

# False Reality: Deepfakes in Terrorist Propaganda and Recruitment

Muhammad N. A. Latif Al Waroi<sup>1</sup>

<sup>1</sup>School of Strategic and Global Studies, Universitas Indonesia, Jakarta, Indonesia

<sup>1</sup>latif.alwaroi46@gmail.com

## Article Info

Received: 03-Jul-2024

Revised: 29-Jul-2024

Accepted: 11-Aug-2024

## Keywords

Deepfake Detection; Deepfake Technology; Policy Frameworks; Psychological Impact; Terrorist Propaganda

## Abstract

Deepfake technology, which leverages artificial intelligence to create hyper-realistic digital fabrications, has emerged as a significant threat across various domains, notably in terrorism. This review critically examines the exploitation of deepfakes in terrorist propaganda and recruitment, presenting a systematic analysis of the technical mechanisms behind their creation and detection, historical and contemporary propaganda methods, and their psychological impacts on audiences. The study identifies key advancements in deepfake detection technologies, such as ensemble learning and convolutional neural networks, which are crucial in distinguishing real from synthetic media. Furthermore, the review highlights the importance of public awareness and psychological resilience as vital countermeasures against deepfake manipulation. Despite technological advancements, significant challenges remain, including the development of real-time detection systems capable of operating in diverse and uncontrolled environments and a comprehensive understanding of the psychological processes affected by deepfake propaganda. The review underscores the urgent need for robust policy frameworks and international cooperation to address the ethical, legal, and security implications of deepfake technology. By integrating technical, psychological, and policy perspectives, this study provides a holistic understanding of deepfake technology's role in modern terrorism and offers insights for developing effective countermeasures. The comprehensive approach aims to contribute to the creation of robust strategies to mitigate the misuse of deepfake technology, ensuring a safer and more trustworthy digital environment.

## 1. Introduction

Deepfake technology, a form of artificial intelligence that creates hyper-realistic digital fabrications, has recently emerged as a powerful tool with significant implications across various domains, including entertainment, politics, and security (Abbas et al., 2023). This technology manipulates video and audio content to create false but convincing representations of real individuals, often leading to severe consequences. The ability to produce realistic deepfakes poses a substantial threat, particularly in the realms of misinformation and propaganda, as it undermines trust in authentic media and communications (Abir et al., 2023; Agarwal et al., 2021).

The intersection of deepfake technology with terrorism is a growing concern. Terrorist organizations have historically used propaganda to recruit and radicalize individuals. With the advent of deepfakes, these efforts can be significantly enhanced, making the dissemination of misleading and harmful content easier and more convincing (Amin, Hu, Li, et al., 2024). The potential of deepfakes to amplify the psychological impact of terrorist propaganda necessitates a thorough examination of this technology's role in modern terrorist activities (Amerini et al., 2021; H. Chen et al., 2023).

The primary research problem addressed in this review is the exploitation of deepfake technology by terrorist organizations to enhance their propaganda and recruitment strategies. This issue is crucial because deepfakes can create highly persuasive and seemingly authentic media that can effectively deceive

\*Corresponding author, email: [latif.alwaroi46@gmail.com](mailto:latif.alwaroi46@gmail.com)

doi: -

This work is licensed under a [Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International \(CC BY-NC-SA 4.0\)](https://creativecommons.org/licenses/by-nc-sa/4.0/)

and manipulate target audiences (Abbas et al., 2023; M. Chen et al., 2022). The integration of deepfakes into terrorist propaganda not only extends the reach of their messages but also increases the psychological impact on individuals, leading to higher levels of radicalization and recruitment (Amin, Hu, & Hu, 2024).

To address this problem, it is essential to systematically examine the technical mechanisms behind the creation and detection of deepfakes. By understanding how deepfakes are generated and identified, researchers and policymakers can develop more effective countermeasures. Additionally, analyzing the historical and contemporary use of propaganda by terrorist organizations, with a focus on the integration of deepfakes, provides insights into the evolving strategies used for psychological manipulation and recruitment (Appel & Prietzel, 2022; Chhabra et al., 2024). This comprehensive approach aims to mitigate the impact of enhanced terrorist activities through technical, psychological, and policy-related interventions.

Recent literature emphasizes the need for robust detection methods to counter the threats posed by deepfakes. Various deep learning techniques have been proposed to enhance deepfake content detection. For example, the use of ensemble learning and convolutional neural networks (CNNs) has shown promise in identifying deepfake manipulations with high accuracy (Abbas et al., 2023; Amerini et al., 2021). These methods leverage unique artifacts and inconsistencies introduced during the deepfake creation process, making them effective tools for distinguishing between real and fake content (M. Chen et al., 2022).

Furthermore, research highlights the importance of integrating temporal and spatial features for comprehensive deepfake detection. Techniques such as optical flow analysis and spatiotemporal convolutional networks have been developed to detect inconsistencies in movement and texture that are characteristic of deepfakes (Amin, Hu, & Hu, 2024; G.-L. Chen & Hsu, 2023). These approaches not only improve detection accuracy but also enhance the system's resilience to various types of deepfake manipulations.

In addition to technical solutions, the literature also emphasizes the importance of psychological resilience and public awareness as countermeasures against deepfake propaganda. Educating the public about the existence and risks of deepfakes can reduce vulnerability to manipulation (Appel & Prietzel, 2022; Chhabra et al., 2024). Psychological interventions aimed at building critical thinking and media literacy skills are crucial components of a comprehensive strategy to combat the influence of deepfake-enhanced terrorist propaganda.

Despite advancements in deepfake detection, several gaps remain in the literature. One major gap is the need for real-time detection systems that can operate effectively in diverse and uncontrolled environments. Most current detection methods are computationally intensive and may not be suitable for real-time applications, which are crucial for mitigating the immediate impact of deepfake propaganda (Abir et al., 2023; Amin, Hu, & Hu, 2024). Developing lightweight and efficient algorithms that can be deployed in real-time remains a significant challenge.

Another gap is the limited understanding of the psychological mechanisms through which deepfake propaganda affects audiences. While there is evidence that deepfakes can enhance the persuasive power of propaganda, more research is needed to elucidate the specific psychological processes involved (Appel & Prietzel, 2022; B. Chen et al., 2022). Understanding these mechanisms is essential for developing targeted interventions that can mitigate the psychological impact of deepfake propaganda.

Finally, there is a need for comprehensive policy frameworks that address the ethical and legal implications of deepfake technology. Current regulatory measures are often inadequate to address the rapid advancements in deepfake creation and dissemination (Abbas et al., 2023; Chhabra et al., 2024). Developing robust policies that balance the benefits of deepfake technology with its potential risks is crucial for mitigating the threat posed by deepfake-enhanced terrorist activities.

The objective of this review is to systematically examine the role of deepfake technology in the realm of terrorist propaganda and recruitment. Specifically, this review aims to understand the technical mechanisms and advancements in deepfake technology, analyze the historical and contemporary use of propaganda in terrorism and the integration of deepfakes, investigate the psychological impact of deepfake propaganda on audiences, examine the use of deepfakes in terrorist recruitment processes, and evaluate

preventive measures and response strategies to mitigate the impact of deepfake propaganda and recruitment.

