors propose a comprehensive strategy that encompasses technological, legal, and policy approaches to tackle the issues raised by deepfakes. To frame the development and implementation of detection systems and provide a fundamental viewpoint for the ongoing study, it is imperative that these ethical and legal ramifications be understood.

Figure 4. Ethics of Deepfake



Source: Mukta et al. [16]

A study placed an emphasis on the necessity of moral standards and conscientious innovation in AI and related domains [17]. These moral issues are important to make sure that deepfake detection technology development follows moral guidelines, encouraging responsible use and reducing any risks.

### Technological Advancements and Future Directions

The research conducted on the "DeepFake Detection Challenge" delineates the outcomes of a substantial competition designed to propel the advancement of deepfake detection methodologies [18]. Through the provision of a standard dataset, this challenge scrutinized various detection approaches. The findings of this challenge furnish crucial insights into the existing capabilities and constraints of deepfake detection technologies, enriching ongoing research by spotlighting efficacious strategies and pinpointing areas necessitating enhancement [19].

A study delivers an exhaustive examination of media forensics strategies for identifying deepfakes [20]. The study delves into diverse methodologies, obstacles, and future pathways in this domain. This comprehensive overview proves pertinent as it amalgamates the present state of media forensics and deepfake detection, offering a comprehensive outlook on the assorted techniques and their efficacy. It functions as a valuable resource for situate the present study within the wider research milieu.

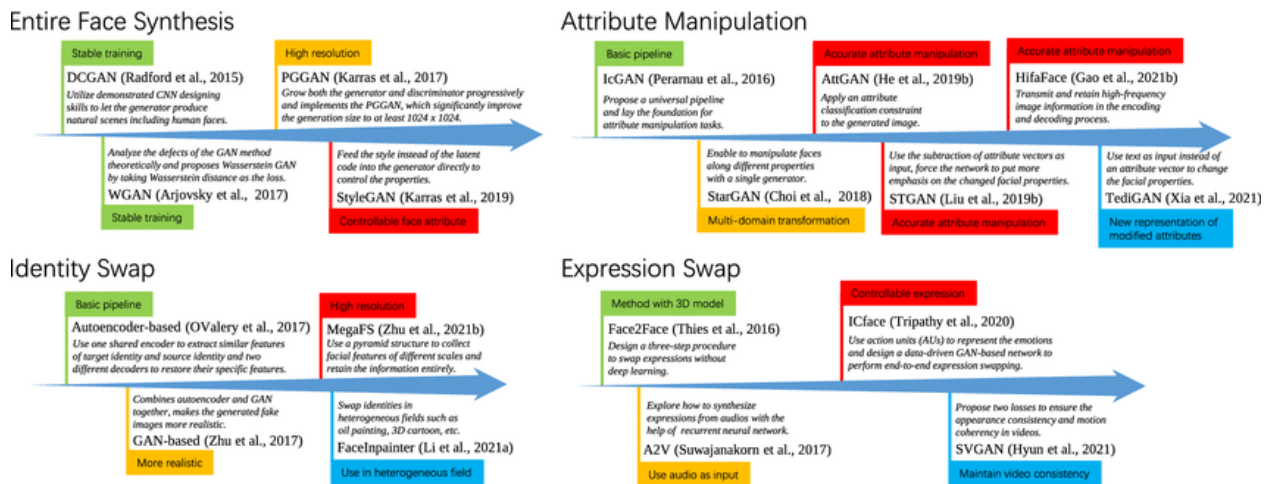
The aforementioned scholarly works furnish a sturdy groundwork for comprehending the current landscape of deepfake detection technologies. They underscore the effectiveness of AI-driven techniques, particularly Convolutional Neural Networks (CNNs) and Generative Adversarial Networks (GANs), the significance of extensive datasets for training and assessment, as well as the ethical and legal ramifications associated with deepfake technology [21]. By leveraging these insights, the current study endeavors to cultivate advanced AI and data analytics methodologies to amplify the identification and alleviation of deepfake risks, thereby enriching the realm of cybersecurity.

### Evolution and Impact of Deepfake Technology

Deepfake technology is one of the most important developments in artificial intelligence, with far-reaching effects on democracy, security, and privacy [15]. In the study, attention is drawn to the dual-use nature of deepfakes and point out that although they can be employed in innovative and advantageous ways, it is impossible to discount the possibility of harm they may cause through identity theft, misinformation, and other nefarious actions [21]. The authors stress that in order to stop this technology from being abused, strong detection methods are required.



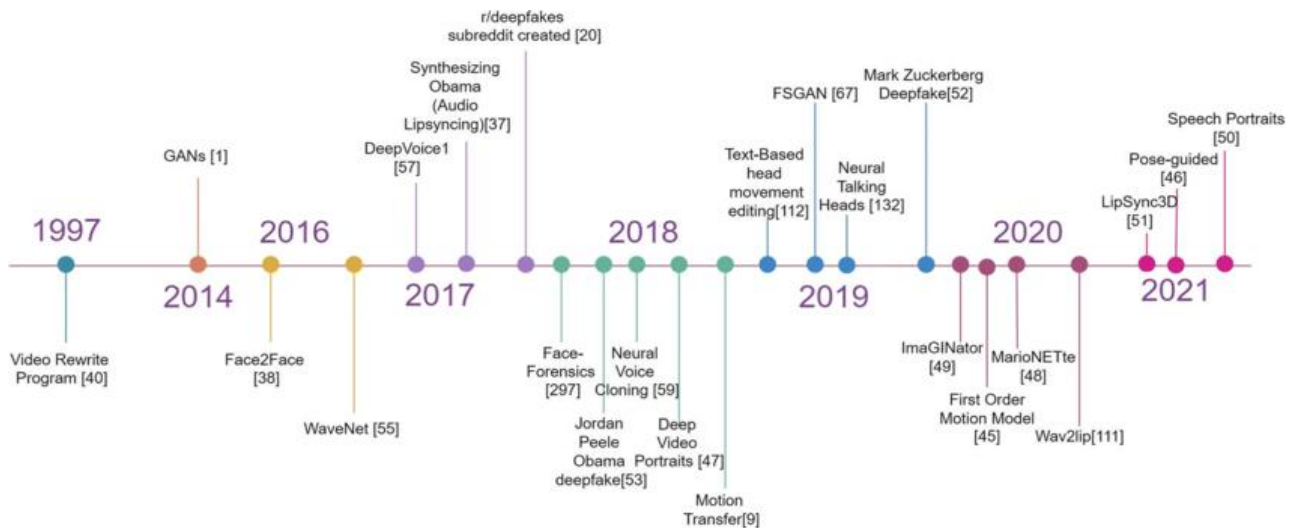
Figure 5. The evolution of DeepFake generation techniques



