# The Age of Synthetic Realities: Challenges and Opportunities

JOÃO PHILLIPE CARDENUTO[1], JING YANG[1], RAFAEL PADILHA[1], RENJIE WAN[2], DANIEL MOREIRA[3], HAOLIANG LI[4], SHIQI WANG[5], FERNANDA ANDALÓ[1], SÉBASTIEN MARCEL[6,7] AND ANDERSON ROCHA[1]

*Synthetic realities are digital creations or augmentations that are contextually generated through the use of Artificial Intelligence (AI) methods, leveraging extensive amounts of data to construct new narratives or realities, regardless of the intent to deceive. In this paper, we delve into the concept of synthetic realities and their implications for Digital Forensics and society at large within the rapidly advancing field of AI. We highlight the crucial need for the development of forensic techniques capable of identifying harmful synthetic creations and distinguishing them from reality. This is especially important in scenarios involving the creation and dissemination of fake news, disinformation, and misinformation. Our focus extends to various forms of media, such as images, videos, audio, and text, as we examine how synthetic realities are crafted and explore approaches to detecting these malicious creations. Additionally, we shed light on the key research challenges that lie ahead in this area. This study is of paramount importance due to the rapid progress of AI generative techniques and their impact on the fundamental principles of Forensic Science.*

## I. INTRODUCTION

In the last decade, there has been a growing expectation that Artificial Intelligence (AI) companies and researchers would dedicate their efforts to integrating humans into the digital realm or the so-called Metaverse. However, while technologies like Augmented Reality (AR) and Virtual Reality (VR) have been topics of discussion for a long time, it is only recently that significant technological advancements have made it possible to materialize such systems.

This version of a *synthetic reality* has been technically discussed at least since the 90s [1], and this was the main vision of where AI and related technologies would take us. However, what came as a less expected outcome is that these very technologies would inundate our physical world with content and creations, profoundly transforming our interactions with the virtual realm and reshaping how we engage with one another.

This more complex notion of *synthetic reality* has been a topic of discussion by the greatest minds of our time [2]. On the one hand, some hold a positive perspective, recognizing the immense advantages it can bring in domains such as automation, healthcare, and innovation. On the other, a group expresses concerns about the potential perils posed by AI, such as the generation of propaganda and untruth.

They even advocate for temporary halts in AI experimentation within laboratories [3] to allow time for legal and ethical considerations to align with the pace of progress.

We can adopt a more pragmatic perspective and carefully embrace this emerging paradigm. This entails rapidly adapting ourselves and our societies and understanding and revitalizing our scientific endeavors. Hence, we redefine the term "synthetic realities" herein as any contextual digital creation or augmentation enabled by artificial intelligence methods. These techniques/models draw upon massive amounts of data leading to a new "reality" or narrative regardless of its intention to deceive the individual interacting with it. When the synthetic creation harms individuals, minorities, human rights, or the rule of law, it is paramount to devise forensic techniques to pinpoint such creations and separate what is real from what is synthetic. As an example, consider the creation of a fake news piece. Someone could fabricate a story from scratch using a chatbot, illustrate it with a synthetic image and a video, and then broadcast it to the world as if it were real via social media.

Consequently, Forensic Science has been continually adapting to these evolving circumstances. Rooted in the foundational principle that "every contact leaves a trace", coined by researcher Edmond Locard [4], Forensic Science asserts that every interaction between individuals, objects, and places leaves behind a trail of evidence. While this concept was initially centered around physical traces like fingerprints, footprints, and blood, it has recently expanded to encompass digital counterparts such as photos, audio, video, and social media posts [5].

The revolving question around digital evidence is: "Are they fake or not?", as manipulating these multimedia assets can be quickly done with simple and inexpensive tools. Moreover, credible manipulations can be used to fabricate more believable multimedia stories. Research in Multimedia Forensics has yielded important approaches to

[1] Artificial Intelligence Lab., Recod.ai, Institute of Computing, University of Campinas, Campinas, SP, Brazil.
[2] Department of Computer Science, Hong Kong Baptist University, Hong Kong.
[3] Department of Computer Science, Loyola University Chicago, Chicago, IL, USA.
[4] Department of Electrical Engineering, City University of Hong Kong, Hong Kong.
[5] Department of Computer Science, City University of Hong Kong, Hong Kong.
[6] Idiap Research Institute, Martigny, Switzerland.
[7] University of Lausanne, Lausanne, Switzerland.

**Corresponding author:** João Phillipe Cardenuto
Email: phillipe.cardenuto@ic.unicamp.br

detecting altered media [6]. More specifically, progress has been made in analyzing digital media (image, video, and audio), identifying manipulations [7–11], tracing provenance [12–16], and establishing links with other digital evidence [17–19].