The novelty of this research lies in its comprehensive approach to addressing the multifaceted threat posed by deepfakes in terrorist activities. By integrating insights from technical, psychological, and policy perspectives, this review aims to provide a holistic understanding of the impact of deepfake technology on the propaganda and recruitment strategies used by terrorist organizations. This integrative approach is crucial for developing effective countermeasures and informing policy decisions.

The scope of this review includes the intersection of deepfake technology, propaganda, and terrorist recruitment. It encompasses a technical review of deepfake creation and detection methods, an analysis of the evolution and effectiveness of deepfake-based propaganda, case studies illustrating the use of deepfakes in terrorist recruitment campaigns, and a discussion of psychological, technological, and policy-related countermeasures to the use of deepfakes. By providing a comprehensive and nuanced understanding of these issues, this review aims to contribute to the development of effective strategies for combating the threat of deepfake-enhanced terrorist activities.

## **2. Literature Review**

### **2.1. Evolution of Deepfake Technology**

#### **2.1.1. Historical Development of Deepfake Technology**

The development of deepfake technology is rooted in advancements in computer graphics and artificial intelligence, with the term "deepfake" emerging around 2017, combining "deep learning" and "fake" (Nguyen et al., 2022). Early iterations involved basic image manipulation, but Generative Adversarial Networks (GANs) significantly advanced the field by enabling highly realistic synthetic media creation (Kietzmann et al., 2020). Introduced by Ian Goodfellow in 2014, GANs consist of a generator and a discriminator working in tandem to produce convincing fake images and videos (Goodfellow et al., 2014), leading to significant improvements in synthetic media realism and complexity (Ferreira et al., 2021b).

#### **2.1.2. Key Advancements and Milestones**

Notable milestones include the introduction of DeepFaceLab, an open-source tool democratizing deepfake creation (Liu et al., 2023), and the application of deepfake technology across various domains such as entertainment and education. For instance, deepfake technology has enhanced visual effects and storytelling in films by creating realistic digital actors (Kietzmann et al., 2020). However, its misuse for fake news and misleading content has raised ethical and legal concerns (Casu et al., 2024a). Ongoing improvements in algorithms and computational power, including adversarial training, transfer learning, and the integration of audio-visual modalities, have increased deepfake realism and detection difficulty (Ding et al., 2022; Nguyen et al., 2022).

## **2.2. Technical Mechanisms of Deepfakes**

### **2.2.1. Algorithms and Techniques Used in Creating Deepfakes**

Deepfakes primarily rely on GANs, where a generator creates fake images and a discriminator distinguishes between real and fake images (Goodfellow et al., 2014). Variations like CycleGAN and StyleGAN have further advanced capabilities, enabling image transformation between domains without paired examples and generating high-resolution images with style and content control (Fu et al., 2019; Karras et al., 2020).

### **2.2.2. Machine Learning and AI Involvement**

Convolutional Neural Networks (CNNs) are extensively used for image and video processing tasks, including facial recognition and manipulation (Caldelli et al., 2021), enabling realistic facial expressions and movements. Transfer learning improves efficiency by leveraging pre-trained weights from large datasets, reducing the need for extensive computational resources and training data (Ferreira et al., 2021a). Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks handle temporal aspects in video deepfakes, ensuring consistency in frame transitions and lip synchronization (B. Chen et al., 2022).

## 2.3. Deepfake Detection Methods

### 2.3.1. Current Methods for Detecting Deepfakes

Detecting deepfakes remains challenging due to their increasing sophistication. Key approaches include:

- (1) Feature-based Methods: Identify artifacts and inconsistencies in synthetic media by analyzing color discrepancies, pixel-level anomalies, and unnatural facial movements (G.-L. Chen & Hsu, 2023; Chhabra et al., 2024).
- (2) Deep Learning-based Methods: Employ models like CNNs and RNNs to detect deepfakes by learning patterns and features specific to synthetic media, with models such as XceptionNet and EfficientNet showing promise (Biswas et al., 2021; Tran et al., 2023).
- (3) Multi-modal Approaches: Combine audio and visual data for enhanced detection accuracy, leveraging correlations between modalities (Yang et al., 2023).
- (4) Explainable AI Methods: Provide insights into detection models' decision-making processes, improving transparency, trust, and reliability (Abir et al., 2023).

### 2.3.2. Challenges and Limitations of Detection Techniques

Despite advancements, deepfake detection faces several challenges:

- (1) Evolving Sophistication: As generation techniques advance, detection methods must adapt to new manipulations, creating an ongoing arms race (Hussain et al., 2022).
- (2) Generalization: Detection models often struggle to generalize across different datasets and techniques, requiring robust and adaptable methods (Coccomini et al., 2023).
- (3) Resource Intensive: Training and deploying detection models demand significant computational resources and large datasets, limiting widespread adoption, especially in resource-constrained environments (Xu et al., 2024).
- (4) False Positives and Negatives: Balancing sensitivity and specificity is challenging, as high false positive rates cause unnecessary alarm while false negatives allow undetected deepfakes, undermining detection systems' effectiveness (Nguyen et al., 2022).

## 2.4. Ethical and Legal Considerations

### 2.4.1. Ethical Implications of Deepfake Usage

Deepfake technology's ethical implications are multifaceted. While it can revolutionize entertainment, education, and creative industries, its misuse poses significant ethical concerns, including:

- (1) Misinformation and Fake News: Deepfakes can create convincing fake news, misleading the public and undermining trust in media and institutions, exacerbating social and political tensions (Hameleers et al., 2024b).

(2) Privacy Violations: Unauthorized deepfakes, especially involving private individuals, constitute severe privacy invasions, causing significant emotional and reputational harm (de Ruiter, 2021).

(3) Identity Theft and Fraud: Deepfakes can facilitate identity theft and financial fraud, posing significant risks to personal and organizational security (Guo et al., 2023).

### **3. Method**

The study involved a systematic review of literature focusing on the intersection of deepfake technology and its use in terrorist propaganda and recruitment. The primary sources included peer-reviewed journal articles, conference papers, and authoritative reports published between 2010 and 2024. A comprehensive search was conducted in databases such as IEEE Xplore, PubMed, Scopus, and Google Scholar using keywords such as "deepfake technology," "terrorist propaganda," "terrorist recruitment," and "deepfake detection." Inclusion criteria were set to select studies that specifically addressed the technological, psychological, and policy aspects of deepfakes in the context of terrorism. Data was collected through a two-stage process. The first stage involved an initial screening of titles and abstracts to identify relevant studies. The second stage involved a full-text review of the selected articles to extract detailed information on the creation, detection, and application of deepfake technology in terrorist activities. Information regarding the impact of deepfakes on public perception, psychological mechanisms of influence, and existing countermeasures was systematically extracted and organized. The extracted data were categorized into themes based on the study objectives: technical mechanisms of deepfakes, psychological impact of deepfake propaganda, historical and contemporary use of deepfakes by terrorist organizations, and countermeasures. A thematic analysis approach was employed to identify patterns and relationships within the data. Key findings were synthesized to provide a comprehensive understanding of the role of deepfakes in terrorist propaganda and recruitment. Given the sensitive nature of the research topic, ethical considerations were prioritized. The review ensured that all secondary data used were appropriately cited and that no personal data were involved. The study adhered to ethical guidelines for conducting systematic reviews, ensuring transparency, and reproducibility of the research process. The study acknowledges potential limitations such as the exclusion of non-English language studies, which may lead to a biased understanding of the global impact of deepfakes. Additionally, the rapidly evolving nature of deepfake technology implies that new developments may emerge that were not covered in the selected literature. This research method outlines a systematic approach to reviewing the literature on deepfake technology's role in terrorist propaganda and recruitment. By integrating technical, psychological, and policy perspectives, the study aims to provide a holistic understanding of the multifaceted threat posed by deepfakes and inform the development of effective countermeasures.

