



INTELLIGENCE BRIEFING

DEEPPFAKE IMAGES
VIDEO & AUDIO

A: Pitt Street, Sydney, NSW,
2000, Australia

P: +61 2 9188 7896
E: contact@cybertrace.com.au

UNCLASSIFIED / OSINT

DEEPPAKE IMAGES, VIDEO & AUDIO

Classification: Unclassified
Distribution: Public
Source: OSINT

Executive Summary:

This intelligence briefing provides an overview of the article, "Investigating Deepfake Images, Video, and Audio" published by Cybertrace, as well as additional insights from various OSINT sources. This briefing highlights the rising threat of deepfake technology and its implications for different sectors. It outlines the characteristics of deepfake media, discusses the challenges faced during the investigation process, and suggests potential solutions for identifying and combating deepfakes.

Key Findings:

1. Deepfake technology has rapidly advanced, enabling the creation of highly realistic and convincing fake images, videos, and audio recordings.
2. Deepfakes can be exploited for malicious purposes, including the spread of disinformation, fraud, blackmail, and reputational damage.
3. Investigating deepfakes presents unique challenges due to their sophisticated techniques, making them difficult to detect, attribute, and debunk.
4. Traditional forensic techniques are often insufficient to identify deepfakes, requiring specialised tools, algorithms, and expertise.
5. Technical approaches to deepfake detection include analysing artifacts, inconsistencies, or anomalies, leveraging AI algorithms, and utilising blockchain technology.
6. Human-based approaches, such as visual analysis by trained experts, can complement technical methods in deepfake identification.
7. Collaboration among law enforcement agencies, technology companies, academia, and other stakeholders is crucial for developing effective countermeasures against deepfakes.

UNCLASSIFIED / OSINT

UNCLASSIFIED / OSINT



A deepfake example that was released to coincide with the expected arrest of Donald Trump in March 2023. The images were reportedly created via OpenAI's DALL-E which is the image version of ChatGPT. Source: <https://www.thetimes.co.uk/article/donald-trump-deepfakes-ai-twitter-g50n7vnbm>.

UNCLASSIFIED / OSINT

UNCLASSIFIED / OSINT



Another Donald Trump deepfake example which at the time of writing had over 1.1 million views on Twitter. Source: <https://twitter.com>.

UNCLASSIFIED / OSINT

UNCLASSIFIED / OSINT

Analysis:

Deepfake technology has rapidly evolved, presenting a significant challenge to the authenticity and trustworthiness of media content. Deepfakes have the potential to impact various sectors, including politics, military journalism, entertainment, and finance. The widespread use of deepfakes for malicious purposes poses a serious threat to public trust, social stability, and national security.

Investigating deepfakes requires a multidimensional approach due to their sophisticated nature. The article highlights the limitations of traditional forensic techniques, as deepfakes are designed to pass visual inspection and avoid easy detection. However, additional sources provide insights into promising methods and technologies to combat deepfakes.

According to a report by the Center for Security and Emerging Technology (CSET), technical approaches to deepfake detection involve analysing artifacts, inconsistencies, or anomalies within the media. AI algorithms can be trained to recognise patterns and identify tell-tale signs of manipulation. Furthermore, blockchain technology can be utilised to create a tamper-proof record of media content, ensuring its authenticity, and preventing unauthorised modifications.

In addition to technical approaches, human-based methods play a crucial role in deepfake investigation. The CSET report suggests that visual analysis by trained experts can help identify subtle visual cues that indicate deepfake manipulation. Human involvement is essential in distinguishing nuanced details that may evade automated algorithms.

Collaboration among stakeholders is imperative in the fight against deepfakes. Law enforcement agencies, military, technology companies, academia, and industry experts must work together to share knowledge, resources, and best practices. Initiatives like the Deepfake Detection Challenge, a collaboration between tech giants and academia, have demonstrated the importance of collective efforts in advancing deepfake detection capabilities.

UNCLASSIFIED / OSINT

UNCLASSIFIED / OSINT

It is imperative to investigate deepfake images, videos, and audio considering the growing prevalence and accessibility of deepfake technology. Exploiting the competitive dynamics of Generative Adversarial Networks (GANs), deepfakes can deceive not only machine learning models but also human perception. While deepfakes were once limited to proficient hackers and computer enthusiasts with high-performance hardware, their creation has become easily accessible to anyone with a basic laptop and internet connection. Moreover, the advent of AI voice cloning enables the manipulation of vocal content using a mere three-second audio clip, thereby exacerbating the potential for misuse.

