Discussion Paper

# Deepfakes: Manipulation on Demand?

Evolution of Deepfake Technology, Societal Impact, and the Path Forward

# Deepfakes: Manipulation on Demand?

Evolution of Deepfake Technology, Societal Impact, and the Path Forward

## Author

Maria Pawelec

## Commissioned by

Heinrich Böll Foundation Tel Aviv
Website: www.il.boell.org
→ info@il.boell.org

## Project Management

Israel Public Policy Institute (IPPI)
Website: www.ippi.org.il
→ office.israel@ippi.org.il

The opinions and views expressed in this publication are solely
those of the author/s and do not necessarily reflect the positions and/or viewpoints of the Heinrich Böll Foundation Tel Aviv and/or the Israel Public Policy Institute.

Date: March 2025

# Contents

# Executive Summary

This paper examines the rapid evolution and growing societal impact of deepfakes—hyperrealistic, synthetic, or manipulated audiovisual media created using advanced artificial intelligence (AI) techniques. Emerging from the intersection of computer graphics and machine learning, tools such as Generative Adversarial Networks (GANs) and GPT-based models now enable the creation and manipulation of lifelike images, audio, and video with minimal expertise. While deepfakes have positive applications in areas like entertainment, education, and marketing, they also present significant ethical and security challenges, including disinformation, propaganda, and image-based abuse.

Deepfakes are frequently exploited for political manipulation, election interference, and propaganda, posing threats to democratic processes and eroding public trust in media. They facilitate the "liar's dividend," which undermines confidence in authentic media, challenges legal norms, and perpetuates harm to marginalized groups, exacerbating gender-based inequalities and societal polarization. Moreover, they increasingly threaten personal privacy and security, with 98% of online deepfakes now consisting of (mostly non-consensual) sexualizing content, primarily targeting women. Additionally, deepfake technology is being leveraged for crimes such as impersonation and fraud, deceiving individuals and organizations alike.

Regulatory frameworks, such as the EU's Digital Services Act and AI Act, address some of these issues, e.g. by mandating content labeling for certain deepfake applications and enhancing social media moderation. However, stronger regulations are required, particularly to criminalize non-consensual sexualizing deepfakes. Encouragingly, the EU and several countries from around the world have recently taken significant steps in this direction. However, enforcement of existing regulations remains inadequate, not least due to insufficient resources for law enforcement and the sheer volume of deepfake content online. Therefore, advancements in detection technologies, including forensic tools for identifying manipulated media, are critical to combating unlabeled or misleading deepfakes.

In addition to technological and regulatory measures, societal strategies are crucial for mitigating the technology's adverse impact. Enhancing media literacy can empower the public to recognize and resist deceptive content, while providing technical and financial support to journalists and media organizations can enhance their ability to verify manipulated media effectively. Ethical education for AI developers is equally important, fostering accountability in the creation and deployment of these tools. Moreover, developers and companies offering deepfake technologies must adopt stricter policies and safeguards to prevent misuse.

Looking ahead, deepfakes are poised to play an increasingly influential role in digital interactions, with applications spanning from lifelike digital assistants to metaverse avatars. However, their potential for malicious use, coupled with significant ethical and political concerns, necessitates balanced societal, regulatory, and technical responses to maximize their benefits while effectively mitigating their risks.

# 1. What are Deepfakes?

Deepfakes are manipulated or synthetically generated images, videos, and audio of human faces, bodies, or voices that appear highly authentic. They are typically produced using deep learning, a method within machine learning. The outcome is a form of synthetic media that can portray people as saying or doing things they never actually said or did.

Deepfakes represent a new technological advancement within the longstanding history of media manipulation. Even before the era of Photoshop, images were altered for a range of purposes—some benign, such as artistic creation[1], and others more malicious, including political manipulation.[2] However, deepfakes amplify many challenges (and potential benefits) associated with manipulated audiovisual media for several key reasons.

> **Deepfakes are now more prevalent and influential than ever before. As their quality and accessibility continue to improve, they have proliferated across all areas of digital life, from politics to pornography.**

Firstly, deepfakes are spreading within a significantly altered and strained information environment. The rise of social media and messaging platforms has eroded the traditional media's "gatekeeping" role, boosting the influence of visual content and increasing the risk of disinformation. Secondly, most viewers tend to perceive audiovisual media as more credible than text-based media, often without an awareness of its potential for manipulation. Finally, advances in artificial intelligence (AI) and computing power have driven rapid improvements in the quality

and accessibility of convincing deepfakes. This trend has accelerated with the development of generative AI, as tools for text-to-image, text-to-speech (and speech-to-speech), and text-to-video generation now allow for the creation of targeted deepfakes using simple text prompts, requiring no input data or technical expertise.

Deepfakes are thus now more prevalent and influential than ever before. As their quality and accessibility continue to improve, they have proliferated across all areas of digital life, from politics to pornography. Some observers even believe that synthetic media will soon dominate the digital sphere, becoming the standard for online content.[3] The ethical and societal implications of deepfakes vary greatly depending on how they are used, which, in turn, calls for tailored political and societal responses. Consequently, it is crucial to understand the current state and future potential of deepfake technology, the diverse ways it is employed by different actors, and its broader impact.

# 2. Deepfakes: Technological Progress and Key Developers

Deepfakes are created using various AI and computer graphics techniques. Their origins can be traced back to 2014, when Google researcher Ian Goodfellow and colleagues published a paper[4] establishing the technological foundation for Generative Adversarial Networks (GANs), a specialized form of deep learning. In a GAN setup, two competing neural networks work together to produce increasingly realistic synthetic media. The first network, known as the "generator," creates an initial output. The second network, the "discriminator," assesses whether the output is genuine or fake, providing feedback to the generator, which then refines its synthesis

process to improve authenticity. Initially, most deepfakes were based on GANs. Following Goodfellow's seminal publication, researchers across universities, start-ups, and companies such as Nvidia, Samsung, and Disney advanced the technology[5] and improved its outputs.

> **The release of DALL-E in 2021, an image generator developed by OpenAI, marked the beginning of a new era in deepfake creation through generative AI.**

Concurrent to the development in research and the private sector, in 2017 the technology became available to the public and deepfakes began to proliferate when a user named u/deepfakes uploaded non-consensual sexualizing fake videos to Reddit, in which the faces of porn actors were replaced by those of female celebrities. This user's profile name led to the term "deepfakes." This same user also established the first online forum dedicated to creating and sharing (sexualizing) deepfakes and improving the underlying technology. Numerous other (code-sharing) fora followed. Since then, technological progress in deepfakes has been fueled not only by scientists and the industry (including a growing number of specialized start-ups and recent major players like OpenAI), but also by individual developers who often collaborate anonymously on code-sharing platforms such as GitHub.[6]

While GANs are capable of producing high-quality images, they have notable limitations, such as limited diversity in generated images and reduced control over specific details. As a result, developers continued to explore alternative models for media synthesis. The release of DALL-E in 2021, an image generator developed by OpenAI, marked the beginning of a new era in deepfake creation through generative AI. DALL-E was based on Generative Pre-Trained Transformers (GPT)[7],

a model initially developed for natural language processing. GPT employs unsupervised learning to detect patterns and relationships within large datasets, enabling it to predict subsequent words based on preceding text. For DALL-E, GPT was retrained using text-image pairs to predict "the next pixel," thus enabling it to generate images. Since then, GPT has been integrated with additional machine learning models, particularly Diffusion Models. Diffusion Models can generate high-definition images and allow for precise control over the output, significantly enhancing the quality and specificity of generated media.[8] This advancement has effectively elevated the potential of image (and other media) generation to new heights.

The technological advancements in deepfake creation have been substantial. Since their inception, researchers and professional developers have focused on enhancing the realism and resolution of deepfakes, reducing the amount of input data and training time needed, and accelerating the generation process. This includes progress toward real-time deepfake video generation[9], which currently relies on extensive pre-training but may, in the future, function without such preparation—potentially enabling seamless use in contexts like video conferences.

By 2019, GAN-based tools were already available online, allowing anyone to create realistic fake portraits with a single click.[10] Today, a wide range of text-to-image generators — such as DALL-E, MidJourney and StableDiffusion — enable users to create photorealistic images, even of specific individuals, in any setting simply by describing the desired outcome. Synthetic audio has also become highly realistic, with text-to-speech and speech-to-speech (voice cloning) applications, like ElevenLabs, now accessible to the public. These tools offer a variety of voices and can generate lifelike synthetic audio from any text prompt.

Deepfake video quality is advancing as well, both in terms of code available for more "traditional" forms of video deepfakes such as face swaps[11] and through generative AI. For instance, in February 2024, OpenAI introduced Sora, a video generator that can reportedly create photorealistic videos from text prompts and is expected to be publicly available soon. Both apps and online tools, as well as commercial services, have steadily expanded, turning deepfake production into a commodified process. This proliferation has been driven by specialized start-ups and tech giants like OpenAI, who are now key players in the field. Current estimates[12] indicate that there are over 2,000 tools available for creating video deepfakes, more than 10,000 tools for image generation, and over 1,000 for voice synthesis.

> **Generative AI has significantly accelerated a trend that has been unfolding since deepfakes were first released in 2017: a continual increase in their quality and accessibility, alongside a steady reduction in the expertise, training data, and computing power needed to create and deploy them.**

Generative AI has significantly accelerated a trend that has been unfolding since deepfakes were first released in 2017: a continual increase in their quality and accessibility, alongside a steady reduction in the expertise, training data, and computing power needed to create and deploy them. Consequently, today, even individuals without technological expertise can create convincing image and audio deepfakes with only a few clicks—and often with minimal or no input data.