However, the emerging concept of *synthetic reality* paints an even more unsettling scenario: around 90% of the digital content will be synthetic in the upcoming years, meaning that almost all content will be generated synthetically by definition [20]. The distinction between what is genuine and what is fake takes on a new meaning. This phenomenon becomes evident in various domains, including movies [21], social media [22, 23], marketing [24], and education [25].

Notably, companies are now exploring the adoption of AI-generated models to promote their products [24] or employing AI to simulate eye contact in video conferencing software, enhancing the sense of connection during remote interactions [23]. Schools worldwide are banning chatbots and the use of generative AI on their networks in response to concerns about students submitting unauthentic and potentially plagiarized work [25]. The examples are many when thinking about how AI is shaping our reality.

Therefore, Forensic Science has to adapt yet again to this new reality. To expose synthetic content and tell apart malicious from harmless manipulations or even positive creations, there is one new key element: context. Contextual information can be leveraged to understand the semantics behind media objects. This can empower fact-checking solutions that mitigate the effect of falsified news, misinformation, and false political propaganda. To that end, and following a cognitive science interpretation, we can view digital objects through three perspectives: technological artifacts, sources of information, and platforms to convey ideas.

Traditional forensic techniques thus far have primarily focused on the first perspective. Analysts examined an asset as a digital signal and aimed to detect any possible artifact related to pixel-level or physical-level inconsistencies (e.g., concerning compression, sensor noise, illumination, shadows) to establish its authenticity. The second perspective pertains to standard fact-checking procedures going beyond multimedia forensics. When treating an object as a source of information, it is essential to identify (or know) the acquisition device and to identify the time and location where it was produced as basic steps towards a fact-checking effort. The third perspective considers a digital object as a platform for conveying ideas. Determining the intentional goal of the asset leads to answering the question of why something happened. The answer to this question and the result of forensic analyses from the other two perspectives can reveal the ultimate goal of the falsified information. For example, it can help identify if there is an ongoing campaign to bias public opinion, influence the mood of a social group, or even incite a group to articulate plans for violent acts.

The latest advancements in AI have compelled forensic techniques to navigate the intricacies between these perspectives. Taking this into consideration, we focus on studying synthetic realities in different forms of media: images, videos, audio, and text. In the remainder of this paper, we discuss how synthetic media is created, considering each of these modalities, and the implications for Digital Forensics when such creations intend to harm third parties in different ways. We explore how to detect such malicious creations and pinpoint key research challenges that lie ahead. This is particularly significant due to the remarkable progress of AI generative techniques in generating realistic content and effectively concealing the typical artifacts left behind during the creation process. Each new generative method aims at creating ever-more-believable realities, thus directly colliding with Locard's principle, the cornerstone of forensics.

## II. SYNTHETIC IMAGES

The proliferation of sophisticated synthetic and manipulated media has captured people's attention worldwide. For forensic researchers, this surge in synthetic reality evokes the daunting scenario reminiscent of the early days of Digital Forensics, where image editing was recognized as a powerful tool capable of altering reality [26]. While traditional manual image-editing software like Photoshop and GIMP continue to improve, they are being overshadowed (or enhanced, in the particular case of Photoshop [27]) by the emergence of powerful AI-based generative techniques. Today, it has become remarkably effortless to transform a simple concept or idea into a realistic image, with no requirement for drawing or painting skills to produce stunning, high-quality results.

Amid this rapidly evolving landscape lies generative models. Generative images have outstanding widespread applications in entertainment, as reviving legendary artists [28]; healthcare, as aiding surgeons in developing new abilities [29]; and accessible tools, as serving people with disabilities [30]. However, many other harmful uses have been reported, such as nonconsensual DeepFake porn [31], misinformation generation [32], and sophisticated types of scams [33]. Scientific integrity researchers are also concerned that such technology would create fraudulent synthetic images in science [34, 35] and the medical area [36, 37], in particular.

Given the alarming potential for misuse of generative models in creating synthetic realities, this section delves into state-of-the-art generative models and the detection of AI-generated images, providing perspectives on the future of synthetic images.

### A) Image Synthesis

The accelerated research on generative approaches in recent years gave birth to a plethora of AI models and techniques that are developed and open-sourced to the community. They are coupled with easy-to-use environments and applications [38, 39], allowing users to freely