## **4. Result and Discussion**

### **4.1. Deepfake as a Tool for Propaganda**

#### **4.1.1. Historical Use of Propaganda in Terrorism**

Terrorist organizations have historically utilized various propaganda techniques to recruit members, disseminate their ideology, and instill fear within the public. Traditional methods included leaflets, posters, and audio recordings, all of which were straightforward yet effective in delivering messages to a specific audience (Giessmann, 2002). These materials frequently featured graphic imagery and provocative language designed to elicit strong emotional responses and incite violence. Religious and cultural symbolism was also heavily used to resonate deeply with the targeted audience's beliefs and values.

The role of charismatic leaders was critical in these early propaganda efforts. Figures such as Osama bin Laden would deliver speeches and make public appearances that were carefully recorded and distributed to garner support and legitimize the group's actions. These leaders often positioned themselves as defenders of a particular faith or culture, framing their violent activities as a necessary reaction to oppression or injustice.

### 4.1.2. Role of Media and Technology in Propaganda Dissemination

With the advent of mass media and the internet, the scope and sophistication of terrorist propaganda significantly expanded. Television, radio, and print media became vital tools for disseminating messages to a broader audience. However, the rise of digital technology and social media platforms like YouTube, Twitter, and Facebook revolutionized the landscape of terrorist propaganda (Abady et al., 2024). These platforms allowed terrorist organizations to bypass traditional media gatekeepers and directly engage with potential recruits and sympathizers globally.

Social media facilitated the rapid spread of propaganda videos, infographics, and memes aimed at attracting and radicalizing individuals. These platforms also enabled the formation of online communities where like-minded individuals could share content, discuss strategies, and plan attacks. The anonymity provided by the internet allowed recruiters to operate more freely, making it challenging for authorities to trace and disrupt their activities.

Table 1. Historical and Modern Propaganda Techniques

Era	Techniques	Key Features	Examples
Traditional Era	Leaflets, Posters, Audio Recordings	Simple, direct communication, graphic imagery, inflammatory language	Leaflets with religious symbols, audio recordings of speeches
Mass Media Era	Television, Radio, Print Media	Broader reach, use of mass communication channels	TV broadcasts of propaganda, radio shows, newspaper articles
Early Digital Era	Internet, Websites, Email Campaigns	Initial online presence, static web pages, basic digital communication	Terrorist websites, email newsletters
Social Media Era	Social Media Platforms (YouTube, Twitter, Facebook)	Rapid dissemination, direct engagement, anonymity, creation of online communities	Propaganda videos on YouTube, Twitter hashtags
Deepfake Era	Deepfake Technology	Highly realistic synthetic media, exploitation of psychological mechanisms, advanced deception	Deepfake videos of leaders making false statements

The use of traditional propaganda by terrorist organizations has long been a cornerstone of their recruitment and ideological dissemination efforts. Historically, methods such as leaflets, posters, and audio recordings were effective in their simplicity and reach within targeted communities (Giessmann, 2002). These materials were designed to provoke emotional responses and incite action through graphic imagery and inflammatory language. Charismatic leaders played a pivotal role in legitimizing these messages and galvanizing support.

In contrast, the digital age has brought a transformation in the tools and methods of propaganda dissemination. The rise of social media has allowed terrorist organizations to expand their reach exponentially, using platforms to spread their messages more quickly and to a broader audience. The ability

to bypass traditional media gatekeepers has enhanced the effectiveness of propaganda efforts (Abady et al., 2024). Platforms such as YouTube, Twitter, and Facebook enable direct engagement with potential recruits and sympathizers, making it easier to spread radical ideologies and coordinate activities.

The introduction of deepfake technology represents a significant leap in the evolution of terrorist propaganda. Deepfakes enable the creation of highly realistic but fake images, videos, and audio recordings, adding a new layer of complexity and deception. This technology has been used to produce convincing propaganda materials that are difficult to distinguish from genuine content, thereby increasing their persuasive power and impact (Abir et al., 2023). The psychological mechanisms exploited by deepfakes, such as confirmation bias and the illusory truth effect, make them particularly effective tools for manipulating public perception and behavior (Albahar & Almalki, 2019; Ali et al., 2021).

Table 2. Comparison of Traditional and Modern Propaganda Techniques

Aspect	Traditional Techniques	Modern Techniques
Materials	Leaflets, posters, audio recordings	Social media posts, videos, memes, deepfakes
Key Figures	Charismatic leaders	Deepfake videos of influential figures
Dissemination Channels	Leaflets, public speeches, print media	Social media platforms (YouTube, Twitter, Facebook)
Reach	Local communities	Global audience
Psychological Mechanisms	Emotional provocation, religious symbolism	Confirmation bias, illusory truth effect
Impact	Provoked emotional responses, incited violence	Increased trust in false information, sophisticated manipulation

The findings highlight the evolution of terrorist propaganda from traditional methods to sophisticated digital strategies, emphasizing the significant role that media and technology play in modern terrorism. The historical use of propaganda techniques such as leaflets, posters, and audio recordings, while effective in their time, now appears rudimentary compared to the capabilities offered by digital and social media platforms (Giessmann, 2002). The ability of these platforms to facilitate rapid and widespread dissemination of propaganda has fundamentally altered the landscape of terrorist communication and recruitment (Abady et al., 2024).

The emergence of deepfake technology adds a troubling dimension to this landscape. Deepfakes' capacity to create highly realistic synthetic media that can deceive viewers poses a significant threat to public trust in information sources. The psychological impact of deepfake propaganda, which can exploit biases and increase the perceived credibility of false information, underscores the potential for deepfakes to incite violence, disrupt political processes, and undermine societal cohesion (Abdulreda & Obaid, 2022; Albahar & Almalki, 2019; Ali et al., 2021).