Source: Juefei-Xu et al. [22]

Investigating how advances in technology have made it possible to produce remarkably lifelike deepfakes is reflected contemporarily in a particular study [23]. They point out that the time and skill needed to create convincing deepfakes have drastically decreased thanks to advancements in generative adversarial networks (GANs) and other AI techniques. Since technology has become more accessible, there is a greater chance of misuse, which calls for the creation of advanced detection techniques [5].

Figure 6. Timeline of the Evolution of Deepfakes



Source: Masood et al. [2]

On the other hand, Vaccari and Chadwick offer a more pessimistic perspective, contending that the perceived danger posed by deepfakes might be exaggerated [4]. Since the majority of deepfakes can still be identified by the typical viewer, there is an argument that the impact of deepfakes on public opinion and democracy has been minimal thus far [24]. The authors advocate for a moderate approach and warn against excessive regulation, which could impede advancements in artificial intelligence and related areas.

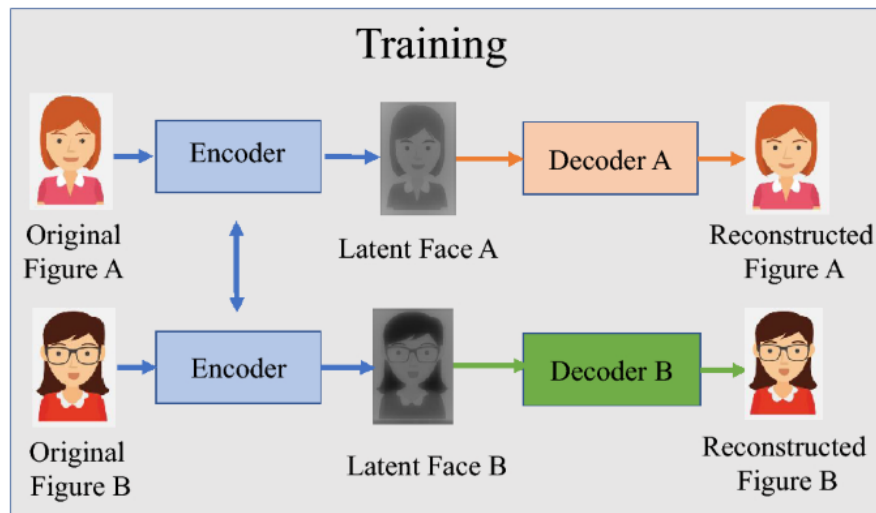
According to Hwang et al., the excitement around deepfakes can cause unwarranted anxiety and disinformation about their true potential [25]. The study contended that although deepfakes are a problem, other types of digital manipulation and false information that may be more common and simple to create should also be given attention.

**Detection Techniques for Deepfake Media**

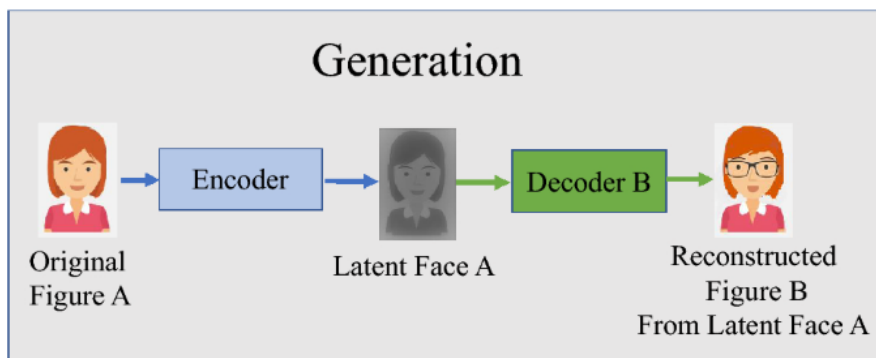
In a study of several AI-based detection methods, some proponents emphasise how successful these methods are at spotting deepfake content [26]. This study portrays techniques like recurrent neural networks (RNNs) and convolutional neural networks (CNNs), which have demonstrated excellent accuracy in identifying deepfakes by examining discrepancies in audio-visual synchronisation, lighting, and face gestures. The authors contend that ongoing advancements in AI detection techniques can match the rapid advancement of deepfake technology.

Thus, studies describe a unique method for deepfake detection utilising capsule networks [27]. This indicated that capsule networks outperform conventional CNN-based techniques in capturing spatial hierarchies and detecting minute anomalies in deepfake films [27].

Figure 7. Detection for Deepfake Media



(a) Training Phase



(b) Generation Phase

Source: Mitra et al. [28]

Although AI-based detection techniques are promising, a study cautioned that they are not infallible and can be thwarted by increasingly complex deepfake generating techniques [28]. The study draw attention to the continuous arms race that exists between deepfake producers and detectors, implying that no detection technique will ever be completely reliable [28].