(a) Text-to-Image

(b) Inpainting

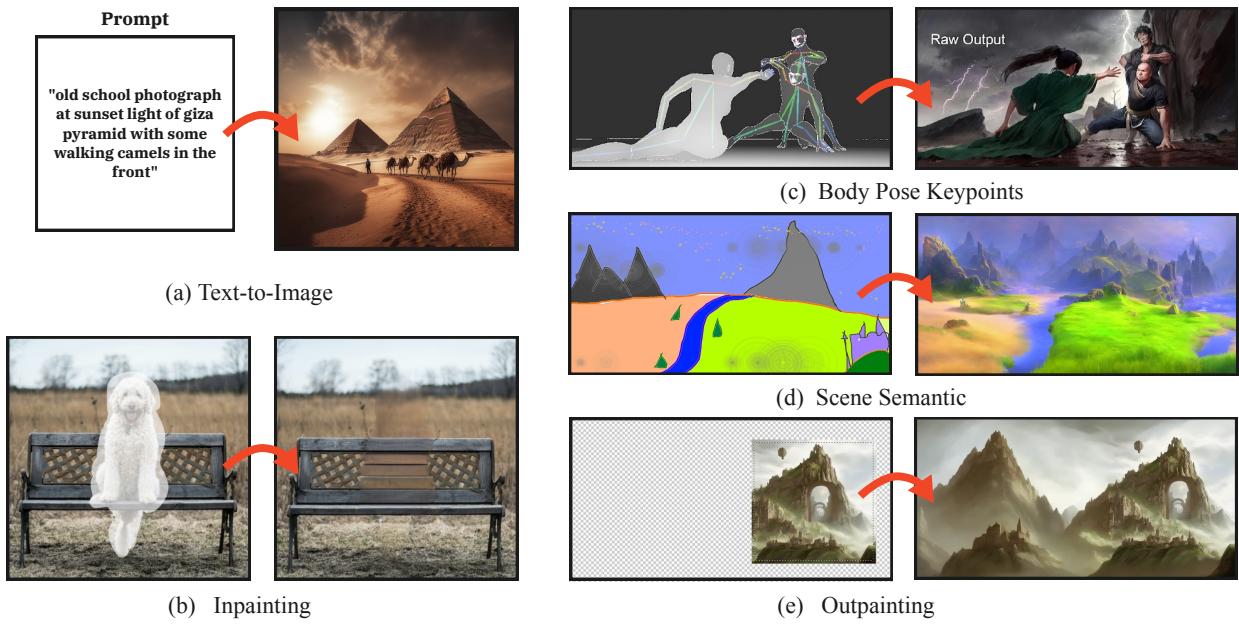(c) Body Pose Keypoints

(d) Scene Semantic

(e) Outpainting

Fig. 1.: Examples of image synthesis conditioning. The input modality and expected output define the type of generation task performed by the model and help express desired characteristics in the synthesized creation. Examples created with and/or reproduced from [38, 43, 44].

explore and share their creations. The increased accessibility further expanded the hype in image synthesis, fostering novel use cases and commercial applications that range from *outpainting* famous art pieces [40] (i.e., synthetically extending the borders of an image) to designing political campaign ads [41]. With the increased interest in the topic from research and industry communities and the rapid development of techniques, one can safely assume that not all synthetic images are born equally. We can categorize existing approaches by how the generation task is conditioned and what family of AI models they rely on.

The generation task defines the goal of the method and, consequently, how it learns to map the expected input to a synthetic output image. Additionally, the input data modality conditions the generation process into expressing particular visual concepts and characteristics desired by the user [42]. The most common tasks fall into *text-to-image* or *image-to-image* generation. Popularized by recent applications such as MidJourney [38] and DreamStudio [39], *text-to-image* generation involves a natural language prompt describing the desired image. This often includes the object or concepts that should be created, the desired artistic style, and the feeling the composition should convey. Whereas, in *image-to-image* generation, guidance may come as visual information, such as a picture, semantic map, or body pose keypoints. These may aid the generation process with information often difficult to express by natural language prompts, such as the relative positioning of the elements. We show examples of conditioning tasks and modalities in Figure 1.

The families of AI models utilized in image generation techniques comprise several types of architectures, with the most common being Generative Adversarial Networks (GANs), Diffusion models, and Variational Autoencoders (VAE).

**GANs** [45, 46] are built on two components: a generator and a discriminator. The generator learns the underlying distribution of real examples to generate new data, while the discriminator decides whether the input is from the real data space. Through an adversarial training process, the discriminator learns to identify synthetic images, while the generator progressively improves its ability to produce high-quality images that can deceive its counterpart. Numerous GAN variants have been developed in recent years to enhance the performance and stability of image generation [47–50]. Among them, StyleGAN [51–53] allowed for intuitive control over the generated image attributes by modulating the convolutional kernels at different levels of the generator, instead of directly controlling the network input. Its successor, StyleGAN-T [54], builds upon its architecture for text-guided image synthesis. It leverages Contrastive Language-Image Pre-Training [55] (CLIP), a powerful text encoder that aligns textual descriptions with corresponding images. Other approaches [56, 57] follow a similar path, relying on CLIP to integrate natural language understanding into the image synthesis process.