# 3. Applications and Impacts of Deepfakes

Deepfake technology has permeated a wide array of fields, leaving a significant mark on politics, personal privacy, crime, entertainment, and cultural practices. While some applications showcase remarkable innovation and potential, others present grave ethical, societal, and legal challenges. This section explores the primary areas where deepfakes are employed, highlighting their transformative possibilities while critically examining the risks and controversies they entail.

## 3.1. Political Deception, Manipulation, and Propaganda

Public debate around deepfakes largely centers on their potential for political manipulation[13], with concerns especially heightened during democratic election cycles. Observers worry that deepfakes could be used to discredit candidates or disseminate false information about election procedures.[14] Such politically motivated uses of deepfakes have indeed been on the rise in recent years. For example, during Argentina's 2023 presidential election campaign, both candidates and their teams extensively disseminated deepfakes.[15] In Slovakia's 2023 parliamentary election campaign, an audio deepfake defaming progressive candidate Michal Šimečka circulated just days before the election, making it difficult to correct it in time.[16] In Turkey, one of the three presidential candidates withdrew from the race after the release of a non-consensual sexualizing deepfake.[17]

Concerns about deepfake-fueled election manipulation were also prominent ahead of the 2024 United States presidential election, where numerous deepfakes have been identified. For example, in January 2024, thousands of "robocalls"

featuring deepfake audio were recorded, posing a significant challenge for detection.[18] These robocalls were specifically designed to spread disinformation in a personalized way to individual voters. This tactic is just one of many examples of deepfakes in the U.S. election landscape. Notably, Donald Trump has even shared deepfakes, including AI-generated images of Taylor Swift fans wearing "Swifties for Trump" T-shirts, adding to the proliferation of politically charged synthetic media.[19]

Nevertheless, it appears that on a global scale, no individual deepfake has yet decisively influenced the outcome of a democratic election.[20] In Argentina, for example, the synthetic or manipulated nature of many shared deepfakes was relatively obvious, with some even being satirical. In Slovakia, it remains uncertain whether the defamatory deepfake contributed to Šimečka's defeat.[21] In Turkey, the candidate who withdrew due to a non-consensual sexualizing deepfake had limited prospects of winning the presidency. In the United States, analysts believe that deepfakes did not decisively sway the outcome of the election. However, their proliferation has further eroded the public's trust in the country's institutions and contributed to societal polarization by strengthening partisan beliefs and sowing distrust in the media and statements by political opponents.[22]

In the future, deepfakes in elections could unfold even greater manipulative and deceptive power, given their advancing quality, accessibility, and potential for strategic use. Malicious actors are likely to deploy high-quality deepfakes to undermine democratic processes and foster division and mistrust. This risk is not confined to national elections; local and regional elections could become prime targets, as actors typically have fewer resources to combat disinformation. Furthermore, candidates and their campaign teams themselves are increasingly leveraging deepfakes to promote their agendas in election

campaigns (see section on softfakes). Additionally, the mere fear of deepfake-based manipulation poses its own threat to democracy by eroding public confidence in the integrity of elections and the democratic representation they uphold.[23]

## It appears that on a global scale, no individual deepfake has yet decisively influenced the outcome of a democratic election.

More broadly, deepfakes erode trust in shared empirical insights, truths, and facts, hampering democratic debate and effective problem-solving. They also enable what is known as the "liar's dividend"—a concept that describes how individuals can deny incriminating evidence, such as recordings of their statements or actions, by asserting that the evidence is fake. By way of example, since 2020, speculation about the authenticity of videos featuring both Donald Trump[24] and Joe Biden[25] has been widespread in the U.S., with supporters and detractors questioning their veracity—demonstrating the far-reaching effects of deepfakes.

## As deepfakes become more widespread and public awareness of them grows, the liar's dividend—and the erosion of trust in general—only deepens.

Moreover, less than a month after George Floyd's murder by a police officer, a Republican politician published a "report" alleging that the video of the murder was a deepfake,[26] illustrating how the liar's dividend can further marginalize disadvantaged communities.[27] In the 2024 U.S. presidential elections, President-elect Donald Trump falsely claimed that pictures showing massive support for his opponent Harris at party rallies were faked.[28] In repressive regimes, claims of "deepfaked evidence"

could also undermine the efforts of human rights activists and political opponents.[29] As deepfakes become more widespread and public awareness of them grows, the liar's dividend—and the erosion of trust in general—only deepens.

Deepfakes can also be weaponized against individuals whose public statements or activities make them vulnerable to retaliation. This includes political opposition members, activists, and critical journalists. One notable example is Rana Ayyub[30], a Muslim Indian journalist who reported on a rape case involving a Hindu man. Following her coverage, a non-consensual, sexualizing deepfake video of Ayyub went viral, resulting in death threats against her. Moreover, deepfakes of lesser-known (or non-existent) people can also be used for targeted political attacks. In 2020, for example, a political activist and her partner were accused in an article in a newspaper in the US of sympathizing with terrorists. It was later discovered that the article was authored by a non-existent individual with a fake profile picture, showing how easily deepfakes can deceive even reputable outlets.[31]

Deepfakes are also employed to spread propaganda and exert foreign influence—aimed at influencing the domestic politics of other countries—both within and outside the context of wars and armed conflicts. A striking example occurred during Russia's war against Ukraine when, in 2022, a deepfake of Ukrainian President Zelensky urging Ukrainians to surrender marked a new form of deepfake-based interference. This video appeared on a hacked Ukrainian TV broadcaster's website and spread on social media. While the deepfake's low quality reduced its effectiveness, it was historic as the first documented attempt to use a deepfake to influence a major international conflict.[32]

Since then, deepfakes have become central to Russian disinformation efforts. Russia has used the technology to "erode support for Ukraine,

discredit [...] democratic institutions and officials, [and] seiz[e] on existing political divides"[33] in countries allied with Ukraine. For example, pro-Russian actors conducted a targeted disinformation campaign against France, French President Emmanuel Macron, and the International Olympic Committee (IOC) in the lead-up to and during the 2024 Paris Olympics.[34] This campaign included a disparaging "documentary" about the IOC featuring a synthetic voice clone of actor Tom Cruise.

> **Deepfakes are also employed to spread propaganda and exert foreign influence—aimed at influencing the domestic politics of other countries—both within and outside the context of wars and armed conflicts.**

In the Gaza conflict, both sides and their supporters deployed deepfakes for propaganda purposes.[35] Examples include pro-Israeli AI-generated images showing people cheering Israeli soldiers from their balconies[36] and pro-Hamas AI-generated images depicting Palestinian babies and children buried under rubble. These images set out to underscore the suffering of Palestinian civilians or to exaggerate public support for one of the warring parties with the aim of evoking strong emotional responses. Other deepfakes falsely depicted celebrities endorsing one side of the conflict,[37] likely spread by both domestic and foreign actors to advance the agendas of the respective warring parties. However, such deepfakes also contribute to growing uncertainty and mistrust in media evidence related to conflict situations.

In recent years, deepfake-based political disinformation, election manipulation, and propaganda have steadily increased. While the fear that a single deepfake might derail a democratic election has not yet materialized, deepfakes have become an integral part of broader disinformation

campaigns that erode trust in the integrity of elections and other democratic institutions, fuel societal polarization, target marginalized groups, and weaken the quality of democratic discourse. In these ways, deepfakes represent a serious threat to democracy.[38]

## 3.2. Non-Consensual Sexualizing Deepfakes and Image-Based Abuse

Empirical data on deepfakes reveals that up to 98% of all deepfakes online are pornographic[39] or non-consensual sexualizing images[40], with 99% of them targeting women. The technology is also used to generate child pornography, though this aspect remains challenging to investigate empirically. Numerous open-source deepfake projects and forums provide users with the tools needed to create deepfake pornography.[41] Apps like DeepNude, which circulate on messaging platforms such as Telegram,[42] can digitally "undress" images of any woman. While major AI image generators typically include filters to block "not safe for work" (NSFW) content, many smaller image generators either enable or specialize in creating such material.
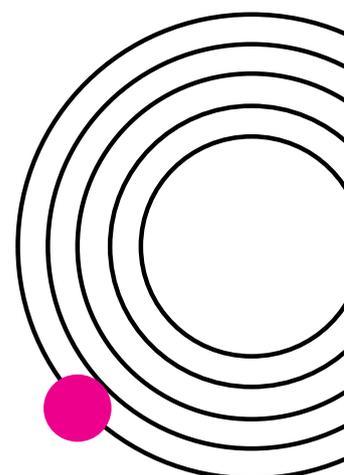
> **Deepfakes have become an integral part of broader disinformation campaigns that erode trust in the integrity of elections and other democratic institutions.**

Non-consensual sexualizing deepfakes initially targeted female celebrities, with the technology itself emerging in this context. To date, 94% of the victims of non-consensual sexualizing deepfake pornography work in the entertainment industry.[43] However, deepfakes' growing accessibility and the reduced amount of training data required now allow users to create non-consensual sexualizing

images of any woman who has shared photos on social media. This has led to an influx of deepfake pornography on both mainstream and dedicated porn sites, along with a rise in cases of revenge porn, sextortion, and blackmail.