The shortcomings of the detection techniques used contemporarily, especially their vulnerability to adversarial assaults, are covered by [29]. They contend that those who create deepfakes can purposefully create content to take advantage of flaws in detection algorithms, making them useless. Hence, to lessen the impact of deepfakes, a more all-encompassing approach that incorporates ongoing monitoring, updating detection algorithms, and public awareness efforts was presented [30].

### Convolutional Neural Networks (CNNs)

One type of deep learning algorithms that is very useful for processing structured grid data, like photographs, is Convolutional Neural Networks (CNNs) [31]. CNN have proven to be very successful at tasks like object and picture recognition and, more recently, anomaly detection like deepfakes [31].

#### Structure and Functioning

**A typical CNN architecture consists of several key layers:**

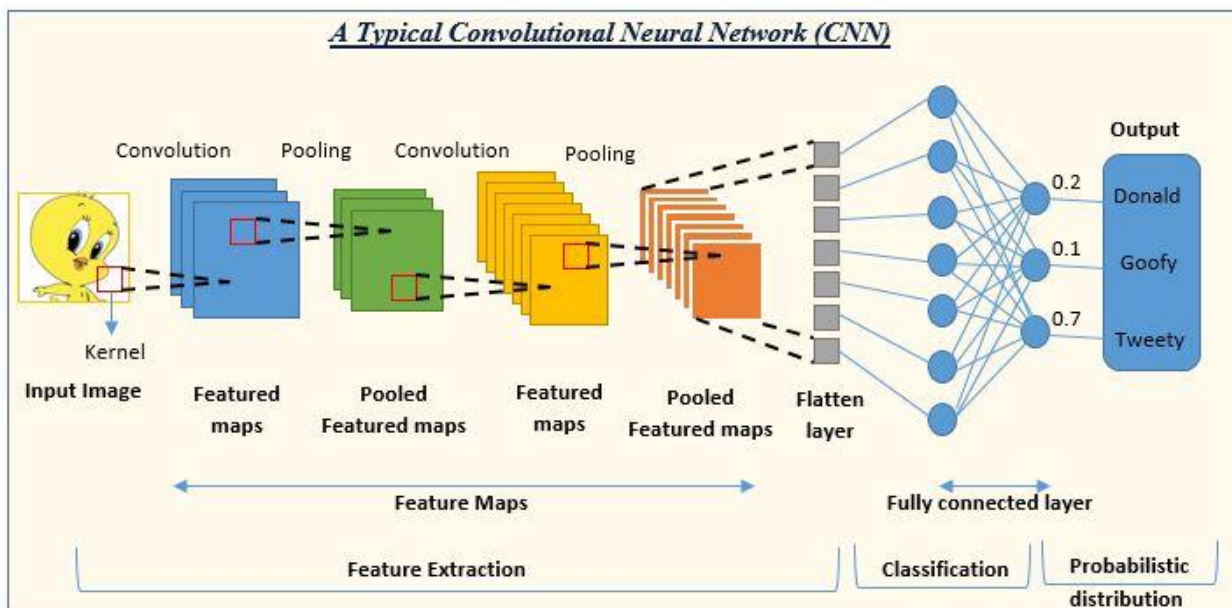
**Convolutional Layers:** These layers create feature maps by applying a collection of filters, sometimes known as kernels, to the input image. Sliding across the input image, the filters carry out convolution operations to pick up details like edges, textures, and shapes that are specific to the area [31].

**Pooling Layers:** In order to minimise computational complexity and preserve important characteristics, these layers downsample the feature maps' spatial dimensions. Max pooling and average pooling are two common pooling methods [31].

**Fully Connected Layers:** The network usually consists of one or more fully connected layers that interpret and make predictions from the retrieved information, following a number of convolutional and pooling layers. They resemble classical neural networks in that they have these layers [32].

**Activation Functions:** Non-linear activation functions, such as ReLU (Rectified Linear Unit), are applied to introduce non-linearity into the model, enabling it to learn complex patterns [32].

Figure 8. Convolutional Neural Networks (CNNs)



Source: Shah [33]

### Applications in Deepfake Detection

CNNs are particularly effective for deepfake detection due to their ability to capture and analyze fine-grained features in images and videos [34]. By training on large datasets of real and fake media, CNNs can learn to identify subtle inconsistencies and artefacts that are characteristic of deepfakes.

**Feature Extraction:** CNNs can detect anomalies in facial movements, lighting, and texture that are often present in deepfake media [34].

**Temporal Analysis:** When applied to video data, CNNs can analyze frame-by-frame consistency, detecting discrepancies in motion and expression [34].

**Generative Adversarial Networks (GANs)**

A type of machine learning frameworks called Generative Adversarial Networks (GANs) is made to produce artificial data that closely mimics real data. They are made up of a discriminator and a generator neural network, which are trained concurrently in a competitive environment [35].

**Structure and Functioning**

**Generator:** Using random noise, the generator network generates synthetic data. Its goal is to generate data that is identical to actual data.

**Discriminator:** The discriminator network analyses the information and makes a distinction between synthetic and real data. It gives the generator suggestions on how to make the synthetic data more realistic.