As an alternative to the min-max optimization game of GANs, **Diffusion models** [58, 59] are trained to revert a stochastic diffusion process that progressively adds noise

to a target image. To generate new images, the model iteratively denoises the perturbed image at each step, until a high-quality picture is reconstructed. By relying on a deterministic denoising function instead of adversarial learning, their training is more stable and easier to control than GANs. On the other hand, diffusion models rely on multiple network passes to reconstruct samples, constituting a considerably more computationally expensive method than adversarial networks. To improve efficiency, Stable Diffusion [60] operates on compressed latent representations instead of pixel space, mapping the denoising function to smaller manifolds. When considering textual prompts, Imagen [61] leverages text encoders, such as CLIP, to combine them with multiple cascaded diffusion models to generate high-resolution outputs from text. Similarly, Ramesh et al. [62] train a diffusion decoder that produces images from CLIP embeddings extracted from textual prompts.

Another prominent family of models is **Variational Autoencoders** (VAE) [63]. Autoencoders follow an encoder-decoder architecture that projects the input data into a low-dimensional latent space and learns to reconstruct the original input from it. VAE, in turn, extends upon autoencoders by adding a probabilistic component to the latent representation. Instead of learning a deterministic encoding for each input, the network learns the parameters of a probabilistic distribution that models the latent space. By doing so, it can sample from the latent space distribution to generate new samples. Constraining the low-dimensional latent space further, VQ-VAE [64, 65] uses vector quantization to learn discrete latent variables, which improves the interpretability of the learned concepts and allows for easily manipulating them when generating new compositions. Building on top of the previous technique, DALL-E [66] and CogView [67] address text-to-image generation by combining the rich representations learned by variations of VQ-VAE with the predictive capabilities of Transformers [68]. Both approaches use Transformer modules to predict the best image tokens from the VQ-VAE codebook, given a textual token and previously selected visual tokens. This results in coherent and semantically meaningful synthesized creations.

Each of these AI model families has its pros and cons. GANs excel at generating sharp and visually compelling images, but they may suffer from mode collapse and training instability issues when used in large and diverse datasets [46]. Diffusion models offer a powerful approach to generating high-quality visual data, but they can be computationally expensive due to their iterative nature. VAEs provide a more straightforward training process with a clear optimization objective, but they may generate less sharp images than GANs and diffusion models. Nonetheless, all of them made significant advancements in the field of image generation, enabling the synthesis of realistic and diverse pictures. As the interest in this area increases, more advances will come, and the realism gap between real and synthetic data will shorten to the point that distinguishing between them will be challenging. This poses numerous problems in assessing the reliability and authenticity of visual content in an increasingly digital world. With this in mind, we discuss existing forensic approaches that may help to identify synthetic creations in the next section.

## B) Synthetic Images Detection

In contrast to old-fashioned types of image manipulation, modern synthetic images take their realism to higher standards. Figure 2 compares classic Digital Forensics cases with those created by generative models, showcasing the remarkable advancement achieved. The level of refinement and potential harm associated with these synthetic images raise concerns about the capability of Digital Forensics to identify such content. However, we anticipate that **Locard's exchange principle still holds for synthetic imagery**, with forensic traces taking the form of visual inconsistencies and artifacts left by the generation process. Nevertheless, as these models continue to evolve and new forms of counter-forensics attacks emerge, the validity of this claim may be challenged.

In this section, our analysis focuses on examining possible traces left by GANs and other generative models. We categorize our study based on the types of forensic evidence utilized for detection, namely **Visual Artifacts** and **Noise Fingerprints**.

### 1) Visual Artifacts

Despite the impressive realism of cutting-edge synthetic images, a closer look reveals various aberrational results and visual inconsistencies. Borji [69] has presented several image failures that can occur when generating synthetic content, even with recent generative models such as DALL-E 2 [62], Midjourney [38], and StableDiffusion [44]. These failures may occur in the background, reflections, lighting, shadows, text, body parts, and objects, as depicted by Figure 3. This figure illustrates clues that have been explored by Digital Forensics researchers to detect synthetic content.

Farid, for instance, analyzed the 3D illumination [70] and geometric [71] consistency of structures and objects in a photograph generated by state-of-the-art generative models. His analysis employed classic digital forensics techniques for on-scene illumination and 3D geometric analysis, similar to those used in previous classic forensic works [72–74]. By doing so, Farid showed that while the local structures in the photo may be globally consistent, they exhibit local inconsistencies that serve as valuable clues for forensic analysts.

However, as generative models continue to evolve, it is expected that the visual inconsistencies and artifacts observed in synthetic images will eventually become rarer or imperceptible, as demonstrated in the case of the AI-synthetic faces examined by Nightingale and Farid [75]. This situation asks for other strategies that rely on other types of forensic clues, such as the noise left by the generation processes.

(a)   Image Manipulation



(b)   Image Synthesis

Fig. 2.: Classic manual image manipulation versus modern image synthesis. The image on the left (a) represents a well-known case of state-level manual image manipulation, which was misleadingly published by multiple news websites as genuine in 2008. In contrast, we generated the image on the right (b) with MidJourney by using a prompt as simple as "missile test". Cases like the former one currently represent an even greater challenge to authenticity verification.