> **Empirical data on deepfakes reveals that up to 98% of all deepfakes online are pornographic or non-consensual sexualizing images, with 99% of them targeting women.**

Creators of deepfake pornography often claim it is a form of free expression, entertainment, or harmless fun.[44] However, the impact on those affected is severe: victims frequently experience humiliation, a loss of control, and intimidation, often leading to psychological distress and setbacks in both their personal and professional lives.[45] Furthermore, deepfake pornography exploits the work of sex workers who created the original content, compounding the usual industry hazards and risks by misappropriating their work without consent or compensation.[46] Since deepfake porn overwhelmingly victimizes women, it poses a threat to gender equality in a digital society, potentially silencing women and discouraging them from participating in public and professional pursuits. Some observers argue that deepfake porn is a more pressing concern[47] than political deepfakes and that it requires greater attention.[48]

## 3.3. Criminal Uses of Deepfakes

Audio and video deepfakes have been increasingly employed in imposter schemes targeting both individuals and companies.[49] One notorious case involved a UK energy firm that lost over $240,000 in 2019 when its CEO was deceived by a synthetic audio deepfake mimicking the voice of the parent company's CEO.[50] A similar incident occurred in 2021 with a Hong Kong-based bank.[51] In 2022, fraudsters used an alleged deepfake video of a cryptocurrency executive to steal $32 million[52], while fake videos of Elon Musk continue to circulate, promoting fraudulent cryptocurrency schemes.[53] Additionally, a large-scale scam using deepfakes of Elon Musk, Justin Trudeau, Ryan Reynolds, and other celebrities has been luring victims into investing in "Quantum AI,"[54] a platform claiming to utilize quantum computing to generate profits.

> **Deepfakes can also be weaponized to sabotage competitors and manipulate stock markets by, for example, publishing fake videos of CEOs making offensive remarks or falsely announcing mergers or major business decisions.**

Deepfakes can also be weaponized to sabotage competitors and manipulate stock markets by, for example, publishing fake videos of CEOs making offensive remarks or falsely announcing mergers or major business decisions. Additionally, deepfakes are capable of circumventing facial recognition technology and remote identification processes. Since 2018, criminals have used image-manipulation apps to breach government facial recognition systems in China, facilitating large-scale tax fraud.[55] In 2022, the FBI warned that cybercriminals were applying for remote IT jobs using deepfakes to gain unauthorized access to corporate systems.[56] Thus, fraudsters have been executing increasingly sophisticated attacks, using real-time deepfakes to bypass facial recognition during online Know Your Customer (KYC) onboarding processes, particularly within the banking sector.[57]

Furthermore, pornographic and other compromising deepfakes can be exploited for blackmail and (s)extortion, as demonstrated in cases involving several male Indian politicians.[58] Cybercriminals are increasingly using deepfakes for such purposes. In the future, they may also employ real-time deepfakes in video calls or real-time voice cloning for phone scams, further intensifying the threat that deepfakes pose in the realm of criminal activity.[59]

## 3.4. Deepfakes in Personal Entertainment and Commercial Applications

Deepfakes are increasingly used in new forms of self-expression and digital engagement. Numerous YouTube and TikTok channels are dedicated to deepfakes, and especially those focused on celebrity deepfakes showcase highly sophisticated content.[60] For instance, the YouTube channel PresidentsUniverse, with 260,000 followers, has produced over 160 deepfake videos purportedly showing Donald Trump, Barack Obama, and Joe Biden playing popular video games.[61] Deepfakes also fuel meme culture[62], with apps enabling users to insert their faces into music or video clips. On TikTok, fake celebrity channels of figures like Margot Robbie[63] and Leonardo DiCaprio[64] have emerged, gaining popularity for their realistic renditions.

The technology also offers numerous commercial applications, especially in marketing and film. With deepfakes, actors can be seamlessly made to appear much younger or older, their voices cloned,

and lip movements synchronized in dubbed films.[65] Celebrities no longer need to travel to sets to participate in films or advertisements[66], as their likenesses can be digitally inserted, while post-editing can alter misspoken or age-inappropriate words without needing reshoots.[67] In 2023, British broadcaster ITVX launched Deep Fake Neighbor Wars,[68] the first fully deepfake-based series, featuring celebrities like Kim Kardashian and Greta Thunberg as if they were neighbors in a small-town setting, humorously portrayed in everyday arguments and interactions.

Deepfakes can even digitally "resurrect" deceased actors and musicians for new performances.[69] In the emerging "digital afterlife industry," deepfake technology is reshaping mourning practices and transforming how we engage with memories of those who have passed. For instance, the genealogy website MyHeritage launched DeepNostalgia[70] in 2021, a service that animates photos of deceased relatives. Voice cloning[71] enables individuals to create AI models of their voices while still alive, allowing their digital likeness to interact with loved ones after their death. This trend could even extend to deepfake avatars of the deceased attending their own funerals, profoundly altering cultural and personal approaches to grief and remembrance.[72]

In the marketing industry, deepfakes are employed to localize campaigns[73]—for instance, to promote local businesses—or even to hyper-personalize ads tailored to individual consumers.[74] The fashion industry is also embracing deepfakes[75], using them to create diverse digital models for online clothing displays and fittings, enhancing representation and personalization in fashion marketing. Additionally, in medical assistance technologies, synthetic audio powered by deepfake technology enables individuals who cannot speak due to medical conditions to communicate using a digital version of their own voice, enhancing personal expression.[76]

The rise of virtual influencers is another significant trend.[77] These digital personas, created by companies, agencies, or individuals, act as virtual models[78] or brand ambassadors, engaging audiences much like human influencers. Some virtual influencers, like Lil Miquela[79], launched in 2016, have amassed millions of followers and even appear to express political views—Lil Miquela, for instance, publicly supported the Black Lives Matter movement.[80] This development illustrates the expanding role of virtual entities in digital marketing, blurring the lines between human and digital influence in economic, social and political spheres.

Both non-commercial uses of deepfakes for personal entertainment and digital self-expression, as well as their commercial applications, raise significant ethical concerns. These range from the risks associated with micro-targeted ads, which are exacerbated by hyper personalized deepfakes, and the potential for virtual influencers to deceive or manipulate consumers, to issues surrounding actors' futures and the lack of consent for certain deepfakes. Additionally, deepfakes of the deceased[81] bring forth complex legal and ethical questions, including the preservation of an individual's rights and dignity, the importance of obtaining consent while they are still alive (or from their relatives), and the potential impact of these deepfakes on the grieving process and cultural practices surrounding mourning in the future.

## 3.5. Deepfakes in Education, Activism, and Satire

Educational deepfakes could allow historical figures to speak directly to students or museumgoers. For example, the Dalí Lives exhibition in St. Petersburg, Florida, allows visitors to "interact" with the deceased artist Salvador Dalí, presenting a personalized and immersive

experience of his persona.[82] Similarly, deepfake technology is used in projects that recreate the speeches of iconic figures such as former U.S. presidents or civil rights leader Martin Luther King Jr.[83], bringing their voices and messages to a contemporary audience in a powerful, memorable way. While these deepfakes enhance learning by creating vivid historical encounters, they inevitably present a singular, and sometimes simplified, interpretation of historical figures and events.[84] This selective portrayal underscores the importance of using and viewing educational deepfakes with discernment to maintain historical accuracy and avoid oversimplification.

> **Even when used for just causes, deepfakes carry the risk of misleading viewers, particularly when such clips are shared without context on social media.**

Political activists have employed deepfakes of murder victims in campaigns advocating for stricter gun laws in the United States[85] and for deeper criminal investigations in Mexico.[86] Artistic applications of deepfakes serve as a powerful means to highlight the dangers of mass data collection[87] and critique patriarchal attitudes.[88] In a landmark use of the technology, the 2020 documentary Welcome to Chechnya[89] used deepfakes to publicize the stories of victims of persecution based on their sexual orientation while protecting their real identities.

Deepfakes also play a role in political satire,[90] such as in web shows about U.S. President Trump or television programs featuring former Italian Prime Minister Matteo Renzi.[91] Individual satirical deepfakes are often used to critique those in positions of power and challenge their policies. For instance, a satirical deepfake in 2023 appeared to show German Federal Chancellor Olaf Scholz announcing a ban on the right-wing

party Alternative für Deutschland (AfD).[92] Since the deepfake was not clearly marked as satire, a German court prohibited its dissemination in February 2024 to prevent potential deception.[93]

Artistic, activist, and satirical deepfakes frequently serve to critique those in power or to spotlight socio-political injustices. The technology allows for the creation of uniquely impactful content, which can be intriguing or even haunting—such as portraying deceased individuals or using deepfakes to help viewers empathize with victims of persecution. This was evident in the aforementioned documentary Welcome to Chechnya, where activists' faces were digitally superimposed over those of persecuted individuals, rather than obscuring their identities through blurring or off-angle filming.

However, even when used for just causes, deepfakes carry the risk of misleading viewers, particularly when such clips are shared without context on social media. Additionally, some creators of harmful or defamatory deepfakes may label them as satire[94] to circumvent legal action or avoid platform-based consequences, such as labeling or algorithmic restrictions. This dual use underscores the ethical complexity of deepfakes in activism and satire, as well as the potential for their abuse under the guise of artistic or satirical intent.

## 3.6. Softfakes

In recent years, a new kind of deepfake application has emerged in the political sphere where politicians and campaign teams employ deepfakes during election campaigns to reach new target groups, disseminate political ideas, craft favorable self-representations, or even imply endorsements from well-known (and sometimes deceased) personalities. These politically driven, non-malicious deepfakes are often referred to as "softfakes."[95]

The first known instances of softfakes appeared in India in 2020 when the Bharatiya Janata Party (BJP) in Delhi released two videos of its leader campaigning in English and Haryanvi—languages he does not actually speak.[96] These videos, which did not indicate they contained synthetically generated content, were distributed across over 5,800 WhatsApp groups and reached more than 15 million voters. In 2022, South Korean President Yoon Suk-Yeol used a deepfake-based avatar of himself, "AI Yoon."[97] This avatar resembled him visually but communicated in a more relaxed, colloquial style, using humor to "respond" to thousands of citizens' questions, particularly to appeal to younger voters. The avatar proved so effective that Yoon's opponent eventually employed a similar deepfake-powered avatar in his own campaign.[98]

In 2024, softfakes gained visibility on an unprecedented scale, particularly during the nationwide elections in India. All major parties employed softfakes to engage citizens directly, reaching them in a highly personalized manner by speaking their native languages and tailoring messages to local contexts. For example, a video of Prime Minister Modi circulated directly to citizens via WhatsApp, addressing them by name in their mother tongues. Moreover, various parties used AI-driven "robocalls" featuring the voices of prominent politicians to convey party positions, often without citizens realizing they were not speaking to a real person. In addition, avatars of BJP representatives were used to promote party messages, appearing in personal settings and delivering content adapted to the recipient's location, as determined by their device.[99] Other softfakes included singing and dancing versions of candidates, designed to make them more relatable[100], as well as deepfakes of deceased politicians expressing support for particular parties.[101] These varied applications reflect the sophisticated integration of softfakes into political campaigns, aiming to build personalized connections and sway public opinion on an extensive scale.