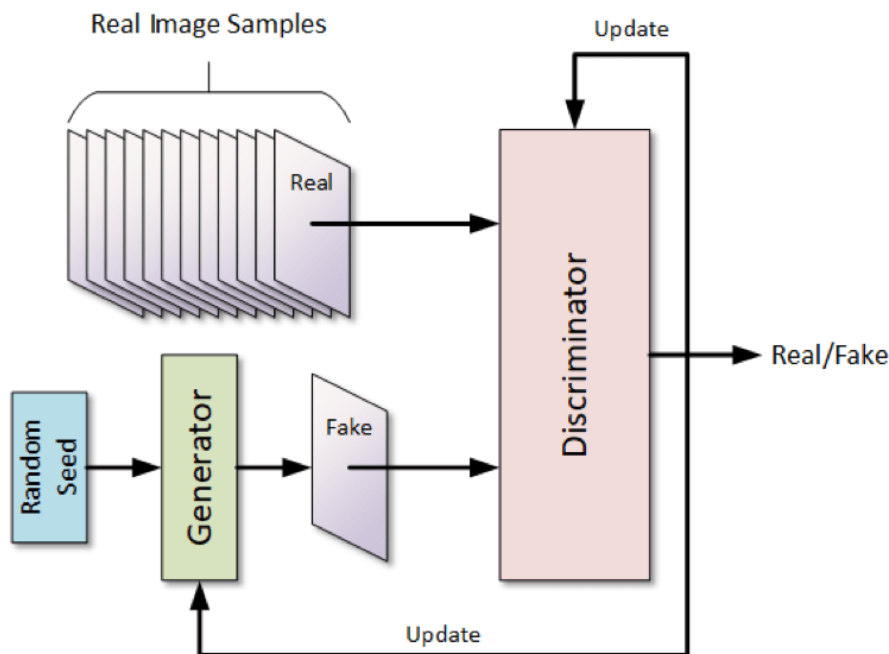
During the training process, the discriminator attempts to correctly distinguish between actual and bogus data, while the generator tries to trick it. The generator keeps going through this adversarial process until it generates extremely realistic data [35].

**Applications in Deepfake Creation and Detection**

The fundamental technology used to create deepfakes is GANs. With extensive training on real-world media datasets, they are able to produce remarkably lifelike photos, movies, and audio samples [36].

Thus, GANs are able to create realistic movements and sounds, interchange faces, and adjust attitudes in order to synthesise photos and videos of people [36]. It's interesting to note that deepfake detection can also be improved by GANs. Through discriminator training on generated and actual data, researchers may create models that are very good at detecting deepfakes [2].

Figure 9. Generative Adversarial Networks (GANs)



Source: Moyer [22221]

In contrast, deepfake media relies heavily on Convolutional Neural Networks (CNNs) and Generative Adversarial Networks (GANs) [2]. While GANs are essential to both the production of deepfakes and the development of sophisticated detection methods, CNNs are critical to the detection of deepfakes through the analysis of complex patterns and inconsistencies in images and videos [35]. Gaining an understanding of and making use of the advantages of both CNNs and GANs is crucial to building resilient systems that combat the growing threat of deepfake media, improving cybersecurity, and upholding public confidence.



### **Integration of Detection Systems in Cybersecurity Frameworks**

Nguyen et al. deliberates on the incorporation of deepfake detection mechanisms within broader frameworks of cybersecurity [36]. The contention posited is the indispensability of such amalgamation in furnishing a holistic defense strategy against digital menaces. The scholars underscore the capacity of artificial intelligence (AI) and data analytics in fortifying the prowess of cybersecurity systems through heightened detection, prevention, and response capabilities, thereby fortifying resilience against emerging perils like deepfakes.

Bajpai et al. elucidate a case study portraying the efficacious assimilation of deepfake detection systems within a corporate cybersecurity architecture [37]. The exposition showcases the utility of sophisticated AI methodologies for real-time monitoring and scrutiny of media content, effectively discerning and alleviating the risks posed by deepfakes [37].

Kumar et al. advise on the challenges accompanying the integration of deepfake detection systems into cybersecurity frameworks [38]. The study draws attention to impediments such as elevated computational expenses, plausible privacy apprehensions, and the imperative of recurrent updates to align with the evolving landscape of deepfake techniques.

Shen et al. explains the likelihood of erroneous identifications in AI-driven detection systems, which could instigate redundant alerts and resource allocation [39]. The study contended that a sole dependence on automated detection is inadequate and underscores the necessity of human proficiency and supervision.

### **Ethical and Legal Considerations**

Scrutinizing the ethical and legal implications linked to deepfakes, emphasizes the urgent need for regulatory frameworks to oversee their potential misuse. They contend that existing legal provisions are insufficient in tackling the unique challenges posed by deepfakes, which encompass concerns regarding consent, defamation, and violations of privacy [40]. On the other hand, exploring the ethical quandaries associated with deepfake technology, specifically its impact on trust and authenticity within the realm of digital media is pertinent. Studies including [37] and [41] argue that deepfakes undermine the reliability of authentic media sources and diminish public trust, highlighting the necessity for ethical guidelines and rules to govern their usage. The studies unilaterally suggest collaborative initiatives involving technology experts, policymakers, and societal entities to establish a responsible framework for deepfake technology.

Overall, the literature on deepfake media detection offers a complicated picture of scientific developments, moral dilemmas, and real-world difficulties. Although perspectives on the efficacy of present approaches, the proper balance between regulation and innovation, and the optimal ways for incorporating detection systems into cybersecurity frameworks are divided [5]; there is general agreement on the necessity of robust detection systems [43]. By creating advanced AI and data analytics methods for deepfake detection, this work seeks to further the current conversation by tackling the ethical and technical aspects of this urgent problem.

## **III. METHODOLOGY**

This study seeks objectivity and dependability in the gathering and analysis of data by employing a quantitative methodology based in positivist philosophy. Twenty professionals with vast expertise and notable contributions to the fields of data analytics, artificial intelligence, and cybersecurity are involved.

A systematic poll will be used to gather information about the experts' opinions of several AI and data analytics methods for deepfake detection. The efficacy, scalability, flexibility, and practicality of these methods will all be covered by the survey questions.