(a) Human Body Exception.



(b)   Object Inconsistency.



(c)   Text Failure.



(d) Reflection Inconsistency.

Fig. 3.: Examples of generated image inconsistencies. All images were generated with version 5.1 of the Midjourney model – the latest one released at the writing of this article. In (a), an unnatural synthetic hand with six fingers. In (b), a synthetic wheelchair with inconsistent design; the seat orientation does not match the wheels' position. In (c), a synthetic billboard with text that makes no sense and presents aberrant letters. In (d), a synthetic paisage with cloudy skies and mountains by a lake; the highlighted cloud is not congruently reflected on the lake's surface. To generate these images, we used the following prompts: (a) "lady's hand with a ring on it", (b) "wheelchair in a hospital", (c) "outdoor sign with a religious statement on it", and (d) "realistic photo; mountains with a lake at the bottom".

## 2)  Noise Fingerprint

As image synthesis is rapidly improving, it is crucial to employ a variety of alternative detectors that explicitly exploit different characteristics of synthetic images. Therefore, as an alternative to visual artifacts, noise-based detectors have been a promising path to expose synthetic content.

In this direction, Marra et al. [76] investigated statistics-based techniques to detect potential noise fingerprints left by GANs on their generated content. By utilizing photo response non-uniformity (PRNU) analysis, similar to camera attribution methods, they discovered a correlation between residual noise patterns and specific GAN models. Furthermore, Marra et al. demonstrated the feasibility of differentiating between distinct GAN models used for image synthesis through residual noise analysis, enabling GAN model attribution.

In a similar study, Mandelli et al. [77] revealed that comparable residual noise patterns could be leveraged to identify GAN-generated scientific images, indicating the potential extension of this approach to other image types beyond natural images. These findings emphasize the applicability of residual noise analysis in detecting and identifying synthetic images, contributing to the field of Digital Forensics. Noise signatures allied with visual clues were explored in tandem by Kong et al. [78], showing that combining different evidence might be the way forward in dealing with the challenges of synthetic realities detection.

In a more recent investigation, Corvi et al. [79] examined the presence of fingerprints left by state-of-the-art generative models, including GAN-based and Diffusion-based models. Their findings reveal that **no generative model appears to be completely artifact-free** at present. Both GAN-generated and Diffusion-generated images exhibit anomalous periodic patterns in the Fourier spatial domain.

However, as highlighted by Gragnaniello et al. [80], such artifacts may be challenging to detect when post-processing operations are applied, such as image resizing and compression. These operations are frequently employed on social media platforms to save storage and speed up sharing, further complicating the identification of these visual irregularities. Besides investigating how post-processing operations impact synthetic image detectors, researchers have also identified the potential for synthetic image detectors to be deceived through counter-forensics attacks.

### 3) Counter Forensics

In a recent study that challenges noise-based detection methods, Osakabe et al. [81] developed a GAN model that can generate images without "checkerboard artifacts", a specific type of artifact in the Fourier domain that is common in synthetic images. They achieved this by incorporating a fixed convolutional layer into every upsampling and downsampling layer of the GAN architecture. Remarkably, their model successfully fooled a detector that previously identified fake images with 92% accuracy, reducing the accuracy to a mere 12%.

Similarly, Cozzolino et al. [82] demonstrated that synthetic image detectors could be fooled by transferring residual noise fingerprints from real cameras onto GAN-generated images. This process produces a spoofed image that can avoid accurate GAN detectors and camera-model identifiers, causing the image to be misattributed as originating from the transferred camera model.

As counter-forensic attacks indicate, synthetic image detection presents multiple research challenges to ensure media integrity and prevent images from being used by a malicious actor.

## C) Challenges and Directions

As with classic edited images, creating synthetic realities using cutting-edge generative models has sparked ethical debates and raised concerns about their use. Once again, Digital Forensics plays a crucial role in this debate by investigating ways to detect the traces of artificial intelligence techniques left in these images. While most synthetic images can be identified through a close look into scene inconsistencies and object aberrations, as depicted in Figure 3, the advancement of image synthesis will inevitably render these visual incongruences invisible to the naked eye.

Consequently, researchers also rely on new types of fingerprints inherently left by the generation process, such as specific patterns on the Fourier spatial spectrum and residual noise analysis. Some of these artifacts can be compared to the PRNU noise left by camera sensors. They can aid forensic analysis not only in detecting fake images but also in identifying the specific generative model used to render them (source attribution). However, traces alone provide a vulnerable target for synthetic image detection, as they can easily be manipulated to deceive accurate fake image detectors. Such attacks may involve common post-processing operations or sophisticated techniques like camera noise transference. Therefore, it is imperative for forensic researchers to develop robust techniques capable of detecting and distinguishing these artifacts, even in the presence of common post-processing operations or more sophisticated attacks.