> **In 2024, softfakes gained visibility on an unprecedented scale, particularly during the nationwide elections in India. All major parties employed softfakes to engage citizens directly.**

In Indonesia, a deepfake of a deceased political leader appeared in 2024[102], seemingly endorsing his former political party. Presidential candidate General Subianto also employed a "cuddly" AI avatar—a softened, non-photorealistic version of himself—to appear more approachable and distract from allegations of past human rights abuses.[103] Meanwhile, in Pakistan, deepfake technology allowed imprisoned former Prime Minister Imran Khan to connect with the public.[104] His written messages from prison were converted into video messages, enabling him to engage with his electorate and maintain a voice in political discourse despite his confinement.

Softfakes allow parties and candidates to reach new target groups, including marginalized communities, engage more directly with citizens, and deliver election messages and promises in a hyper-personalized manner. They can enhance candidates' public image and create innovative messages of support, potentially mobilizing greater citizen participation in political issues and discussions. Due to their cost-effectiveness and impact, softfakes are increasingly part of the toolkit for political campaign teams and specialized PR firms. While softfakes have primarily been used in Asia, political parties in countries like France[105] and Germany[106] are beginning to explore generative AI for political messaging. This trend suggests that softfakes may become an integral

part of election strategies worldwide. However, as the novelty of this technology fades, its influence may also wane over time.

> **Softfakes allow parties and candidates to reach new target groups, including marginalized communities, engage more directly with citizens, and deliver election messages and promises in a hyper-personalized manner.**

Despite their seemingly harmless nature, softfakes carry significant risks as well: they can mislead viewers, intentionally or unintentionally, about candidates' backgrounds, abilities, or positions. Therefore, it is essential to label all softfakes clearly to promote transparency—a requirement that, unfortunately, is often not practiced. However, even clear labeling does not eliminate all concerns. For instance, softfake-generated campaign content may still convey false or misleading information. During the 2024 Indian elections, for example, robocallers "hallucinated," falsely attributing positions to candidates.[107]

In addition to such unintentional errors, even labeled softfakes can be purposefully deceptive. The "AI Yoon" avatar and the "cuddly" version of General Subianto, for instance, served to amuse the public, present the candidates as younger and more approachable, and distract from past scandals. Like other forms of disinformation, such portrayals can leave a lasting impact, even if transparently labeled. Moreover, softfakes of deceased politicians introduce further ethical complexities, as previously discussed in the section on personal and commercial uses. Finally, hyper-personalized softfakes may also fragment political discourse, contributing to societal division and undermining cohesion.[108]

# 4. Regulatory, Technological and Societal Responses

As deepfake technology advances, its applications raise profound questions about digital trust, authenticity, and societal adaptation. Addressing these challenges requires coordinated regulatory, technological, and societal responses. This section explores the emerging solutions designed to mitigate risks while fostering responsible innovation.

## 4.1. Regulatory Responses

Deepfakes cannot be categorically outlawed due to their numerous legitimate applications across fields like education, commerce, and entertainment. Many deepfakes also fall under the protection of freedom of expression and artistic freedom. At the same time, as demonstrated, deepfake-based disinformation and hate speech pose serious threats to democracy and can inflict significant harm on individuals. Policymakers are therefore faced with the complex task of finding balanced approaches to address the potential harms of deepfakes while preserving their beneficial uses.

> **Deepfakes cannot be categorically outlawed due to their numerous legitimate applications across fields like education, commerce, and entertainment.**

The growing awareness that self- and co-regulation by social media platforms has proven insufficient to combat the escalating threat of deepfake-based disinformation, coupled with rapid advancements in AI, has led the European

Union (EU) to adopt stronger regulatory measures targeting deepfakes. Platform regulation plays a central role in these efforts, as social media platforms are often the primary channels for the dissemination of deepfakes.

In 2022, the EU introduced the Digital Services Act (DSA)[109], which primarily targets large social media platforms with over 45 million users within the EU. Its objective is to enhance the safety and reliability of online interactions. The DSA mandates that major platform operators establish transparent content moderation rules, detailing enforcement measures to ensure compliance. For example, platforms can implement a "notice-and-takedown" procedure, allowing victims of deepfake-related harm to report harmful content directly to operators. Upon receiving a report, platform operators are then required to review the content, inform the complainant of the outcome, and potentially take action such as blocking or deleting the material. Individuals impacted by such moderation actions also have the right to appeal.

The DSA also allows for the registration of "trusted flaggers" — accredited institutions whose content reports must be prioritized by platform operators. Platforms may also use (partially) automated tools, such as upload filters, to identify and remove harmful content proactively. Furthermore, the DSA obligates platforms to cooperate with law enforcement on deepfakes that could be subject to criminal prosecution, enhancing the regulatory framework around digital accountability and security. However, national laws ultimately define which deepfakes are considered criminal offenses, such as those used for defamation or fraud. Additionally, many deepfakes infringe upon privacy rights or copyright protections. Beyond platform regulation, the EU's 2024 AI Act (coming into effect in 2026) is another central framework for regulating deepfakes on the European level.[110] Often seen as a potential model for global AI

governance (the so-called "Brussels Effect"[111]), the AI Act introduces a risk-based, context-sensitive approach, regulating AI applications by their associated risks based on context.

## Critics argue that categorizing deepfakes as low-risk AI downplays their potential for political manipulation and abuse.

Under this system, deepfakes are classified as low or minimal risk rather than prohibited or high-risk AI, meaning the AI Act primarily enforces transparency measures. Providers of AI capable of generating deepfakes are required to embed mechanisms (e.g., watermarking or metadata labeling) so that manipulated or generated content can be clearly identified. Those who create deepfakes must disclose that their creations are based on AI, ideally allowing viewers to understand that the content is synthetic or manipulated. However, deepfakes used for artistic or satirical purposes only need to be labeled in ways that do not impair the work's enjoyment, potentially leaving gaps for malicious uses.

Critics argue that categorizing deepfakes as low-risk AI downplays their potential for political manipulation and abuse.[112] Lastly, concerns exist around practical enforcement, with transparency requirements' implementation details yet to be clarified. This has led to doubts about the AI Act's effectiveness in preventing misuse and ensuring the responsible use of deepfake technology.

Other countries are also introducing specialized legislation to address the challenges posed by AI and deepfakes in particular. In January 2023, China enacted the "Provisions on Deep Synthesis Technology,"[113] a specialized law addressing deepfake technology. Unlike regulations in other regions that focus on platforms disseminating deepfakes, this legislation targets the creators

and providers of synthetic media technologies. The law mandates that all synthetic content carry clear labels, requires consent from all individuals depicted in such media, and bans numerous applications of deepfake technology, including sexualizing deepfakes, content infringing on intellectual property, and any material deemed counter to China's "national interests."[114] This legislation has sparked concerns over potential state censorship. Also, the expansive requirement for consent limits various applications, such as satire, and could stifle creative expression. By imposing broad constraints on deepfake content, the law reflects a regulatory approach that prioritizes control over synthetic media sources, setting China apart in its restrictive stance on deepfake technology.

Regulators worldwide are intensifying efforts to address the issue of non-consensual sexualizing deepfakes, which are often not adequately covered by existing criminal laws, such as those targeting image-based abuse. As a result, victims frequently lack sufficient legal protection.[115] To combat this gap, several regions have enacted new legislation within the last two years, criminalizing the creation, distribution, or even consumption of non-consensual sexualizing deepfakes.

> **In the European Union, the Digital Services Act (DSA) supports the reporting and removal of non-consensual sexualizing deepfakes from major social media platforms.**

In the United States, multiple states have recently passed laws specifically targeting non-consensual deepfake pornography, aiming to provide victims with legal recourse and hold offenders accountable.[116] Similarly, the UK[117] and South Korea[118] have introduced regulations that broaden protections against non-consensual sexualizing

deepfakes, strengthening penalties and increasing prosecutorial focus on such offenses.

In the European Union, the Digital Services Act (DSA) supports the reporting and removal of non-consensual sexualizing deepfakes from major social media platforms. Additionally, the "Directive on combating violence against women and domestic violence," enacted in June 2024, explicitly bans the production and dissemination of intimate or manipulated images without consent, marking a significant step toward safeguarding individuals—particularly women—against sexualizing violence through deepfakes.[119] That these legal frameworks are taking shape shows that politicians, lawmakers, and civil society organizations are increasingly focused on curbing the spread and impact of non-consensual sexualizing deepfakes, advancing legal protections for victims and establishing stronger deterrents.