Table 3. Implications of Deepfake Technology on Terrorist Propaganda

Implications	Details
Psychological Impact	Exploits confirmation bias, increases perceived credibility of false information
Public Trust	Erodes trust in legitimate news sources and public figures

Social Cohesion	Undermines societal cohesion, leading to increased polarization and radicalization
Political Processes	Potential to disrupt political processes and incite violence
Countermeasures Needed	Advanced detection technologies, legislative measures, public awareness, and education

Recognizing the transformative impact of deepfake technology on terrorist propaganda necessitates a multi-faceted response. Governments and organizations must develop advanced detection technologies and implement legislative measures to mitigate the spread and influence of deepfake propaganda (Ascott, 2020; Asha et al., 2024). Public awareness and education campaigns are essential to build resilience against such manipulative content, teaching individuals to critically evaluate the information they consume (Barabanshchikov & Marinova, 2022; Cafiero, 2023).

## 4.2. Deepfake in Terrorist Recruitment

### 4.2.1. Traditional Recruitment Strategies

Historically, terrorist organizations have relied on direct interpersonal interactions and the exploitation of existing social networks for recruitment. Recruiters often targeted individuals within specific communities, leveraging familial, religious, or ideological connections to establish trust and influence (Giessmann, 2002). Propaganda materials, including pamphlets, videos, and audio recordings, were widely used to spread extremist ideologies and attract new members. These methods aimed to manipulate individuals' beliefs and convince them to join terrorist organizations.

### 4.2.2. Role of Online Platforms in Recruitment

The advent of the internet and social media has significantly expanded terrorist recruitment strategies. Online platforms provide global reach, enabling recruiters to access and influence potential recruits across geographical boundaries. Social media platforms, in particular, have become pivotal in disseminating propaganda, recruitment messages, and the radicalization process. Extremist groups use these platforms to create and distribute content appealing to various psychological and social vulnerabilities (Lakhani, 2023). This shift to online recruitment has increased the speed and efficiency of spreading extremist ideologies and identifying potential recruits.

Table 4. Traditional and Modern Recruitment Strategies

Aspect	Traditional Methods	Modern Methods
Interaction Type	Direct, interpersonal interactions	Online platforms, social media
Targeting Strategy	Leveraging familial, religious, ideological connections	Using online behavior and interests
Materials Used	Pamphlets, videos, audio recordings	Social media posts, videos, memes, deepfakes
Reach	Specific communities	Global audience
Key Features	Establishing trust through direct contact	Personalized and emotionally charged content, anonymity

The integration of deepfake technology into terrorist recruitment represents a significant advancement in their strategies. Traditional recruitment methods involved direct interactions and the use of propaganda materials to manipulate beliefs (Giessmann, 2002). These methods, while effective, were limited by their reliance on physical presence and specific community ties.

In contrast, modern recruitment methods leverage online platforms to reach a global audience rapidly. Social media enables the dissemination of radical content to a broader audience, enhancing the recruitment process's efficiency (Lakhani, 2023). The use of deepfake technology adds another layer of sophistication to these methods. Deepfakes allow for the creation of highly realistic synthetic media, making it challenging for individuals to distinguish between genuine and fake content (Abir et al., 2023).

Deepfakes are particularly effective in manipulating perceptions and inciting action. For instance, deepfake videos of influential religious leaders endorsing extremist ideologies have been used to manipulate and radicalize viewers (Naskar et al., 2024). Additionally, deepfake videos depicting false narratives of government atrocities aim to incite anger and recruit individuals seeking retribution (Ganguly et al., 2022). These sophisticated techniques increase the likelihood of successful radicalization by creating a more personalized and persuasive recruitment experience (Dong et al., 2023).

Table 5. Traditional vs. Modern Recruitment Methods

Aspect	Traditional Methods	Modern Methods with Deepfakes
Interaction Type	Direct, interpersonal interactions	Online platforms, social media with deepfake content
Targeting Strategy	Familial, religious, ideological connections	Online behavior and interests, personalized deepfake content
Materials Used	Pamphlets, videos, audio recordings	Social media posts, deepfake videos
Reach	Specific communities	Global audience, anonymous reach
Key Features	Trust through direct contact	Highly realistic synthetic media, emotional manipulation

The findings underscore the evolution of terrorist recruitment from traditional direct methods to sophisticated digital strategies involving deepfake technology. This shift highlights the increased effectiveness and reach of modern recruitment tactics. The use of online platforms and social media allows terrorist organizations to bypass geographical limitations, reaching a global audience quickly and efficiently (Lakhani, 2023).

The integration of deepfake technology into recruitment processes adds a troubling dimension to these strategies. Deepfakes' ability to create highly realistic synthetic media significantly enhances the persuasive power of recruitment materials. This technology exploits psychological mechanisms such as confirmation bias and the illusory truth effect, making it particularly effective in manipulating perceptions and inciting action (Albahar & Almalki, 2019; Ali et al., 2021).

Governments and international organizations must invest in advanced detection technologies and develop legislative measures to counter the use of deepfakes in terrorist recruitment (Agarwal et al., 2021; Ascott, 2020). Public awareness campaigns are crucial in educating individuals about the existence and risks of deepfake technology, reducing its impact by promoting critical evaluation of information (Barabanshchikov & Marinova, 2022; Cafiero, 2023).

Furthermore, technology plays a vital role in counter-recruitment efforts. Machine learning models for deepfake detection and collaborative efforts between technology companies and law enforcement agencies are essential in tracking and mitigating the spread of malicious content (Gupta et al., 2024). These strategies must be comprehensive, combining technological solutions with public education to safeguard

against the sophisticated manipulative tools used in modern terrorist recruitment.

Table 6. Implications of Deepfake Technology in Terrorist Recruitment

Implications	Details
Psychological Impact	Exploits confirmation bias, increases perceived credibility of false information
Public Trust	Erodes trust in legitimate news sources and public figures
Social Cohesion	Undermines societal cohesion, leading to increased polarization and radicalization
Recruitment Efficiency	Enhances personalization and emotional manipulation, increasing recruitment success
Countermeasures Needed	Advanced detection technologies, legislative measures, public awareness, and education

### 4.3. Case Studies and Real-World Examples

#### 4.3.1. Notable Incidents of Deepfake Propaganda

Deepfake technology has been utilized in various propaganda campaigns, significantly impacting the socio-political landscape. One notable incident involved a deepfake video of Nancy Pelosi, where her speech was altered to appear as if she were intoxicated. This video, widely circulated on social media, led to discussions on the authenticity of political communications and the potential for misinformation (Hameleers et al., 2024a). Another significant case was the creation of a deepfake video depicting Ukrainian President Volodymyr Zelenskyy urging his troops to surrender. This video, though quickly debunked, demonstrated the potential of deepfakes to influence public opinion and destabilize governments (de Ruiter, 2021).

The impact of these deepfake propaganda instances is substantial, ranging from undermining public trust in media to creating political unrest. The Pelosi video, for instance, fueled partisan divides and heightened skepticism towards news outlets. The Zelenskyy incident, while not successful in its immediate goal, highlighted the vulnerability of societies to digital deception, necessitating advancements in deepfake detection and media literacy (Hameleers et al., 2024b).