The analysis of responses will use the Relative Importance Index (RII) method. To calculate the RII, we use a straightforward formula:

$$RII = \frac{\text{Sum of Scores}}{\text{Highest Factor} \times \text{Total Respondents}}$$

Eq:1

The RII will assist us in prioritising the elements according to our assessment of their significance, giving us a clear image of which methods are deemed most successful by professionals. These realisations will direct the analysis and interpretation of the findings, ultimately aiding in the creation of reliable deepfake detection systems.

### Analysis and Interpretation

The replies were analysed using RII, and the results show that the methodologies of CNNs, GANs, SVMs, and RNNs all received excellent marks for practicality, scalability, efficacy, and adaptability—with CNNs and RNNs performing especially well in a number of categories.

Large dataset accessibility, simplicity of integration, and an intuitive user interface were among the other elements that received excellent marks, demonstrating their significance in the creation of reliable deepfake detection systems.

Figure 10. Response on Effectiveness of Deepfake Frameworks

Effectiveness: The technique accurately detects deepfakes.

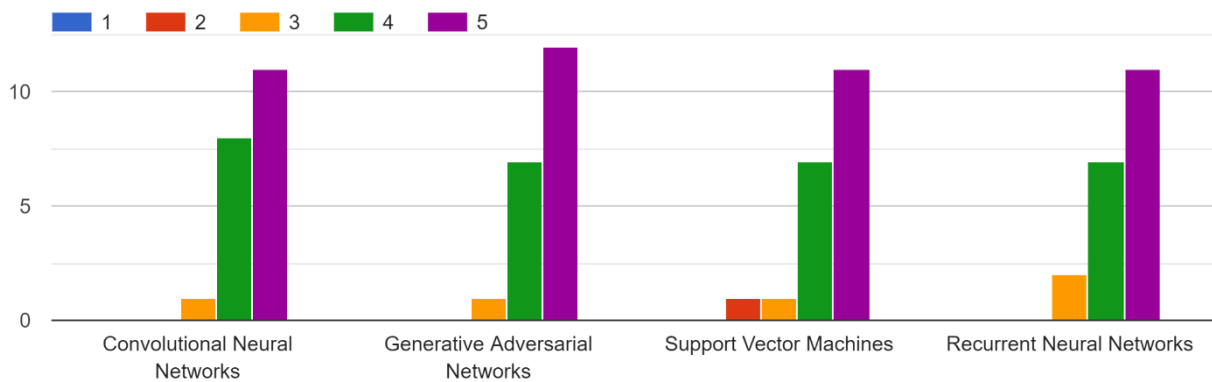


Figure 11. Response on Scalability of Deepfake Frameworks

Scalability: The technique can be applied to large datasets in real time.

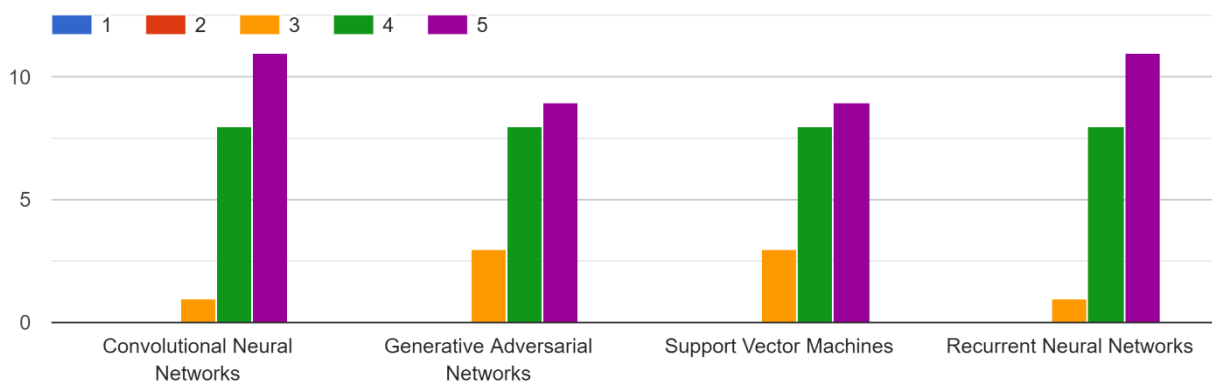


Figure 12. Response on Adaptability of Deepfake Frameworks

Adaptability: The technique can be easily updated to counter new deepfake generation methods.