## 4.2. Technological and Societal Responses

While regulation is a crucial step in addressing the risks associated with deepfakes, it alone would be insufficient, particularly given that malicious actors are likely to disregard it. Technical advancements in deepfake detection are therefore essential to identify unlabeled or misleading deepfakes, enabling society to respond more effectively. Significant strides in image and audio forensics are underway, with detection methods focusing on various distinctive traits of deepfakes. However, this progress can paradoxically fuel improvements in deepfake technology itself, turning the effort into a continuous "cat-and-mouse game."[120]

Besides deepfake detection, one promising area of development is the encoding of images to prevent their unauthorized use in AI-based image synthesis.[121] Additionally, verification systems

for authenticating untampered media may provide a supplementary approach to combat deepfake challenges. These systems, if effectively implemented and widely adopted, could help establish the credibility of genuine content. Yet, they are not insurmountable and their limitations and potential biases must be considered.[122] For instance, if content authentication becomes a standard or "watermark" for major media outlets, evidence captured by marginalized groups—such as citizen journalists or activists documenting human rights abuses with basic mobile devices that lack such technology—may be unfairly discounted or deemed unreliable.[123]

In addition to regulatory and technical measures, societal strategies are essential for effectively managing the impact of deepfakes. First and foremost, enhancing citizens' media literacy[124] and raising awareness about the risks posed by deepfakes can help reduce the influence and spread of deceptive or harmful content.[125] Empowering journalists and media organizations with both technical and financial resources to collaboratively verify increasing volumes of potentially manipulated media is also crucial, as it reinforces their critical role in upholding information integrity.

> **While regulation is a crucial step in addressing the risks associated with deepfakes, it alone would be insufficient, particularly given that malicious actors are likely to disregard it.**

Furthermore, ethical considerations should be woven into the educational pathways of future developers—such as through curricula in computer science and related fields—promoting a sense of responsibility among those who may create or work with deepfake technology.

Companies that develop and offer deepfake tools must also take proactive steps to prevent abuse, adopting stronger policies and technologies to safeguard against potential abuses.[126]

> **Ethical considerations should be woven into the educational pathways of future developers— such as through curricula in computer science and related fields—promoting a sense of responsibility among those who may create or work with deepfake technology.**

Together, these combined efforts can help mitigate the adverse effects of deepfakes while fostering legitimate and beneficial applications of the technology.

# 5. Where are Deepfakes Headed?

Since the technology first emerged publicly in 2017, deepfakes have been on a rapid rise. The launch of image generator DALL-E in 2021 marked a major turning point, as generative AI significantly accelerated the accessibility and quality of deepfakes. Now, convincing audio and image deepfakes can be created with a simple text description, and similar advancements in video deepfakes are anticipated soon. As a result, increasingly realistic deepfakes featuring both celebrities and everyday individuals are becoming woven into our digital interactions, signaling a shift in everyday digital interactions and how we experience and interpret online content.

While deepfakes are still overwhelmingly present in the realm of image-based abuse, they are now permeating nearly every area

of digital life, from politics and cybercrime to entertainment, advertising, and education. In politics, deepfakes have long been regarded as a powerful tool for manipulation and deception. Yet, new applications—such as political art, education, satire, and even election campaigning (so-called "softfakes")—are also emerging, offering potential avenues for increased participation, mobilization, justice, and enhanced accountability for politicians. However, these developments also bring new risks, as these tools can also be misused or have unintended consequences.

> **While deepfakes are still overwhelmingly present in the realm of image-based abuse, they are now permeating nearly every area of digital life, from politics and cybercrime to entertainment, advertising, and education.**

Deepfake technologies are increasingly combined with Large Language Models (LLMs), another influential form of generative AI, to create chatbots with human-like faces, realistic digital twins or avatars, and advanced digital assistants. While these integrations offer innovative applications, such as lifelike customer support or interactive learning experiences, LLMs also facilitate large-scale disinformation. When coupled with deepfakes, this allows for sophisticated multimodal disinformation, heightening concerns about AI-powered manipulation and deception on an unprecedented scale. This convergence intensifies fears around the potential abuse of AI for misleading and persuasive content that can easily blur the line between reality and fabrication.

Generative AI's influence on both work and personal life is growing rapidly. Looking to the future, proponents of the "metaverse"[127] foresee a digital evolution where hyperrealistic synthetic versions of ourselves—fully realized digital clones— or other highly realistic avatars interact within an immersive virtual world. Deepfake technology is anticipated to play a central role in this vision, enabling the creation of highly realistic avatars that bring a new level of lifelike presence to these digital spaces. As this technology advances, it could make virtual interactions in the metaverse almost visually indistinguishable from real-world encounters, reshaping how people connect, work, and socialize in digital environments.

Deepfakes hold substantial promise across sectors, from commercial gain and entertainment to education and cross-cultural communication. They offer the potential for meaningful interactions across linguistic and geographic divides—and even beyond death, as illustrated by the emerging digital afterlife industry. However, their rise also amplifies serious risks, including disinformation, manipulation, fraud, blackmail, and intimidation. This rapid development raises complex ethical and societal questions concerning individual rights, freedom of expression, and cultural practices related to death and mourning.

> **Deepfake technologies are increasingly combined with Large Language Models (LLMs), another influential form of generative AI, to create chatbots with human-like faces, realistic digital twins or avatars, and advanced digital assistants.**

In response, scientists, policymakers, and communities are actively exploring approaches to address these challenges and to manage the multifaceted impacts of deepfakes in different contexts. Regulatory frameworks, such as the EU's Digital Services Act, AI Act, and Directive on combating violence against women and domestic

violence aim to promote transparency and accountability and to protect individuals from abuse, e.g. by enforcing content moderation on social media, mandating the labeling of synthetic media, or outlawing non-consensual sexualizing deepfakes. Advances in detection technologies and content authentication systems offer promising solutions for combating disinformation and preserving trust in digital media. Additionally, societal efforts, including educational initiatives to improve media literacy and empower content verifiers, are crucial to strengthening societal resilience against manipulation. These combined approaches represent essential steps toward fostering an ethical and responsible integration of deepfake technology.

# 6. Outlook

As deepfake technology continues to evolve, it embodies both the potential and the formidable risks of artificial intelligence in the digital age. While deepfakes have promising applications in various fields, their capacity for harm—ranging from privacy violations and fraud to political manipulation and image-based abuse—demands careful consideration. The widespread accessibility of deepfake tools, coupled with advancements in generative AI, places these capabilities in the hands of both creators and malicious actors alike, raising urgent questions around digital trust, ethical usage, and the boundaries of synthetic media.
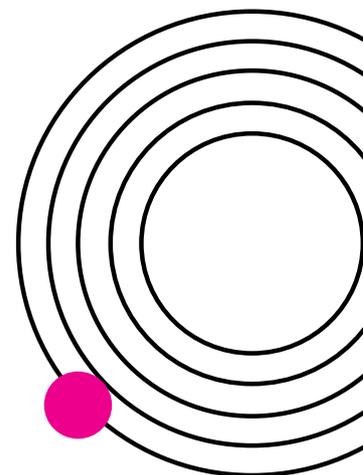
Regulatory frameworks are emerging worldwide to address the challenges of deepfakes, with the European Union, the UK, South Korea and other regions and countries leading efforts to promote transparency, protect individual rights, and mitigate disinformation. Yet, regulation alone cannot fully resolve the unique challenges posed by this technology. A multifaceted approach—combining policy, technological safeguards, ethical

development standards, and societal awareness—will be necessary to responsibly integrate deepfakes into society.

> **Deepfakes represent a pivotal moment in the intersection of technology, society, and ethics, challenging us to shape a future that balances innovation with responsibility and accountability in a rapidly changing digital landscape.**

Looking ahead, deepfakes are poised to become more deeply embedded in digital life, from virtual assistants and avatars in the metaverse to increasingly realistic applications in cross-cultural and commercial settings. Maximizing their potential while addressing inherent risks requires a forward-looking strategy that continuously adapts to technological developments. The path forward involves not only harnessing the positive uses of deepfake technology but also cultivating a resilient digital ecosystem where authenticity is preserved, individuals' rights are protected, and public trust remains intact.

Deepfakes represent a pivotal moment in the intersection of technology, society, and ethics, challenging us to shape a future that balances innovation with responsibility and accountability in a rapidly changing digital landscape.

# Endnotes

[1] Burkell, J., & Gosse, C. (2019). Nothing new here: Emphasizing the social and cultural context of deepfakes. First Monday, 24(12). https://doi.org/10.5210/fm.v24i12.10287

[2] J. Sharma and R. Sharma, "Analysis of Key Photo Manipulation Cases and their Impact on Photography," IIS University Journal of Arts, vol. 6, no. 1, pp. 88-99, 2017. ISSN 2319-5339

[3] Deepfake-Expertin Nina Schick: Jede Identität ist heute gefährdet," Der Standard, July 1, 2022. Available: https://www.derstandard.de/story/2000137059050/deepfake-expertin-nina-schick-jede-identitaet-ist-heute-gefaehrdet

[4] I.J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative Adversarial Networks," arXiv:1406.2661 [stat.ML], June 10, 2014. Available: https://doi.org/10.48550/arXiv.1406.2661

[5] KI in der Praxis: Deepfakes: Wie alles begann - und wohin es führen könnte," The Decoder, April 27, 2022. Available: https://the-decoder.de/geschichte-der-deepfakes-so-rasant-geht-es-mit-ki-fakes-voran/

[6] Winter, R., & Salter, A. (2019). DeepFakes: uncovering hardcore open source on GitHub. Porn Studies, 7(4), 382–397. https://doi.org/10.1080/23268743.2019.1642794

[7] M. Chruściński, "A brief history of AI-powered image generation," Sii Blog, April 5, 2023. Available: https://sii.pl/blog/en/a-brief-history-of-ai-powered-image-generation/

[8] S. Ben Souissi, "Diffusion Models: Ein neuer Horizont in der Bilderzeugung," SocietyByte, August 9, 2023. Available: https://www.societybyte.swiss/2023/08/09/diffusion-modells-ein-neuer-horizont-in-der-bilderzeugung/