A more challenging and socially responsible aspect of forensics involves preventing the harmful applications of synthetic images. Given the ease with which such content can be created and shared, it is essential for forensics researchers to design traceable techniques whose synthetic images can be readily distinguished. Traceable evidence would assist analysts in swiftly identifying the source and author of such content, thereby preventing its widespread dissemination. In this vein, researchers have developed deep learning-based watermarking approaches [83–85] to identify synthesized content. These methods use encoder layers to imbue watermark information in the image pixels without perceptually altering its content. For successfully marking creations, researchers aim to be robust to most online alterations, such as compression, cropping, and intensity changes. Unfortunately, most approaches act from the generator's perspective, either adapting existing models or adding external components in the generation process to enable watermarking. This might be viable for well-established and commercial applications (e.g., Midjourney [38]) but will be hardly used with open-sourced models that are trained and distributed by the community.

## III. SYNTHETIC VIDEOS

The ability to generate realistic and useful videos holds immense value across various application domains such as entertainment, virtual reality, and education [86, 87]. Undoubtedly, video generation techniques have made significant positive contributions in these domains. The advancements have opened up new avenues for creativity, synthetic realities, and immersive experiences. However, it is essential to acknowledge that along with their benefits, these techniques also raise potential security concerns. Synthesize realistic videos can be exploited for malicious purposes, such as financial fraud and the dissemination of fake news. Consequently, ensuring the integrity and authenticity of digital content becomes increasingly critical.

Generally speaking, video synthesis can be divided into video generation and text-to-video synthesis. Previous methods for video generation mainly employ GANs [87–89] and VAEs [90, 91] to generate videos. But with the advent of diffusion models, recent methods explore them to generate more realistic videos [62, 92–95]. On the other hand, text-to-video synthesis incorporates text information to guide the model in generating video content that is responsive to specific demands. Analogously, previous text-to-video synthesis methods have resorted to

GANs [96, 97] and diffusion models [98–102], achieving exceptional generation quality in terms of fidelity, resolution, and temporal consistency.

In this section, we provide a review of the methods for video generation, text-to-video synthesis, synthetic video detection and discuss their challenges. We further outline possible future research directions for synthetic video generation and detection techniques.

## A) Video Generation

In the pursuit of advancing video synthesis, previous research has extensively explored diverse generative models, including Generative Adversarial Networks (GANs) [87, 103], autoregressive models [90, 104], and implicit neural representations [105, 106]. However, recent attention has been drawn to the exceptional achievements of diffusion models in visual data synthesis. Several notable works propose outstanding video generation methods and investigate their practical applications [92, 93, 107, 108]. For instance, a pioneering work on diffusion video generation [109] primarily focuses on network architecture modifications to extend image synthesis to video. The 3D U-Net is adopted [110] and achieves outstanding generation results in two cases, including unconditional and text-conditional video generation. For longer video generation, they apply an autoregressive approach, where subsequent video segments are conditioned on the preceding ones. Another example of a diffusion video generation work [94] adopts frame-by-frame video generation models. To evaluate different prediction strategies, the authors conduct an ablation study to determine whether predicting the residual of the next frame yields superior results compared to predicting the actual frame.

Furthermore, Hoppe *et al*. [107] introduced the Random Mask Video Diffusion (RaMViD) technique, which can be utilized for both video generation and infilling tasks. The unmasked frames are used to enforce conditions on the diffusion process, while the masked frames undergo diffusion through the forward process. By employing this training strategy, RaMViD demonstrates outstanding video generation quality. These recent advancements in diffusion-based video generation highlight the potential of this family of models to push the boundaries of video synthesis, addressing the challenges of generating realistic and diverse video content in synthetic realities.

## B) Text-to-video Synthesis

Text-to-video models are highly data-hungry, which require massive amounts of data to learn caption relatedness, frame photorealism, and temporal dynamics [111]. However, video data resources are comparatively more limited in terms of style, volume, and quality. This scarcity of video data poses significant challenges for training text-to-video generation models. To overcome these challenges, additional controls are often incorporated to enhance the responsiveness of generated videos to user demands [97, 112, 113].

Early text-to-video generation models heavily relied on convolutional GAN models combined with Recurrent Neural Networks (RNNs) to capture temporal dynamics [96, 97]. Despite the introduction of complex architectures and auxiliary losses, GAN-based models exhibit limitations in generating videos beyond simplistic scenes involving digit movements or close-up actions. To that end, recent advancements in the field have aimed to extend text-to-video generation to more diverse domains using large-scale transformers [114] or diffusion models [98]. These approaches provide promising directions for generating more complex and realistic video content by leveraging the expressive power of these advanced network architectures.