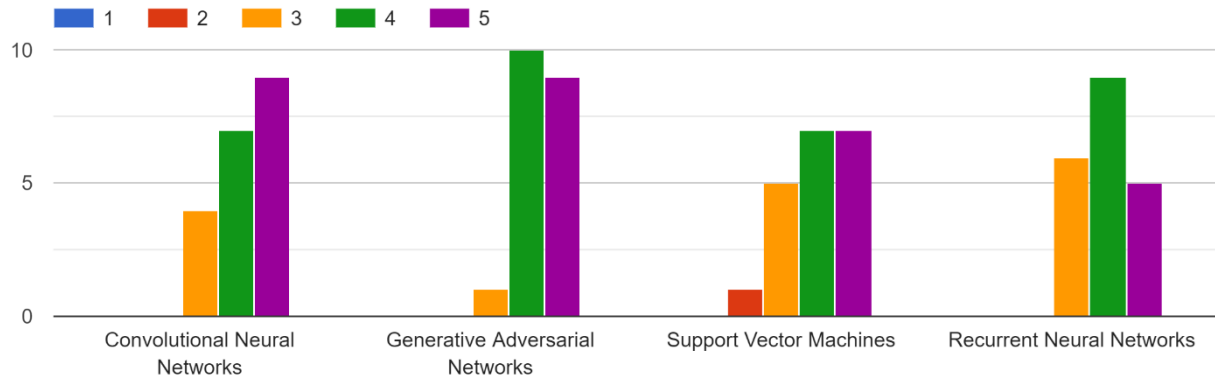
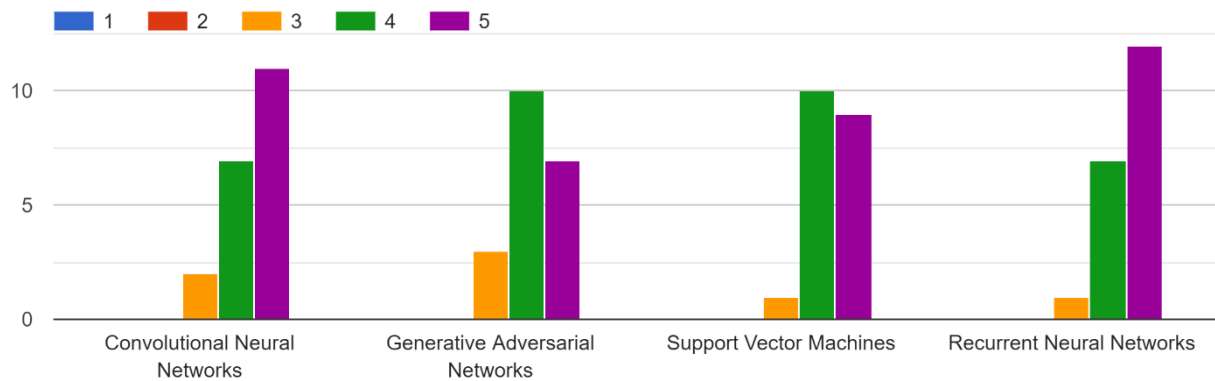


Figure 13. Response on Practicability of Deepfake Frameworks

Practicability: The technique is feasible for integration into existing cybersecurity frameworks.



Hence:

Table 1. Rate of Effectiveness, Scalability, Adaptability and Practicability of Deepfake Frameworks

|                |            |            |            |
|----------------|------------|------------|------------|
| Effectiveness  |            |            |            |
| CNN = 0.90     | GAN = 0.91 | SVM = 0.87 | RNN = 0.89 |
| Scalability    |            |            |            |
| CNN = 0.90     | GAN = 0.86 | SVM = 0.86 | RNN = 0.90 |
| Adaptability   |            |            |            |
| CNN = 0.85     | GAN = 0.88 | SVM = 0.79 | RNN = 0.79 |
| Practicability |            |            |            |
| CNN = 0.89     | GAN = 0.84 | SVM = 0.88 | RNN = 0.91 |

Table 2. Rating the Additional Factors

| Factor                        | Rating |
|-------------------------------|--------|
| Availability of large Dataset | 0.86   |
| Compositional Efficiency      | 0.82   |
| Ease of Integration           | 0.86   |
| Real time Processing          | 0.84   |
| Cost Implementations          | 0.82   |
| Accuracy in Diverse Scenarios | 0.86   |
| Update Frequency              | 0.82   |
| User-Friendly Interface       | 0.88   |

**Designing a Distinctive Method for Deepfake Detection**

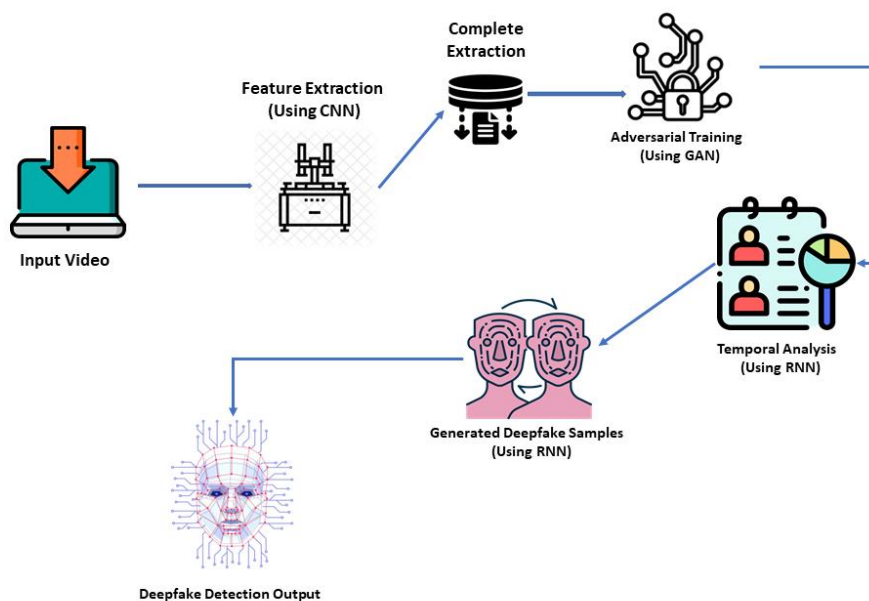
Given the high RII scores for CNNs and GANs across various criteria, a hybrid approach that combines the strengths of both techniques can be designed for deepfake detection. This can be labelled Hybrid Deepfake Detection System (HDDS).

**Hybrid Deepfake Detection System (HDDS) Framework**

- Step 1: Initial feature extraction and categorisation using convolutional neural networks (CNNs). CNNs are useful for spotting the visual irregularities typical of deepfakes since they are good at finding patterns and features in photos and movies.
- Step 2: Adversarial training using Generative Adversarial Networks (GANs). CNNs can continuously enhance their detection capabilities by using the synthetic deepfakes that GANs can produce. Through the use of an adversarial strategy, the detection system is kept resilient and adaptable to new techniques for creating deepfakes.
- For temporal analysis, use recurrent neural networks (RNNs) in step three. RNNs can be used to examine video sequences and make sure that temporal irregularities are found, which gives the system an extra degree of resilience.