Generative Adversarial Networks (GANs):

Generative Adversarial Networks (GANs) are a class of machine learning models that consist of two components: a generator and a discriminator. GANs are designed to generate synthetic data that closely resembles real data by leveraging a competitive process between the two components.

The generator component of a GAN is responsible for creating synthetic data samples. It takes random noise as input and generates data, such as images, videos, or audio, based on patterns it has learned from training on real data. The goal of the generator is to produce samples that are indistinguishable from real data.

The discriminator component, on the other hand, acts as a binary classifier that distinguishes between real and synthetic data. It is trained using both real and generated samples, learning to differentiate between the two. The discriminator's objective is to correctly classify the origin of the data, determining whether it is real or generated by the generator.

During training, the generator and discriminator are pitted against each other in a competitive fashion. The generator aims to generate data that can fool the discriminator into classifying it as real, while the discriminator strives to accurately discriminate between real and fake samples. This adversarial training process drives both components to improve their performance over time.

UNCLASSIFIED / OSINT

UNCLASSIFIED / OSINT

As the training progresses, the generator becomes more adept at generating realistic samples that are difficult for the discriminator to discern as synthetic. Conversely, the discriminator becomes more skilled at distinguishing between real and generated data. This back-and-forth training dynamic results in the generator continually refining its ability to produce high-quality synthetic samples.

Overall, GANs are powerful models for generating synthetic data that closely resembles real data. They have found applications in various domains, including image synthesis, video generation, text generation, and audio synthesis. However, the same underlying principles that make GANs effective for generating realistic data also contribute to their potential misuse in creating deceptive deepfakes.

Recommendations:

1. **Invest in Research and Development:** Allocate resources to advance research and development efforts in deepfake detection and attribution methods. Foster collaboration between academic institutions, technology companies, and law enforcement agencies to expedite progress in this field.
2. **Foster Public-Private Partnerships:** Establish partnerships with technology companies, social media platforms, and content creators to facilitate the detection and removal of deepfake content. Encourage information sharing, joint research, and development of effective countermeasures against deepfakes.
3. **Train Law Enforcement and Judicial Personnel:** Provide specialised training to law enforcement and judicial personnel on deepfake investigation techniques, forensic analysis, and legal considerations. Equip them with the necessary tools and knowledge to identify, investigate, and prosecute deepfake-related crimes.
4. **Develop Legislative Frameworks:** Collaborate with policymakers to develop legislation and regulations that address deepfake-related crimes. Establish legal frameworks that deter the creation, distribution, and malicious use of deepfake content.
5. **Promote Media Literacy and Awareness:** Launch public awareness campaigns to educate individuals about the existence and risks associated with deepfakes.

UNCLASSIFIED / OSINT

UNCLASSIFIED / OSINT

Encourage media literacy, critical thinking, and responsible consumption of digital content.

6. Establish International Cooperation: Foster international cooperation and information sharing among governments, law enforcement agencies, and cybersecurity organisations to effectively combat cross-border deepfake threats.

Summary:

Deepfake technology presents a significant and evolving threat to the authenticity and trustworthiness of media content. Investigating and countering deepfakes require a comprehensive approach that combines technical methods, human expertise, and collaboration among stakeholders. By investing in research and development, fostering partnerships, and raising public awareness, it is possible to mitigate the risks posed by deepfakes and safeguard the integrity of media in the digital age.

Please note that the information provided in this briefing includes insights from the article "Investigating Deepfake Images, Video, and Audio" by Cybetrace and additional sources. It is essential to further verify and corroborate the information through multiple intelligence sources to ensure accuracy and validity.

DISCLAIMER

This report's contents are offered by Cybetrace Pty Ltd. solely for informational purposes and should not be construed as advice or justification for any action taken based on the information and analysis provided by Cybetrace Pty Ltd. The information and analysis presented in this report are based only on data available at the time of drafting. Without the express written consent of Cybetrace Pty Ltd, this report and its contents may not be copied, reproduced, or distributed further. All rights Reserved Cybetrace 2023.

UNCLASSIFIED / OSINT

UNCLASSIFIED / OSINT

ENQUIRIES

For sales, further information, or feedback, please contact our team at contact@cybertrace.com.au.



CYBERTRACE™

UNCLASSIFIED / OSINT