However, modeling high-dimensional videos and addressing the scarcity of text-video datasets present considerable challenges in training text-to-video generation models from scratch. To tackle this issue, most approaches adopt a transfer learning paradigm, leveraging pre-trained text-to-image models to acquire knowledge and improve performance. For instance, CogVideo [115] builds upon the pre-trained text-to-image model CogView2 [116], while Imagen Video [98] and Phenaki [117] employ joint image-video training techniques to leverage pre-existing visual representations. In contrast, Make-A-Video [95] focuses on learning motion solely from video data, reducing the reliance on text-video pairs for training.

Another key consideration in video synthesis is the high computational cost associated with generating high-quality videos. To mitigate this issue, latent diffusion has emerged as a popular technique for video generation, as it offers a computationally efficient alternative [92, 99–101]. Various powerful but computational-efficient methods, such as MagicVideo [118], which introduces a simple adaptor after the 2D convolution layer, and Latent-Shift [99], which incorporates a parameter-free temporal shift module, have successfully utilized latent diffusion for video synthesis. Additionally, PDVM [119] adopts a novel approach of projecting the 3D video latent space into three 2D image-like latent spaces, further optimizing the computational cost of the video generation process.

Despite the active research in text-to-video generation, existing studies have predominantly overlooked the interplay and intrinsic correlation between spatial and temporal modules. These modules play crucial roles in understanding the complex dynamics of videos and ensuring coherent and realistic video generation.

## C) Synthetic Video Detection

Again, detecting synthetic videos relies on the fingerprints left by generative models, which have been explored extensively in the context of synthetic image detection. Existing research in synthetic image detection has shown promise by identifying inconsistencies in illumination and geometric structure [70, 71], as well as specific noise patterns in

the Fourier domain [79]. However, the extension of image-based detection techniques to videos is still in its early stages.

One straightforward approach for extending image-based detection to videos is through frame-level voting, where each frame is individually analyzed and classified as real or synthetic. However, exploiting temporal information, such as temporal coherence, presents a significant challenge. The temporal domain contains valuable cues that can aid in distinguishing synthetic videos from real ones. For instance, temporal coherence refers to the consistent motion and flow of objects across frames in a real video. Detecting such temporal fingerprints could provide valuable insights into the authenticity of a video. The temporal coherence, which is a significant challenge in the field of video synthesis [120], shall also be vital for synthetic video detection. Consequently, exploiting the temporal inconsistency can be utilized to identify generated videos. Currently, there is no existing method specifically designed for detecting video synthesis. However, there have been relevant work on detecting Deepfake videos using generic architectures. For instance, 3DCNN [121], RNN [122], LSTM [123], and temporal transformer [124] have been widely employed for video-level Deepfake detection. Additionally, some studies focus on detecting deepfake videos by examining specific temporal artifacts like lip movement [125], rPPG artifacts [126], and head pose inconsistency [127]. Thus, how to explore the temporal artifacts (*e.g.*, unnatural motions) in generated videos to detect video synthesis is an intriguing task.

This temporal aspect of video analysis introduces additional complexities compared to image analysis. Generated video content can exhibit various artifacts due to the processing techniques employed, including both handcrafted designs and deep neural networks. These artifacts, including blur, compression artifacts, and noise, can be intentionally or unintentionally injected during the video generation process. Consequently, these artifacts may pose significant obstacles to the detection of synthetic videos, particularly when the detection model is trained solely on high-quality video data.

To address this challenge, detection models need to be highly generalized to handle data from different domains. The models must be capable of recognizing and adapting to various levels of quality, distortion types, and content sources. This requirement calls for the application of domain generalization techniques, which enable the model to generalize well beyond the training data distribution. By training the model on a diverse range of video data, encompassing different quality levels, distortion types, and content sources, the detection system can become more robust and effective in identifying synthetic videos across a wide range of scenarios.

Overall, the extension of image-based detection techniques to videos presents a meaningful yet challenging direction for research. Leveraging temporal fingerprints and addressing the presence of artifacts in generated videos require novel approaches and further exploration. Developing detection models that can effectively analyze and distinguish synthetic videos while being adaptable to various domains will play a crucial role in combating the increasing threat of synthetic videos in today's digital landscape.

## D) Challenges and Directions

Recent advancements in diffusion models have revolutionized text-to-video synthesis, achieving remarkable capabilities that surpass previous state-of-the-art approaches and deliver unprecedented generative performance. This breakthrough has significantly enhanced the quality and fidelity of generated videos. However, as we delve deeper into this domain, it becomes evident that there is a need for further research and development.

Despite the demonstrated success in generated image content in the past few years, video generation is still in its infancy. As we have discussed, the challenges of video generation mainly lie in the following three aspects: (1) lacking large-scale, diverse, and in-the-wild video datasets; (2) demanding computational costs; and (3) unstable in synthesizing coherent content from both spatial and temporal perspectives.