**Integration and Real-time Processing**

Figure 14. HDDS Framework





With the use of parallel processing techniques and optimised algorithms, the system can process massive datasets in real-time with ease. Additionally, it features an intuitive UI that was designed to make it simple to use and integrate with current cybersecurity frameworks.

As a result of the HDDS's diversification, the system is updated frequently with fresh deepfake samples produced by GANs to guarantee its continued efficacy against new threats. Because of its scalability, this system can handle growing data volumes without experiencing performance issues. Therefore, a strong and unique deepfake detection system can be created to support cybersecurity efforts against the growing threat of deepfake media by utilising the strengths of CNNs, GANs, and RNNs and concentrating on important factors like real-time processing, ease of integration, and continuous updating.

#### IV. CONCLUSION

To sum up, this paper offers a thorough method for creating a reliable deepfake detection system using cutting-edge AI and data analytics approaches. A hybrid model is created to improve the scalability, adaptability, accuracy, and practicality of deepfake detection by combining Convolutional Neural Networks (CNNs) for feature extraction, Generative Adversarial Networks (GANs) for adversarial training, and Recurrent Neural Networks (RNNs) for temporal analysis. The model's efficacy in real-time processing and its viability for inclusion into current cybersecurity frameworks are highlighted by the Relative Importance Index (RII) analysis of expert survey responses. The suggested hybrid approach offers a dependable way to protect digital information and improve cybersecurity measures in response to the growing threat of deepfake media. Subsequent investigations ought to concentrate on the model's ongoing upgrading and enhancement.

#### REFERENCES

- [1] Nnamdi, N., Oniyinde, O. A., & Abegunde, B. (2023). An Appraisal of the Implications of Deep Fakes: The Need for Urgent International Legislations. *American Journal of Leadership and Governance*, 8(1), 43-70.
- [2] Masood, M., Nawaz, M., Malik, K. M., Javed, A., Irtaza, A., & Malik, H. (2023). Deepfakes generation and detection: State-of-the-art, open challenges, countermeasures, and way forward. *Applied Intelligence*, 53(4), 3974-4026.
- [3] Ness, S., & Khinvasara, T. (2024). Emerging Threats in Cyberspace: Implications for National Security Policy and Healthcare Sector. *Journal of Engineering Research and Reports*, 26(2), 107-117.
- [4] Vaccari, C., & Chadwick, A. (2020). Deepfakes and disinformation: Exploring the impact of synthetic political video on deception, uncertainty, and trust in news. *Social media+ society*, 6(1), 2056305120903408.
- [5] Amodu, O. (2024). Defending the defenseless in cyber-war. *Journal of Multidisciplinary Engineering Science and Technology*, 11(5), 2458-9403
- [6] Tufchi, S., Yadav, A., & Ahmed, T. (2023). A comprehensive survey of multimodal fake news detection techniques: advances, challenges, and opportunities. *International Journal of Multimedia Information Retrieval*, 12(2), 28.
- [7] Xia, Z., Qiao, T., Xu, M., Wu, X., Han, L., & Chen, Y. (2022). Deepfake video detection based on MesoNet with preprocessing module. *Symmetry*, 14(5), 939.
- [8] Javed, M., Zhang, Z., Dahri, F.H. and Laghari, A.A., 2024. Real-Time Deepfake Video Detection Using Eye Movement Analysis with a Hybrid Deep Learning Approach. *Electronics*, 13(15), p.2947.
- [9] Raza, A., Munir, K., & Almutairi, M. (2022). A novel deep learning approach for deepfake image detection. *Applied Sciences*, 12(19), 9820.
- [10] Weerawardana, M. C., & Fernando, T. G. I. (2021, August). Deepfakes detection methods: A literature survey. In *2021 10th International Conference on Information and Automation for Sustainability (ICIAfS)* (pp. 76-81).
- [11] Wang, R., Yang, Z., You, W., Zhou, L., & Chu, B. (2022). Fake face images detection and identification of celebrities based on semantic segmentation. *IEEE Signal Processing Letters*, 29, 2018-2022.
- [12] Rossler, A., Cozzolino, D., Verdoliva, L., Riess, C., Thies, J., & Nießner, M. (2019). Faceforensics++: Learning to detect manipulated facial images. In *Proceedings of the IEEE/CVF international conference on computer vision* (pp. 1-11).
- [13] Taeb, M., & Chi, H. (2022). Comparison of deepfake detection techniques through deep learning. *Journal of Cybersecurity and Privacy*, 2(1), 89-106.
- [14] Islam, H. M. (2022). *Detecting Deepfakes Using Convolutional Neural Networks* (Master's thesis, Texas A&M University-Kingsville).
- [15] Chesney, B., & Citron, D. (2019). Deep fakes: A looming challenge for privacy, democracy, and national security. *California Literature Review*, 107, 1753.
- [16] Mukta, M. S. H., Ahmad, J., Raiaan, M. A. K., Islam, S., Azam, S., Ali, M. E., & Jonkman, M. (2023). An investigation of the effectiveness of deepfake models and tools. *Journal of Sensor and Actuator Networks*, 12(4), 61.