[9] J. Kahn, "Deepfakes are stealing the show on 'America's Got Talent.' Will they soon steal a lot more too?" Fortune, September 6, 2022. Available: https://fortune.com/2022/09/06/deepfakes-americas-got-talent-metaphysic-fraud-metaverse/

[10] "This Person Does Not Exist," Available: https://this-person-does-not-exist.com/en

[11] L. Finger, "Overview Of How To Create Deepfakes - It's Scarily Simple," Forbes, September 8, 2022. Available: https://www.forbes.com/sites/lutzfinger/2022/09/08/overview-of-how-to-create-deepfakesits-scarily-simple/

[12] Sensity AI, "The State of Deepfakes 2024". Available: https://5865987.fs1.hubspotusercontent-na1.net/hubfs/5865987/SODF%202024.pdf

[13] C. Gosse and J. Burkell, "Politics and porn: how news media characterizes problems presented by deepfakes," Critical Studies in Media Communication, vol. 37, no. 5, pp. 497–511, 2020. Available: https://doi.org/10.1080/15295036.2020.1832697

[14] N. Diakopoulos and D. Johnson, "Anticipating and Addressing the Ethical Implications of Deepfakes in the Context of Elections," New Media & Society, October 23, 2019. Available: SSRN: https://ssrn.com/abstract=3474183 or http://dx.doi.org/10.2139/ssrn.3474183

[15] "Is Argentina the First A.I. Election?" The New York Times, November 15, 2023. Available: https://www.nytimes.com/2023/11/15/world/americas/argentina-election-ai-milei-massa.html

[16] "Slovakia's Election Deepfakes Show AI Is a Danger to Democracy," Wired, October 2023. Available: https://www.wired.com/story/slovakias-election-deepfakes-show-ai-is-a-danger-to-democracy/

[17] D. Meyer, "Turkey's deepfake-influenced election spells trouble," Fortune, May 15, 2023. Available: https://fortune.com/europe/2023/05/15/turkeys-deepfake-influenced-election-spells-trouble/

**18** E. Steck and A. Kaczynski, "Fake Joe Biden robocall urges New Hampshire voters not to vote in Tuesday's Democratic primary," CNN, January 22, 2024. Available: https://edition.cnn.com/2024/01/22/politics/fake-joe-biden-robocall/index.html

**19** N. Robins-Early, "How did Donald Trump end up posting Taylor Swift deepfakes?" The Guardian, August 26, 2024. Available: https://www.theguardian.com/technology/article/2024/aug/24/trump-taylor-swift-deepfakes-ai

**20** M. Łabuz and C. Nehring, "On the way to deep fake democracy? Deep fakes in election campaigns in 2023," European Political Science, 2024. Available: https://doi.org/10.1057/s41304-024-00482-9

**21** M. Łabuz and C. Nehring, "On the way to deep fake democracy? Deep fakes in election campaigns in 2023," European Political Science, 2024. Available: https://doi.org/10.1057/s41304-024-00482-9

**22** P. Verma, W. Oremus and C. Zakrzewski, „AI didn't sway the election, but it deepened the partisan divide", The Washington Post, November 9, 2024. Available: https://www.washingtonpost.com/technology/2024/11/09/ai-deepfakes-us-election/

**23** M. Pawelec, "Deepfakes and Democracy (Theory): How Synthetic Audio-Visual Media for Disinformation and Hate Speech Threaten Core Democratic Functions," Digital Society, vol. 1, p. 19, 2022. Available: https://doi.org/10.1007/s44206-022-00010-6

**24** P. Beaumont, "Donald Trump fans cry betrayal as he rebukes Capitol violence," The Guardian, January 8, 2021. Available: https://www.theguardian.com/us-news/2021/jan/08/trump-incites-anger-among-acolytes-let-down-by-lack-of-support

**25** J. Horton and S. Sardarizadeh, "False claims of 'deepfake' President Biden go viral," BBC News, July 28, 2022. Available: https://www.bbc.com/news/62338593

**26** Z. Budryk, "GOP House candidate publishes 23-page report claiming George Floyd death was deepfake video," The Hill, June 24, 2020. Available: https://thehill.com/homenews/house/504429-gop-house-candidate-publishes-23-page-report-claiming-george-floyd-death-was/

**27** "The threat posed by deepfakes to marginalized communities," Brookings, February 2021. Available: https://www.brookings.edu/techstream/the-threat-posed-by-deepfakes-to-marginalized-communities/

**28** P. Verma, W. Oremus and C. Zakrzewski, „AI didn't sway the election, but it deepened the partisan divide", The Washington Post, November 9, 2024. Available: https://www.washingtonpost.com/technology/2024/11/09/ai-deepfakes-us-election/

**29** S. Gregory, "Authoritarian Regimes Could Exploit Cries of 'Deepfake'," Wired, February 14, 2021. Available: https://www.wired.com/story/opinion-authoritarian-regimes-could-exploit-cries-of-deepfake/

**30** N. Jankowicz, "The threat from deepfakes isn't hypothetical. Women feel it every day," The Washington Post, March 25, 2021. Available: https://www.washingtonpost.com/opinions/2021/03/25/threat-deepfakes-isnt-hypothetical-women-feel-it-every-day/

**31** R. Satter, "Deepfake used to attack activist couple shows new disinformation frontier," Reuters, July 15, 2020. Available: https://www.reuters.com/article/us-cyber-deepfake-activist/deepfake-used-to-attack-activist-couple-shows-new-disinformation-frontier-idUSKCN24G15E

**32** "KI und Gesellschaft: Möglicher Selenskyj-Deepfake: Miserabel und dennoch historisch," The Decoder, March 17, 2022. Available: https://the-decoder.de/selenskyj-deepfake-miserabel-und-dennoch-historisch/

**33** S. Bond, "This is what Russian propaganda looks like in 2024," NPR, June 6, 2024. Available: https://www.npr.org/2024/06/06/g-s1-2965/russia-propaganda-deepfakes-sham-websites-social-media-ukraine

**34** D. Milmo, "Russia targets Paris Olympics with deepfake Tom Cruise video," The Guardian, June 3, 2024. Available: https://www.theguardian.com/technology/article/2024/jun/03/russia-paris-olympics-deepfake-tom-cruise-video

**35** D. Klepper, "Fake babies, real horror: Deepfakes from the Gaza war increase fears about AI's power to mislead," Associated Press, November 28, 2023. Available: https://apnews.com/article/artificial-intelligence-hamas-israel-misinformation-ai-gaza-a1bb303b637ffbbb9cbc3aa1e000db47

**36** Reuters Fact Check, „Fact Check: Photo of cheering crowds waving Israeli flags at soldiers is AI-generated", Reuters, October 30, 2023. Available: https://www.reuters.com/fact-check/photo-cheering-crowds-waving-israeli-flags-soldiers-is-ai-generated-2023-10-30/

**37** V. Strauss, "Fake celebrity statements and videos on Israel-Hamas war — and other news literacy lessons," The Washington Post, November 10, 2023. Available: https://www.washingtonpost.com/education/2023/11/10/fake-celebrity-statements-war-news-literacy/

**38** Pawelec, M. "Deepfakes and Democracy (Theory): How Synthetic Audio-Visual Media for Disinformation and Hate Speech Threaten Core Democratic Functions." DISO 1, 19 (2022). https://doi.org/10.1007/s44206-022-00010-6

**39** "State of Deepfakes", Home Security Heroes, Available: https://www.securityhero.io/state-of-deepfakes/#overview-of-current-state

**40** I refer to non-consensual sexual deepfake depictions as "sexualizing" or "image-based abuse" rather than "pornographic", since the negative connotations of the term "pornography" only refer to the public or private thematization of sexuality as obscene, rather than to the severe violation of victim's human rights (see https://www.nomos-elibrary.de/10.5771/9783748930297-42.pdf).

**41** Winter, R., & Salter, A. (2019). "DeepFakes: uncovering hardcore open source on GitHub." Porn Studies, 7(4), 382–397. https://doi.org/10.1080/23268743.2019.1642794

**42** E.J. Dickson, "TikTok Stars Are Being Turned Into Deepfake Porn Without Their Consent," Rolling Stone, October 14, 2020. Available: https://www.rollingstone.com/culture/culture-features/tiktok-creators-deepfake-pornography-discord-pornhub-1078859/

**43** "State of Deepfakes", Home Security Heroes, Available: https://www.securityhero.io/state-of-deepfakes/#overview-of-current-state

**44** J. Burkell and C. Gosse, "Nothing new here: Emphasizing the social and cultural context of deepfakes," First Monday, vol. 24, no. 12, December 2, 2019. Available: https://firstmonday.org/ojs/index.php/fm/article/download/10287/8297. DOI: http://dx.doi.org/10.5210/fm.v24i12.10287

**45** K. Hao, "Deepfake porn is ruining women's lives. Now the law may finally ban it," MIT Technology Review, February 12, 2021. Available: https://www.technologyreview.com/2021/02/12/1018222/deepfake-revenge-porn-coming-ban/

**46** M. MacDonald, "The Double Exploitation of Deepfake Porn," The Walrus, June 10, 2021 (Updated November 9, 2023). Available: https://thewalrus.ca/the-double-exploitation-of-deepfake-porn/

**47** M. MacDonald, "The Double Exploitation of Deepfake Porn," The Walrus, June 10, 2021 (Updated November 9, 2023). Available: https://thewalrus.ca/the-double-exploitation-of-deepfake-porn/

**48** K. Hao, "Deepfake porn is ruining women's lives. Now the law may finally ban it," MIT Technology Review, February 12, 2021. Available: https://www.technologyreview.com/2021/02/12/1018222/deepfake-revenge-porn-coming-ban/