Table 7. Notable Incidents of Deepfake Propaganda

Incident	Description	Impact
Nancy Pelosi Deepfake	Altered video to appear intoxicated	Discussions on political communication authenticity, increased mistrust in media
Volodymyr Zelenskyy Deepfake	Fake video urging Ukrainian troops to surrender	Highlighted societal vulnerability to digital deception, emphasized need for detection tech

### 4.3.2. Examples of Deepfake Recruitment Campaigns

Terrorist organizations have exploited deepfake technology in recruitment campaigns. The Islamic State (ISIS), for instance, has used deepfake videos to create persuasive propaganda that glorifies their cause and depicts a distorted reality of life under their control. These videos often feature manipulated speeches from influential figures or fabricate endorsements from well-known personalities, increasing their appeal and credibility (Albahar & Almalki, 2019).

The success of these campaigns varies but has notable consequences. Deepfake technology enhances the sophistication of these videos, making them more convincing and difficult to debunk. This increases the likelihood of recruitment, as potential recruits are more easily swayed by seemingly credible sources. The consequences extend beyond immediate recruitment, as these videos also serve to radicalize viewers and strengthen ideological adherence among existing members (Giessmann, 2002).

Table 8. Deepfake Recruitment Campaigns

Organization	Deepfake Usage	Impact
ISIS	Glorification of cause, distorted reality of life under their control	Increased recruitment likelihood, difficulty in debunking, radicalization of viewers
General Trend	Manipulated speeches, fabricated endorsements	Enhanced sophistication and credibility, extended influence beyond immediate recruitment

### 4.3.3. Lessons Learned from Case Studies

The examination of deepfake incidents reveals several key takeaways. Firstly, the rapid dissemination of deepfakes underscores the need for robust detection technologies and swift response mechanisms. Enhancements in deepfake detection methods, such as those employing siamese-based verification systems, are critical (Abady et al., 2024). Secondly, the societal impact of deepfakes, particularly in propaganda and recruitment, highlights the importance of media literacy programs aimed at educating the public on identifying manipulated content (Brashier, 2024).

Furthermore, the legal and ethical frameworks surrounding deepfakes need to evolve. Current regulations often lag behind technological advancements, necessitating updates to effectively address the challenges posed by deepfakes. Policymakers must consider the implications of deepfakes on privacy, security, and public trust, and develop comprehensive strategies to mitigate these risks (de Ruiter, 2021).

Table 9. Lessons Learned from Deepfake Case Studies

Lesson	Details
Need for Detection Technologies	Robust systems like siamese-based verification are essential
Importance of Media Literacy	Programs to educate the public on identifying manipulated content

Evolution of Legal Frameworks

Regulations must adapt to address privacy, security, and public trust issues

Policy Development

Comprehensive strategies needed to mitigate deepfake risks

#### 4.3.4. Comparative Analysis

A comparative analysis of deepfake usage across different terrorist organizations reveals distinct trends and patterns. ISIS, for instance, employs deepfakes extensively in their recruitment videos, focusing on creating an illusion of legitimacy and widespread support. Other groups, such as Al-Qaeda, have also begun to explore deepfake technology, though their use is less sophisticated and widespread (Albahar & Almalki, 2019).

Trends indicate a growing adoption of deepfake technology among terrorist organizations, driven by its effectiveness in enhancing the credibility of propaganda. The patterns suggest that as detection technologies improve, these groups are likely to refine their techniques, making deepfakes harder to identify. This ongoing technological arms race necessitates continuous advancements in detection methods and international cooperation to address the global threat posed by deepfakes (Abbas et al., 2023; Abdullah & Ali, 2023). The comparative analysis also underscores the importance of understanding the unique tactics employed by different organizations. While ISIS focuses on recruitment and radicalization, other groups might leverage deepfakes for disinformation campaigns aimed at destabilizing regions or undermining political opponents. Recognizing these patterns allows for the development of targeted countermeasures, enhancing the overall effectiveness of anti-terrorism strategies (Abdulreda & Obaid, 2022).

Table 10. Comparative Analysis of Deepfake Usage

Organization	Focus	Trends and Patterns
ISIS	Recruitment and radicalization	Extensive use of deepfakes to enhance credibility and legitimacy
Al-Qaeda	Emerging use	Less sophisticated and widespread use of deepfake technology
General Trends	Enhanced credibility of propaganda	Growing adoption among terrorist groups, ongoing refinement of techniques

These findings underscore the critical need for robust technological solutions and comprehensive public education to combat the sophisticated threat posed by deepfakes.

#### 4.4. Technological and Psychological Countermeasures

##### 4.4.1. Technological Solutions to Combat Deepfakes

Advances in detection technology have significantly contributed to combating the proliferation of deepfakes. Various detection methods have emerged as robust countermeasures. Recent developments include the use of machine learning models such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs), which have been effective in identifying deepfakes by analyzing inconsistencies in image and video data (Abdullah & Ali, 2023). Specifically, the use of Local Binary Pattern (LBP) in multiple-channel color spaces has shown improvements in image-based deepfake detection. Siamese-based verification systems for open-set architecture attribution of synthetic images enhance detection accuracy by comparing pairs of images to discern real from synthetic content (Abady et al., 2024).

Additionally, ensemble learning approaches have been utilized to reduce dataset specificity, thereby improving the generalization and robustness of detection systems across various datasets (Abbas et al., 2023).

Artificial Intelligence (AI) plays a pivotal role in developing sophisticated countermeasures against deepfakes. AI-driven systems like Vision Transformer Networks have empowered deepfake detection by analyzing complex visual patterns and discrepancies (Arshed et al., 2023). These networks leverage the power of attention mechanisms to focus on crucial image features that are often manipulated in deepfakes. Explainable AI (XAI) methods are also gaining traction as they provide transparency in the decision-making process of AI models, making it easier to understand and trust the detection results (Abir et al., 2023). Deep learning techniques, such as the use of deep dual-level networks, have demonstrated effectiveness in detecting deepfakes by analyzing both the temporal coherence and spatial features of videos (Pu et al., 2022).

Table 11. Technological Solutions for Deepfake Detection

Technology	Description	Impact
Convolutional Neural Networks (CNNs)	Machine learning models analyzing inconsistencies in image/video data	Improved accuracy in identifying deepfakes
Local Binary Pattern (LBP)	Analyzes multiple-channel color spaces for image-based detection	Enhanced image-based deepfake detection
Siamese-based Verification Systems	Compares pairs of images to discern real from synthetic content	Increased detection accuracy
Vision Transformer Networks	AI systems focusing on crucial image features using attention mechanisms	Empowered deepfake detection through complex visual pattern analysis
Explainable AI (XAI)	Provides transparency in AI decision-making processes	Improved understanding and trust in detection results
Deep Dual-Level Networks	Analyzes temporal coherence and spatial features of videos	Effective in detecting frame sequence inconsistencies in deepfakes

#### 4.4.2. Psychological Resilience Building

Building psychological resilience among the public is crucial in mitigating the impact of deepfakes. Strategies include increasing awareness about the existence and potential threats of deepfakes through public education campaigns. Informing individuals about the signs and characteristics of deepfake content can make them more critical of the media they consume (Newman & Schwarz, 2024). Developing cognitive tools that help people critically assess the veracity of information is another key strategy. This can be achieved by integrating media literacy programs into educational curriculums, teaching students how to evaluate sources and cross-check information. These programs can be tailored to different age groups and learning levels to maximize their effectiveness.