- [17] Ashok, M., Madan, R., Joha, A., & Sivarajah, U. (2022). Ethical framework for Artificial Intelligence and Digital technologies. *International Journal of Information Management*, 62, 102433.
- [18] Dolhansky, B., Bitton, J., Pflaum, B., Lu, J., Howes, R., Wang, M., & Ferrer, C. C. (2020). The deepfake detection challenge (dfdc) dataset. *arXiv preprint arXiv:2006.07397*.
- [19] Alhaji, H. S., Celik, Y., & Goel, S. (2024). An Approach to Deepfake Video Detection Based on ACO-PSO Features and Deep Learning. *Electronics*, 13(12), 2398.
- [20] Verdoliva, L. (2020). Media forensics and deepfakes: an overview. *IEEE journal of selected topics in signal processing*, 14(5), 910-932.
- [21] Shree, M. S., Arya, R., & Roy, S. K. (2024). Investigating the Evolving Landscape of Deepfake Technology: Generative AI's Role in it's Generation and Detection. *International Research Journal on Advanced Engineering Hub (IRJAEH)*, 2(05), 1489-1511.
- [22] Juefei-Xu, F., Wang, R., Huang, Y., Guo, Q., Ma, L., & Liu, Y. (2022). Countering malicious deepfakes: Survey, battleground, and horizon. *International journal of computer vision*, 130(7), 1678-1734.
- [23] Kietzmann, J., Lee, L. W., McCarthy, I. P., & Kietzmann, T. C. (2020). Deepfakes: Trick or treat?. *Business Horizons*, 63(2), 135-146.
- [24] Dobber, T., Metoui, N., Trilling, D., Helberger, N., & De Vreese, C. (2021). Do (microtargeted) deepfakes have real effects on political attitudes?. *The International Journal of Press/Politics*, 26(1), 69-91.
- [25] Sheehy, B., Choi, S., Khan, M. I., Arnold, B. B., Sang, Y., & Lee, J. J. (2024). Truths and Tales: Understanding Online Fake News Networks in South Korea. *Journal of Asian and African Studies*, 00219096231224672.
- [26] Sandotra, N., & Arora, B. (2024). A comprehensive evaluation of feature-based AI techniques for deepfake detection. *Neural Computing and Applications*, 36(8), 3859-3887.
- [27] Wani, T. M., Gulzar, R., & Amerini, I. (2024). ABC-CapsNet: Attention based Cascaded Capsule Network for Audio Deepfake Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 2464-2472).
- [28] Mitra, A., Mohanty, S. P., Corcoran, P., & Kougianos, E. (2021). A machine learning based approach for deepfake detection in social media through key video frame extraction. *SN Computer Science*, 2(2), 98.
- [29] Girdhar, M., Hong, J., & Moore, J. (2023). Cybersecurity of autonomous vehicles: A systematic literature review of adversarial attacks and defense models. *IEEE Open Journal of Vehicular Technology*, 4, 417-437.
- [30] Gambín, Á. F., Yazidi, A., Vasilakos, A., Haugerud, H., & Djenouri, Y. (2024). Deepfakes: current and future trends. *Artificial Intelligence Review*, 57(3), 64.
- [31] Khan, A., Sohail, A., Zahoora, U., & Qureshi, A. S. (2020). A survey of the recent architectures of deep convolutional neural networks. *Artificial intelligence review*, 53, 5455-5516.
- [32] Alzubaidi, L., Zhang, J., Humaidi, A. J., Al-Dujaili, A., Duan, Y., Al-Shamma, O., ... & Farhan, L. (2021). Review of deep learning: concepts, CNN architectures, challenges, applications, future directions. *Journal of big Data*, 8, 1-74.
- [33] Shah, S. (2022). Convolution Neural Networks: An Overview. Analytics Vidhya [Online]. Available from: <https://www.analyticsvidhya.com/blog/2022/01/convolutional-neural-network-an-overview/>
- [34] Tran, V. N., Lee, S. H., Le, H. S., & Kwon, K. R. (2021). High performance deepfake video detection on cnn-based with attention target-specific regions and manual distillation extraction. *Applied Sciences*, 11(16), 7678.
- [35] Baowaly, M. K., Liu, C. L., & Chen, K. T. (2019, June). Realistic data synthesis using enhanced generative adversarial networks. In *2019 IEEE Second International Conference on Artificial Intelligence and Knowledge Engineering*, pp. 289-292.
- [36] Nguyen, T. T., Nguyen, Q. V. H., Nguyen, D. T., Nguyen, D. T., Huynh-The, T., Nahavandi, S., ... & Nguyen, C. M. (2022). Deep learning for deepfakes creation and detection: A survey. *Computer Vision and Image Understanding*, 223, 103525.
- [37] Anand, A., Madaan, A., & Danielsson, A. (Eds.). (2024). *Intersections Between Rights and Technology*. IGI Global.
- [38] Rana, M.S., Nobi, M.N., Murali, B. and Sung, A.H., 2022. Deepfake detection: A systematic literature review. *IEEE Access*, 10, pp.25494-25513.
- [39] Khan, I. U., Ouaisa, M., Ouaisa, M., Fayaz, M., & Ullah, R. (Eds.). (2024). *Artificial Intelligence for Intelligent Systems: Fundamentals, Challenges, and Applications*.
- [40] Kalpokas, I., & Kalpokiene, J. (2022). *Deepfakes: a realistic assessment of potentials, risks, and policy regulation*. Springer Nature.
- [41] Gregory, S. (2022). Deepfakes, misinformation and disinformation and authenticity infrastructure responses: Impacts on frontline witnessing, distant witnessing, and civic journalism. *Journalism*, 23(3), 708-729.
- [42] Hossain, S. T., Yigitcanlar, T., Nguyen, K., & Xu, Y. (2024). Local government cybersecurity landscape: A systematic review and conceptual framework. *Applied Sciences*, 14(13), 5501.
- [43] Hossain, S. T., Yigitcanlar, T., Nguyen, K., & Xu, Y. (2024). Local government cybersecurity landscape: A systematic review and conceptual framework. *Applied Sciences*, 14(13), 5501.