**49** L. Edmonds, "Scammer used deepfake video to impersonate U.S. Admiral on Skype chat and swindle nearly $300,000 out of a California widow," Daily Mail, October 24, 2020. Available: https://www.dailymail.co.uk/news/article-8875299/Scammer-uses-deepfake-video-swindle-nearly-300-000-California-widow.html

**50** J. Damiani, „A Voice Deepfake Was Used To Scam A CEO Out Of $243,000" , Forbes, September 3, 2019. Available at: https://www.forbes.com/sites/jessedamiani/2019/09/03/a-voice-deepfake-was-used-to-scam-a-ceo-out-of-243000/

**51** T. Brewster, "Fraudsters Cloned Company Director's Voice In $35 Million Heist, Police Find," Forbes, October 14, 2021 (Updated May 2, 2023). Available: https://www.forbes.com/sites/thomasbrewster/2021/10/14/huge-bank-fraud-uses-deep-fake-voice-tech-to-steal-millions/

**52** D.A.J. Sokolov, "Krypto-Diebstahl mit Deepfake eines Hologramms beschert 32 Millionen Dollar," Heise Online, August 25, 2022. Available: https://www.heise.de/news/Krypto-Diebstahl-mit-Deepfake-eines-Hologramms-beschert-32-Millionen-Dollar-7242635.html

**53** J. Tidy, "YouTube accused of not tackling Musk Bitcoin scam streams," BBC News, June 10, 2022. Available: https://www.bbc.com/news/technology-61749120

**54** B. Potaracke, "Everything You Need to Know About Quantum AI Scams," Locknet Managed IT, August 19, 2024. Available: https://www.locknetmanagedit.com/blog/cybersecurity/quantum-ai-scams

**55** M. Borak, "Chinese government-run facial recognition system hacked by tax fraudsters: report," South China Morning Post, March 31, 2021. Available: https://www.scmp.com/tech/tech-trends/article/3127645/chinese-government-run-facial-recognition-system-hacked-tax

**56** FBI Internet Crime Complaint Center (IC3), "Deepfakes and Stolen PII Utilized to Apply for Remote Work Positions," IC3, June 28, 2022. Available: https://www.ic3.gov/Media/Y2022/PSA220628

**57** Sensity AI, "The State Of Deepfakes", 2024. Available: https://5865987.fs1.hubspotusercontent-na1.net/hubfs/5865987/SODF%202024.pdf

**58** "Deepfake clips: Sextortionists target celebs," The Times of India, February 22, 2021. Available: https://timesofindia.indiatimes.com/city/mumbai/deepfake-clips-sextortionists-target-celebs/articleshow/81162493.cms

**59** Sensity AI, "The State Of Deepfakes", 2024. Available: https://5865987.fs1.hubspotusercontent-na1.net/hubfs/5865987/SODF%202024.pdf

**60** "Deepfake-Tom Cruise: Ist die Dystopie jetzt Wirklichkeit?" MIXED, February 27, 2021. Available: https://mixed.de/deepfake-tom-cruise-bei-tiktok-ist-die-dystopie-jetzt-wahr/

**61** "Presidents Universe," YouTube Channel. Available: https://www.youtube.com/@PresidentsUniverse

**62** K. Hao, "Memers are making deepfakes, and things are getting weird," MIT Technology Review, August 28, 2020. Available: https://www.technologyreview.com/2020/08/28/1007746/ai-deepfakes-memes/

**63** @fake_robbie, TikTok Profile. Available: https://www.tiktok.com/@fake_robbie

**64** @fake_dicaprio, TikTok Profile. Available: https://www.tiktok.com/@fake_dicaprio

**65** S. Glass, "Plötzlich kann Morgan Freeman Deutsch," Tagesschau, May 1, 2021. Available: https://www.tagesschau.de/wirtschaft/technologie/synchronisation-film-deep-fake-stimmen-101.html

**66** "Hulu: Deepfake-Spot statt Corona-Pause," MIXED, August 23, 2020. Available: https://mixed.de/hulu-deepfake-spot-statt-corona-pause/

**67** J. Porter, "Fork me: 'Fall' movie removed more than 30 F-bombs with deepfake dub technology," The Verge, August 10, 2022. Available: https://www.theverge.com/2022/8/10/23299565/fall-2022-movie-fucks-gone-ai-technology-flawless-f-bomb

**68** "Deep Fake Neighbour Wars," ITV. Available: https://www.itv.com/watch/deep-fake-neighbour-wars/10a2895/10a2895a0001

**69** P. Han-sol, "Holographic performances of dead stars welcomed, with caution," The Korea Times, January 21, 2021. Available: https://www.koreatimes.co.kr/www/nation/2021/01/703_302548.html

**70** "Deep Nostalgia," MyHeritage. Available: https://www.myheritage.de/deep-nostalgia

**71** E.H. Schwartz, "Voice Clone Platform Debuts 'Voice Bank' for Future Speech Synthesis," Voicebot.ai, September 1, 2022. Available: https://voicebot.ai/2022/09/01/voice-clone-platform-debuts-voice-bank-for-future-speech-synthesis/

**72** E.H. Schwartz, "StoryFile CEO Brings AI Version of His Mother to Converse With Guests At Her Own Funeral," Voicebot.ai, August 16, 2022. Available: https://voicebot.ai/2022/08/16/storyfile-ceo-brings-ai-version-of-his-mother-to-converse-with-guests-at-her-own-funeral/

**73** A. Ramesh, "Zomato embraces deepfake technology and AI in personalized ads starring Hrithik Roshan," afaqs!, July 4, 2022. Available: https://www.afaqs.com/news/advertising/zomato-embraces-deepfake-technology-and-ai-in-personalised-ads-starring-hrithik-roshan

**74** K. Chitrakorn, "How deepfakes could change fashion advertising," Vogue Business, January 11, 2021. Available: https://www.voguebusiness.com/companies/how-deepfakes-could-change-fashion-advertising-influencer-marketing

**75** K. Chitrakorn, "How deepfakes could change fashion advertising," Vogue Business, January 11, 2021. Available: https://www.voguebusiness.com/companies/how-deepfakes-could-change-fashion-advertising-influencer-marketing

**76** "Deepfake Voice Technology: The Good. The Bad. The Future," EconoTimes, February 1, 2021. Available: https://www.econotimes.com/Deepfake-Voice-Technology-The-Good-The-Bad-The-Future-1601278

**77** K. Molenaar, "Discover The Top 12 Virtual Influencers for 2024 – Listed and Ranked!" Influencer Marketing Hub, October 7, 2024. Available: https://influencermarketinghub.com/virtual-influencers/

**78** "Shudu: The World's First Digital Supermodel," Virtual Humans. Available: https://www.virtualhumans.org/human/shudu

**79** M. Klein, "The Problematic Fakery Of Lil Miquela Explained—An Exploration Of Virtual Influencers and Realness," Forbes, November 17, 2020. Available: https://www.forbes.com/sites/mattklein/2020/11/17/the-problematic-fakery-of-lil-miquela-explained-an-exploration-of-virtual-influencers-and-realness/

**80** Christopher, "Virtual Influencers Speak Out in Response to the Murder of George Floyd," Virtual Humans, June 5, 2020. Available: https://www.virtualhumans.org/article/virtual-influencers-speak-out-in-response-to-the-murder-of-george-floyd

**81** F.J. McEvoy, "Deepfaking the Deceased: Is it Ever Okay?" Youth Data, January 23, 2021. Available: https://youthedata.com/2021/01/23/deepfaking-the-deceased-is-it-ever-okay/

**82** "Dali Lives," The Dalí Museum. Available: https://thedali.org/exhibit/dali-lives/

**83** JFK Unsilenced is an audio and video deepfake of former US-President Kennedy delivering the speech that he was to have given on the day of his murder. Similarly, a deepfake shows an altered version of Nixon's speech, prepared for the eventuality of a failed Apollo-11 moon landing. Deepfakes also power an immersive experience of Martin Luther King's famous "I Have a Dream" speech.

**84** N. Eisikovitz, „The slippery slope of using AI and deepfakes to bring history to life", The Conversation, November 2, 2021. Available: https://theconversation.com/the-slippery-slope-of-using-ai-and-deepfakes-to-bring-history-to-life-166464

**85** A.-C. Diaz, "Parkland victim Joaquin Oliver comes back to life in heartbreaking plea to voters," Ad Age, October 2, 2020. Available: https://adage.com/article/advertising/parkland-victim-joaquin-oliver-comes-back-life-heartbreaking-plea-voters/2285166

**86** S. Maheshwari, "Your Loved Ones, and Eerie Tom Cruise Videos, Reanimate Unease With Deepfakes," The New York Times, March 10, 2021. Available: https://www.nytimes.com/2021/03/10/technology/ancestor-deepfake-tom-cruise.html

**87** "Gallery: 'Spectre' Launches," Bill Posters, Available: https://billposters.ch/spectre-launch/

**88** S. Cole, "A Site Faking Jordan Peterson's Voice Shuts Down After Peterson Decries Deepfakes," Vice, August 26, 2019. Available: https://www.vice.com/en/article/not-jordan-peterson-voice-generator-shut-down-deepfakes

**89** "Welcome to Chechnya: Inside the Russian Republic's Deadly War on Gays," A Film by David France, Available: https://www.welcometochechnya.com/

**90** W. Knight, "AI deepfake satire from South Park creators tests the limits of media mimicry," MIT Technology Review, October 28, 2020. Available: https://www.technologyreview.com/2020/10/28/1011336/ai-deepfake-satire-from-south-park-creators/

**91** J.P. Meneses, "Deepfakes and the 2020 US elections: what (did not) happen," CECS, Portugal, Available: https://arxiv.org/pdf/2101.09092

**92** "Deepfake-Scholz verkündet AfD-Verbot," ZDF, November 27, 2023. Available: https://www.zdf.de/nachrichten/politik/aktion-gefaengnis-afd-verbot-100.html

**93** M. Reuter, "Kunst-Aktion für AfD-Verbot: Landgericht Berlin verbietet Kanzler-Deepfake," Netzpolitik.org, February 19, 2024. Available: https://netzpolitik.org/2024/kunst-aktion-fuer-afd-verbot-landgericht-berlin-verbietet-kanzler-deep-fake/

**94** "Deepfakery: A Livestream Talk Series and Exploration of Critical Questions," WITNESS. Available: https://lab.witness.org/deepfakery/

**95** Chowdhury R. "AI-fuelled election campaigns are here - where are the rules?" Nature. 2024 Apr;628(8007):237. doi: 10.1038/d41586-024-00995-9.