Education plays a fundamental role in fostering critical thinking skills essential for identifying and resisting deepfake content. Integrating AI literacy into school curricula can prepare students to navigate a digital landscape increasingly dominated by AI-generated content. This involves teaching the basics of AI,

its applications, and its potential for misuse in creating deepfakes (Kong et al., 2024). Additionally, fostering an environment that encourages questioning and skepticism towards digital content can help individuals develop a critical mindset. Workshops and seminars that simulate the creation and detection of deepfakes can provide practical experience, enhancing individuals' ability to identify manipulated media (Lee & Park, 2023).

Table 12. Psychological Resilience Building Strategies

Strategy	Description	Impact
Public Education Campaigns	Increasing awareness about deepfakes and their potential threats	Greater public awareness and critical assessment of media
Media Literacy Programs	Teaching students to evaluate sources and cross-check information	Enhanced critical thinking skills and media literacy
AI Literacy in Curricula	Educating students on AI, its applications, and misuse	Preparedness for navigating AI-generated content
Workshops and Seminars	Practical experience in creating and detecting deepfakes	Improved ability to identify manipulated media
Encouraging Questioning and Skepticism	Fostering a critical mindset towards digital content	Greater resilience against deepfake influence

#### 4.4.3. Policy and Regulatory Approaches

Current policies on deepfake regulation vary significantly across different jurisdictions. Some countries have implemented stringent laws aimed at criminalizing the production and distribution of malicious deepfakes, while others are still in the process of developing comprehensive regulatory frameworks. For instance, the European Union's General Data Protection Regulation (GDPR) has provisions that indirectly address the issues posed by deepfakes, particularly in terms of data privacy and consent (van der Sloot & Wagenveld, 2022). In the United States, certain states have enacted laws specifically targeting deepfakes used for malicious purposes, such as non-consensual pornography and electoral interference. However, there is a need for federal-level legislation to provide a uniform approach to tackling this issue (Kirchengast, 2020).

To enhance the effectiveness of policies combating deepfakes, several recommendations can be made. Firstly, there should be a harmonization of laws across jurisdictions to ensure consistency and avoid loopholes that perpetrators could exploit. International cooperation is crucial in this aspect, as deepfakes often transcend national borders (Cafiero, 2023). Secondly, regulatory bodies should consider the rapid evolution of deepfake technology and adopt a proactive approach in updating policies to keep pace with technological advancements. This includes setting up dedicated task forces to monitor and research emerging trends in deepfake technology (de Ruiter, 2021). Furthermore, policies should mandate the development and deployment of technical solutions by online platforms to detect and remove deepfake content promptly. Platforms should also be required to provide users with tools to report and verify the authenticity of media content (Xu et al., 2024).

#### 4.4.4. International Collaboration

Global cooperation is essential in combating the widespread issue of deepfakes. The transnational nature of digital media means that deepfake content can quickly spread across borders, making it a global concern. International collaboration can facilitate the sharing of knowledge, resources, and best practices, thereby enhancing the overall effectiveness of countermeasures (Amaizu et al., 2024). Collaborative efforts can also lead to the development of standardized detection tools and protocols, which can be used worldwide to identify and mitigate deepfake threats. Joint research initiatives and international conferences can provide platforms for experts to exchange ideas and innovations in deepfake detection and prevention (Casu et al., 2024b).

Several international initiatives have demonstrated success in addressing deepfake challenges. The Global Partnership on Artificial Intelligence (GPAI), for instance, brings together experts from various countries to collaborate on AI ethics and governance, including the ethical implications of deepfake technology. This partnership aims to promote the responsible development and use of AI technologies globally (Cowles et al., 2024). Another example is INTERPOL's efforts to combat digital forgery, which includes training law enforcement agencies worldwide on advanced detection techniques and promoting international legal frameworks to address cybercrime involving deepfakes (Giessmann, 2002).

Table 13. Policy and Regulatory Approaches

Policy/Approach	Description	Impact
Harmonization of Laws	Ensuring consistency across jurisdictions to avoid legal loopholes	Effective and uniform regulation
Proactive Policy Updates	Adopting a proactive approach to keep pace with technological advancements	Timely and relevant regulation
Mandating Technical Solutions	Requiring platforms to detect and remove deepfake content	Enhanced detection and mitigation of deepfake content
International Cooperation	Facilitating global sharing of knowledge, resources, and best practices	Strengthened global countermeasures
Global Partnership on AI (GPAI)	Collaboration on AI ethics and governance, including deepfake implications	Responsible development and use of AI technologies globally
INTERPOL Training	Training law enforcement on advanced detection techniques	Improved global capacity to combat digital forgery and cybercrime

In conclusion, technological and psychological countermeasures against deepfakes require a multifaceted approach that includes advancements in detection technology, building public resilience, implementing robust policies, and fostering international cooperation. By leveraging these strategies, it is possible to mitigate the risks posed by deepfakes and protect the integrity of information in the digital age.

## 5. Conclusion

In this comprehensive review, we have thoroughly examined the multifaceted implications of deepfake technology, particularly focusing on its role in terrorist propaganda and recruitment, as well as its broader societal impacts and the necessary countermeasures. The key findings highlight significant technological advancements, including techniques such as Generative Adversarial Networks (GANs), which enable the creation of highly realistic synthetic media. These advancements present both opportunities and threats across various domains, including entertainment, security, and misinformation. Despite the development of various detection methodologies, such as those focusing on spatial-temporal inconsistencies and machine learning techniques, the ongoing arms race between deepfake creation and detection underscores the need for continuous research and development.

Deepfakes have profound implications for trust and authenticity in digital communications, potentially eroding public trust in media and institutions and exacerbating misinformation and manipulation. Additionally, regulatory frameworks and ethical guidelines are still catching up with the technology, necessitating a balance between protecting individuals' rights and fostering innovation. Technological solutions to combat deepfakes include advancements in AI and machine learning, such as convolutional neural networks (CNNs), Local Binary Pattern (LBP), and Vision Transformer Networks, which are crucial in detecting and combating deepfakes. Explainable AI (XAI) methods also enhance transparency and trust in these technologies.

Building public resilience through education and media literacy programs is essential for mitigating the influence of deepfakes. Teaching critical thinking and AI literacy can empower individuals to better navigate and assess digital content, reducing susceptibility to manipulation. Policy and regulatory approaches need to evolve in tandem with technological advancements. Some jurisdictions have implemented stringent laws against malicious deepfakes, but there is a need for harmonized global regulations. The EU's GDPR and certain state laws in the U.S. are steps in this direction. Policymakers should develop comprehensive frameworks, set up dedicated task forces, and mandate the development of technical solutions by online platforms to detect and remove deepfake content promptly.

International collaboration is crucial in combating deepfakes. Sharing knowledge, resources, and best practices can enhance the overall effectiveness of countermeasures. Initiatives like the Global Partnership on Artificial Intelligence (GPAI) and INTERPOL's efforts exemplify successful international cooperation. Future deepfake technology will be shaped by advancements in AI, machine learning, and neural networks. The integration of blockchain and collective intelligence models presents a promising avenue for enhancing detection and prevention. Understanding the evolution of deepfake technology involves analyzing current trends and capabilities, with predictive models and detection algorithms being critical in preemptively identifying potential threats.

Despite significant strides in research, several gaps remain. There is a need for improved detection methods that can handle diverse datasets and varied techniques. Research should also focus on audio deepfakes and the integration of Explainable AI to enhance transparency and trust in detection systems. Ethical considerations of deepfake technology involve issues related to consent, privacy, and authenticity. Addressing these challenges requires robust regulatory frameworks, public education, and fostering a culture of digital literacy and critical thinking.

In conclusion, while deepfake technology presents significant challenges, a proactive and collaborative approach can mitigate their adverse effects and harness their potential for positive impact. Continued vigilance, innovation, and cooperation are essential to navigate the complexities of this evolving technological landscape. By addressing the research gaps, ethical considerations, and leveraging emerging technologies, we can develop robust systems to detect and counteract the misuse of deepfakes, ensuring a safer and more trustworthy digital environment.

## References

- Abady, L., Wang, J., Tondi, B., & Barni, M. (2024). A siamese-based verification system for open-set architecture
- Abbas, Q., Alghamdi, T., Alsaawy, Y., Alyas, T., Alzahrani, A., Malik, K. I., & Bibi, S. (2023). Reducing Dataset Specificity for Deepfakes Using Ensemble Learning. *Computers, Materials and Continua*, 74(2), 4261–4276.
- Abdullah, M. T., & Ali, N. H. M. (2023). DeepFake Detection Improvement for Images Based on a Proposed Method for Local Binary Pattern of the Multiple-Channel Color Space. *International Journal of Intelligent Engineering and Systems*, 16(3), 92–104.
- Abdulreda, A. S., & Obaid, A. J. (2022). A landscape view of deepfake techniques and detection methods. *International Journal of Nonlinear Analysis and Applications*, 13(1), 745–755.
- Abir, W. H., Khanam, F. R., Alam, K. N., Hadjouni, M., Elmannai, H., Bourouis, S., Dey, R., & Khan, M. M. (2023). Detecting Deepfake Images Using Deep Learning Techniques and Explainable AI Methods. *Intelligent Automation and Soft Computing*, 35(2), 2151–2169.
- Agarwal, A., Singh, R., Vatsa, M., & Noore, A. (2021). MagNet: Detecting Digital Presentation Attacks on Face Recognition. *Frontiers in Artificial Intelligence*, 4.
- Albahar, M., & Almalki, J. (2019). Deepfakes: Threats and countermeasures systematic review. *Journal of Theoretical and Applied Information Technology*, 97(22), 3242–3250.
- Ali, S., DiPaola, D., Lee, I., Sindato, V., Kim, G., Blumofe, R., & Breazeal, C. (2021). Children as creators, thinkers and citizens in an AI-driven future. *Computers and Education: Artificial Intelligence*, 2.
- Amaizu, G. C., Njoku, J. N., Lee, J. M., & Kim, D. S. (2024). Metaverse in advanced manufacturing: Background, applications, limitations, open issues & future directions. In *ICT Express* (Vol. 10, Issue 2, pp. 233–255). Korean Institute of Communications and Information Sciences.
- Amerini, I., Anagnostopoulos, A., Maiano, L., & Celsi, L. R. (2021). Deep learning for multimedia forensics. *Foundations and Trends in Computer Graphics and Vision*, 12(4), 309–457.
- Amin, M. A., Hu, Y., & Hu, J. (2024). Analyzing temporal coherence for deepfake video detection. *Electronic Research Archive*, 32(4), 2621–2641.
- Amin, M. A., Hu, Y., Li, C.-T., & Liu, B. (2024). Deepfake detection based on cross-domain local characteristic analysis with multi-domain transformer. *Alexandria Engineering Journal*, 91, 592–609.
- Appel, M., & Prietzel, F. (2022). The detection of political deepfakes. *Journal of Computer-Mediated Communication*, 27(4). <https://doi.org/10.1093/jcmc/zmac008>
- Arshed, M. A., Alwadain, A., Faizan Ali, R., Mumtaz, S., Ibrahim, M., & Muneer, A. (2023). Unmasking Deception: Empowering Deepfake Detection with Vision Transformer Network. *Mathematics*, 11(17).
- Ascott, T. (2020). Microfake: How small-scale deepfakes can undermine society. *Journal of Digital Media and Policy*,
- Asha, S., Vinod, P., Amerini, I., & Menon, V. G. (2024). D-Fence layer: an ensemble framework for comprehensive deepfake detection. *Multimedia Tools and Applications*.
- Barabanshchikov, V. A., & Marinova, M. M. (2022). Deepfake as the basis for digitally collaging “impossible faces.” *Journal of Optical Technology (A Translation of Opticheskii Zhurnal)*, 89(8), 448–453.
- Biswas, A., Bhattacharya, D., & Kumar, K. A. (2021). DeepFake Detection using 3D-Xception Net with Discrete Fourier Transformation. *Journal of Information Systems and Telecommunication*, 9(35), 161–168.
- Brashier, N. M. (2024). Fighting misinformation among the most vulnerable users. In *Current Opinion in Psychology* (Vol. 57). Elsevier B.V.
- Cafiero, F. (2023). Datafying diplomacy: How to enable the computational analysis and support of international negotiations. *Journal of Computational Science*, 71.
- Caldelli, R., Galteri, L., Amerini, I., & Del Bimbo, A. (2021). Optical Flow based CNN for detection of unlearned deepfake manipulations. *Pattern Recognition Letters*, 146, 31–37.
- Casu, M., Guarnera, L., Caponnetto, P., & Battiato, S. (2024a). GenAI mirage: The impostor bias and the deepfake detection challenge in the era of artificial illusions. In *Forensic Science International: Digital Investigation* (Vol. 50). Elsevier Ltd.