**96** N. Christopher, "We've Just Seen the First Use of Deepfakes in an Indian Election Campaign," Vice, February 18, 2020. Available: https://www.vice.com/en/article/the-first-use-of-deepfakes-in-indian-election-by-bjp/

**97** "Deepfake democracy: South Korean candidate goes virtual for votes," France 24, February 14, 2022. Available: https://www.france24.com/en/live-news/20220214-deepfake-democracy-south-korean-candidate-goes-virtual-for-votes

**98** H. Shin and H.Y. Yi, "South Korea candidates woo young voters with 'deepfakes,' hair insurance," Reuters, March 7, 2022. Available: https://www.reuters.com/world/asia-pacific/skorea-candidates-woo-young-voters-with-deepfakes-hair-insurance-2022-03-03/

**99** "Indien-Wahl als 'Testlabor' für Künstliche Intelligenz," Zeit Online, June 1, 2024. Available: https://www.zeit.de/news/2024-06/01/indien-wahl-als-testlabor-fuer-kuenstliche-intelligenz

**100** "Dance videos of Modi, rival turn up AI heat in India election," Channel News Asia, May 16, 2024. Available: https://www.channelnewsasia.com/asia/india-modi-election-ai-deepfake-dance-video-misinformation-4340341

**101** "Indien-Wahl als 'Testlabor' für Künstliche Intelligenz," Zeit Online, June 1, 2024. Available: https://www.zeit.de/news/2024-06/01/indien-wahl-als-testlabor-fuer-kuenstliche-intelligenz

**102** K. Hiebert, "Deepfakes Will Break Democracy," The Walrus, February 26, 2024. Available: https://thewalrus.ca/deepfakes-will-break-democracy/?utm_source=substack&utm_medium=email

**103** "Prabowo Subianto: Indonesia's 'cuddly grandpa' with a bloody past," BBC News, February 7, 2024. Available: https://www.bbc.com/news/world-asia-68028295

[104] S. Ray, "Imran Khan—Pakistan's Jailed Ex-Leader—Uses AI Deepfake To Address Online Election Rally," Forbes, December 18, 2023. Available: https://www.forbes.com/sites/siladityaray/2023/12/18/imran-khan-pakistans-jailed-ex-leader-uses-ai-deepfake-to-address-online-election-rally/

[105] M. Schueler, S. Romano, N. Stanusch, R.B. Cetin, S. Tabti, M. Faddoul, and I. Lilley, "Artificial Elections: Exposing the Use of Generative AI Imagery in the Political Campaigns of the 2024 French Elections," AI Forensics, 2024. Available: https://aiforensics.org/work/french-elections-2024

[106] J. Kerzig, "Erste Parteien machen in Sachsen Politik mit künstlicher Intelligenz - Experten warnen," Freie Presse, 2024. Available: https://www.freiepresse.de/nachrichten/sachsen/erste-parteien-machen-in-sachsen-politik-mit-kuenstlicher-intelligenz-experten-warnen-artikel13216449

[107] "Indien-Wahl als 'Testlabor' für Künstliche Intelligenz," Zeit Online, June 1, 2024. Available: https://www.zeit.de/news/2024-06-01/indien-wahl-als-testlabor-fuer-kuenstliche-intelligenz

[108] K. Muñoz, "Gegen den Strich: Künstliche Intelligenz und Wahlen," Internationale Politik, June 24, 2024. Available: https://internationalepolitik.de/de/gegen-den-strich-kuenstliche-intelligenz-und-wahlen

[109] "Legal Documents: Digital Services Act," European Commission. Available: https://commission.europa.eu/publications/legal-documents-digital-services-act_en

[110] "Artificial Intelligence Act," Available: https://artificialintelligenceact.eu/

[111] "The Brussels effect? Impact of the EU's AI Act – in the EU and beyond," Linklaters Tech Insights, July 11, 2024. Available: https://techinsights.linklaters.com/post/102jcir/the-brussels-effect-impact-of-the-eus-ai-act-in-the-eu-and-beyond

[112] "KI als neues Wahlkampfinstrument: Offene Flanke der KI-Regulierung?" Verfassungsblog, May 3, 2024. Available: https://verfassungsblog.de/ki-als-neues-wahlkampfinstrument/

[113] "Provisions on the Administration of Deep Synthesis Internet Information Services," China Law Translate, December 11, 2022. Available: https://www.chinalawtranslate.com/en/deep-synthesis/

[114] M. Holland, "Mit Deepfakes für den Sozialismus: China formuliert Regeln und Verbote," Heise Online, December 12, 2022. Available: https://www.heise.de/news/Deepfakes-China-formuliert-Regeln-und-verbietet-Einsatz-fuer-unerwuenschte-Zwecke-7373232.html

[115] Moreover, prosecuting perpetrators is difficult as many remain anonymous and are located abroad, and as law enforcement agencies often lack the technical resources, know-how, and personnell needed for prosecution.

[116] Davis Jr., Elliott. "These States Have Banned the Type of Deepfakes That Targeted Taylor Swift." U.S. News & World Report, January 30, 2024. Available: https://www.usnews.com/news/best-states/articles/2024-01-30/these-states-have-banned-the-type-of-deepfake-porn-that-targeted-taylor-swift

[117] PA Media, "Creating sexually explicit deepfake images to be made offence in UK," The Guardian, April 16, 2024. Available: https://www.theguardian.com/technology/2024/apr/16/creating-sexually-explicit-deepfake-images-to-be-made-offence-in-uk

[118] "South Korea to criminalize watching or possessing sexually explicit deepfakes," CNN (Story by Reuters), September 26, 2024. Available: https://edition.cnn.com/2024/09/26/asia/south-korea-deepfake-bill-passed-intl-hnk/index.html

[119] "Directive (EU) 2024/1385 of the European Parliament and of the Council of 14 May 2024 on combating violence against women and domestic violence," Official Journal of the European Union, May 24, 2024. Available: http://data.europa.eu/eli/dir/2024/1385/oj

[120] A. Engler, "Fighting deepfakes when detection fails," Brookings Institution, November 14, 2019. Available: https://www.brookings.edu/research/fighting-deepfakes-when-detection-fails/

[121] M. Anderson, "Encoding Images Against Use in Deepfake and Image Synthesis Systems," Unite.AI, December 9, 2022. Available: https://www.unite.ai/encoding-images-against-use-in-deepfake-and-image-synthesis-systems/

[122] T. Shane, E. Saltz and Leibowicz, C., „From deepfakes to TikTok filters: How do you label AI content?", First Draft, May 12, 2021. Available: https://www.niemanlab.org/2021/05/from-deepfakes-to-tiktok-filters-how-do-you-label-ai-content/

[123] R. Pfefferkorn, "The threat posed by deepfakes to marginalized communities," Brookings Institution, April 21, 2021. Available: https://www.brookings.edu/articles/the-threat-posed-by-deepfakes-to-marginalized-communities/

[124] N. Diakopoulos and D. Johnson, "Anticipating and Addressing the Ethical Implications of Deepfakes in the Context of Elections," New Media & Society, October 23, 2019. Available: http://dx.doi.org/10.2139/ssrn.3474183 or https://ssrn.com/abstract=3474183

[125] "Celebrity deepfakes are all over TikTok. Here's why they're becoming common – and how you can spot them," The Conversation, July 18, 2022. Available: https://theconversation.com/celebrity-deepfakes-are-all-over-tiktok-heres-why-theyre-becoming-common-and-how-you-can-spot-them-187079

[126] M. Pawelec, "Decent deepfakes? Professional deepfake developers' ethical considerations and their governance potential," AI Ethics, 2024. Available: https://doi.org/10.1007/s43681-024-00542-2

[127] B. Marr, "Can A Metaverse AI Win America's Got Talent? (And What That Means For The Industry)," Forbes, August 30, 2022. Available: https://www.forbes.com/sites/bernardmarr/2022/08/30/can-a-metaverse-ai-win-americas-got-talent-and-what-that-means-for-the-industry/

## Heinrich Böll Foundation Tel Aviv

The Heinrich Böll Foundation is an independent global think-and-do-tank for green visions. With its international network of 33 international offices, the foundation works in more than 60 countries. The foundation's work in Israel focuses on fostering democracy, promoting environmental sustainability, advancing gender equality, and promoting dialog and exchange of knowledge between public policy experts and institutions from Israel and Europe.

## Israel Public Policy Institute (IPPI)

The Israel Public Policy Institute (IPPI) is an independent think-and-do-tank and a multi-stakeholder dialog platform at the intersection of technology, society, and the environment. It collaborates with a global network of experts and partners that spans government, academia, civil society, and the private sector to foster international and interdisciplinary exchanges of ideas and best practices. Through its research, knowledge dissemination, networking, and public engagement, IPPI works to guide society's transition towards a sustainable future.