- Casu, M., Guarnera, L., Caponnetto, P., & Battiato, S. (2024b). GenAI mirage: The impostor bias and the deepfake detection challenge in the era of artificial illusions. In *Forensic Science International: Digital Investigation* (Vol. 50). Elsevier Ltd.
- Chen, B., Li, T., & Ding, W. (2022). Detecting deepfake videos based on spatiotemporal attention and convolutional LSTM. *Information Sciences*, 601, 58–70.
- Chen, G.-L., & Hsu, C.-C. (2023). Jointly Defending DeepFake Manipulation and Adversarial Attack Using Decoy Mechanism. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(8), 9922–9931.
- Chen, H., Li, Y., Lin, D., Li, B., & Wu, J. (2023). Watching the BiG artifacts: Exposing DeepFake videos via Bi-granularity artifacts. *Pattern Recognition*, 135.
- Chen, M., Liao, X., & Wu, M. (2022). PulseEdit: Editing Physiological Signals in Facial Videos for Privacy Protection. *IEEE Transactions on Information Forensics and Security*, 17, 457–471.
- Chhabra, S., Thakral, K., Mittal, S., Vatsa, M., & Singh, R. (2024). Low-Quality Deepfake Detection via Unseen Artifacts. *IEEE Transactions on Artificial Intelligence*, 5(4), 1573–1585.
- Coccomini, D. A., Caldelli, R., Falchi, F., & Gennaro, C. (2023). On the Generalization of Deep Learning Models in Video Deepfake Detection. *Journal of Imaging*, 9(5).
- Cowles, K., Miller, R., & Suppok, R. (2024). When Seeing Isn't Believing: Navigating Visual Health Misinformation through Library Instruction. *Medical Reference Services Quarterly*, 43(1), 44–58.
- de Ruiter, A. (2021). The Distinct Wrong of Deepfakes. *Philosophy and Technology*, 34(4), 1311–1332.
- Ding, F., Zhu, G., Li, Y., Zhang, X., Atrey, P. K., & Lyu, S. (2022). Anti-Forensics for Face Swapping Videos via Adversarial Training. *IEEE Transactions on Multimedia*, 24, 3429–3441.
- Dong, J., Wang, Y., Lai, J., & Xie, X. (2023). Restricted Black-Box Adversarial Attack Against DeepFake Face Swapping. *IEEE Transactions on Information Forensics and Security*, 18, 2596–2608.
- Ferreira, S., Antunes, M., & Correia, M. E. (2021a). A dataset of photos and videos for digital forensics analysis using machine learning processing. *Data*, 6(8).
- Ferreira, S., Antunes, M., & Correia, M. E. (2021b). Exposing manipulated photos and videos in digital forensics analysis. *Journal of Imaging*, 7(7).
- Fu, H., Gong, M., Wang, C., Batmanghelich, K., Zhang, K., & Tao, D. (2019). Geometry-Consistent Generative Adversarial Networks for One-Sided Unsupervised Domain Mapping. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2422–2431.
- Ganguly, S., Ganguly, A., Mohiuddin, S., Malakar, S., & Sarkar, R. (2022). ViXNet: Vision Transformer with Xception Network for deepfakes based video and image forgery detection. *Expert Systems with Applications*, 210.
- Giessmann, H. J. (2002). Media and the Public Sphere: Catalyst and Multiplier of Terrorism? *Media Asia*, 29(3), 134–136.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative Adversarial Networks. *Science Robotics*, 3(January), 2672–2680.
- Guo, J., Zhao, Y., & Wang, H. (2023). Generalized Spoof Detection and Incremental Algorithm Recognition for Voice Spoofing. *Applied Sciences (Switzerland)*, 13(13).
- Gupta, P., Ding, B., Guan, C., & Ding, D. (2024). Generative AI: A systematic review using topic modelling techniques. *Data and Information Management*, 8(2).
- Hameleers, M., van der Meer, T. G. L. A., & Dobber, T. (2024a). Distorting the truth versus blatant lies: The effects of different degrees of deception in domestic and foreign political deepfakes. *Computers in Human Behavior*, 152.
- Hameleers, M., van der Meer, T. G. L. A., & Dobber, T. (2024b). They Would Never Say Anything Like This! Reasons To Doubt Political Deepfakes. *European Journal of Communication*, 39(1), 56–70.
- Harbinja, E., Edwards, L., & McVey, M. (2023). Governing ghostbots. *Computer Law and Security Review*, 48.
- Hussain, S., Neekhara, P., Dolhansky, B., Bitton, J., Ferrer, C. C., Mcauley, J., & Koushanfar, F. (2022). Exposing Vulnerabilities of Deepfake Detection Systems with Robust Attacks. *Digital Threats: Research and Practice*, 3(3).
- Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., & Aila, T. (2020). Analyzing and Improving the Image Quality of StyleGAN. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 8107–8116.

- Kietzmann, J., Lee, L. W., McCarthy, I. P., & Kietzmann, T. C. (2020). Deepfakes: Trick or treat? *Business Horizons*, 63(2), 135–146.
- Kirchengast, T. (2020). Deepfakes and image manipulation: criminalisation and control. *Information and Communications Technology Law*, 308–323.
- Kong, S. C., Cheung, M. Y. W., & Tsang, O. (2024). Developing an artificial intelligence literacy framework: Evaluation of a literacy course for senior secondary students using a project-based learning approach. *Computers and Education: Artificial Intelligence*, 6.
- Lakhani, S. (2023). When Digital and Physical World Combine: The Metaverse and Gamification of Violent Extremism. *Perspectives on Terrorism*, 17(2), 108–125.
- Lee, J., & Park, J. (2023). AI as “Another I”: Journey map of working with artificial intelligence from AI-phobia to AI-preparedness. *Organizational Dynamics*, 52(3).
- Liu, H., Zhou, W., Chen, D., Fang, H., Bian, H., Liu, K., Zhang, W., & Yu, N. (2023). Coherent adversarial deepfake video generation. *Signal Processing*, 203.
- Naskar, G., Mohiuddin, S., Malakar, S., Cuevas, E., & Sarkar, R. (2024). Deepfake detection using deep feature stacking and meta-learning. *Heliyon*, 10(4).
- Newman, E. J., & Schwarz, N. (2024). Misinformed by images: How images influence perceptions of truth and what can be done about it. In *Current Opinion in Psychology* (Vol. 56). Elsevier B.V.
- Nguyen, T. T., Nguyen, Q. V. H., Nguyen, D. T., Nguyen, D. T., Huynh-The, T., Nahavandi, S., Nguyen, T. T., Pham, Q.-V., & Nguyen, C. M. (2022). Deep learning for deepfakes creation and detection: A survey. *Computer Vision and Image Understanding*, 223.
- Pu, W., Hu, J., Wang, X., Li, Y., Hu, S., Zhu, B., Song, R., Song, Q., Wu, X., & Lyu, S. (2022). Learning a deep dual-level network for robust DeepFake detection. *Pattern Recognition*, 130.
- Tran, V.-N., Kwon, S.-G., Lee, S.-H., Le, H.-S., & Kwon, K.-R. (2023). Generalization of Forgery Detection With Meta Deepfake Detection Model. *IEEE Access*, 11, 535–546.
- Van der Sloot, B., & Wagenveld, Y. (2022). Deepfakes: regulatory challenges for the synthetic society. *Computer Law and Security Review*, 46.
- Vizoso, Á., Vaz-álvarez, M., & López-García, X. (2021). Fighting deepfakes: Media and internet giants’ converging and diverging strategies against hi-tech misinformation. *Media and Communication*, 9(1), 291–300.
- Xu, P., Ma, Z., Mei, X., & Shen, J. (2024). Detecting facial manipulated images via one-class domain generalization. *Multimedia Systems*, 30(1).
- Yang, W., Zhou, X., Chen, Z., Guo, B., Ba, Z., Xia, Z., Cao, X., & Ren, K. (2023). AVoid-DF: Audio-Visual Joint Learning for Detecting Deepfake. *IEEE Transactions on Information Forensics and Security*, 18, 2015–2029.