In the realm of synthetic video detection, existing methods predominantly focus on identifying face forgery, where the manipulation targets primarily involve facial features and expressions. However, to address the evolving landscape of synthetic videos and their potential threats, it is imperative to explore detection techniques that encompass a broader range of scenes and contexts. Detecting synthetic videos with diverse scenes, objects, and backgrounds poses an interesting avenue for future advancements in forensic research. By expanding the scope of detection techniques, we can develop robust and comprehensive methods that effectively identify and mitigate the risks associated with the increasing sophistication of synthetic videos.

## IV. SYNTHETIC AUDIO

Synthetic realities in audio are an emerging technology transforming how we experience sound. From augmented and virtual reality to interactive audio installations, synthetic realities offer a new dimension to our auditory senses. These immersive audio experiences create a simulated environment that can transport listeners to different worlds, trigger emotions, and enhance storytelling. With the advancements in audio technology and the increasing demand for immersive experiences, synthetic realities are poised to revolutionize the entertainment, gaming, and education industries. While the rise of synthetic audio technology has brought about significant benefits to various fields, it also presents a considerable threat to the integrity of our society. One of the most concerning implications is the potential misuse of this technology by malicious actors, who can exploit it for nefarious purposes such as telecommunication fraud. Cai *et al.* [128] have highlighted the

dangers of using generative models to create fake audio that can deceive individuals and organizations, leading to financial losses and reputational damage. The use of synthetic audio in such fraudulent activities underscores the urgent need for developing robust and reliable methods for detecting and mitigating the harms of this technology. In this section, we propose to survey synthetic realities in audio and dive into the possibilities and challenges of this rapidly evolving field.

## A) Synthetic Audio Generation

Audio synthesis is a vital and rapidly evolving research area with a wide range of applications, including text-to-speech (TTS), speech enhancement, voice conversion, and binaural audio synthesis. In the field of TTS, previous works have extensively utilized deep learning-based architectures such as WaveNet [129] and Clarinet [130], as well as transformer models like FastSpeech [131] and Neural TTS [132], and variational autoencoder (VAE) approaches such as MultiSpeech [133] and Hierarchical VAE [134]. Recently, diffusion models have gained prominence in addressing TTS problems, with notable contributions from WaveGrad [135], DiffWave [136], Gradient Flow [137], and Diffusion TTS [138].

Speech enhancement techniques aim to improve speech recognition system performance by mitigating the impacts of ambient noise. The advancement of generative models has led to the development of various approaches for speech enhancement. These include GAN-based methods like MetricGAN [139] and Speech Enhancement GAN [140], as well as diffusion-based models such as Storm [141], Conditional Diffusion [142], and Cold Filter [143]. These models have exhibited general and robust speech enhancement performance.

Voice conversion, another critical task in speech synthesis, aims to transform the voice of one speaker into that of another. Different approaches have been explored for voice conversion, including transformer-based models (VoiceFilter [144]), GAN-based methods (CycleGAN-VC [145], VoiceGAN [146]), and VAE-based techniques (ACVAE [147] and Neural Voice Cloning [148]). These methods facilitate the manipulation of speaker characteristics while maintaining the linguistic content of the speech.

Lastly, binaural audio synthesis [149, 150] focuses on transforming mono audio signals into binaural audio, which enables accurate sound localization and immersive auditory experiences. By simulating the perception of sound through two ears, binaural audio synthesis contributes to creating a more realistic and interactive auditory environment.

Overall, the continuous advancements in deep learning, generative models, and various synthesis techniques have significantly expanded the possibilities and applications of audio synthesis, enhancing the quality, naturalness, and versatility of synthesized speech and audio.

## B) Synthetic Audio Detection

Existing methods for synthetic audio detection can be categorized into two different streams: feature-based and image-based methods [151]. Feature-based approaches describe the audio through signal features such as Mel frequency cepstral coefficient (MFCC) and constant Q cepstral coefficient (CQCC) [152, 153]. These features are then fed into typical classifiers (e.g., support vector machines) [154] and deep neural networks [155]), which are trained to detect synthetic audio. Image-based methods, on the other hand, utilize either spectrogram images [156, 157] computed from the audio signal and use them as the input for deep neural networks to extract discriminative information for synthetic audio detection. The aforementioned techniques have also been widely applied to deepfake detection related to synthetic audio [158] (more detail can be found in Section VII).

## C) Challenges and Directions

Despite the progress made in synthetic audio detection, there are still challenges to overcome. One of the critical challenges is the availability of large datasets of synthetic audio that can be used to train detection models effectively [159]. Generating a large dataset of synthetic audio can be time-consuming and resource-intensive. Additionally, the increasing complexity and sophistication of synthetic audio algorithms may require more advanced detection methods that can keep up with these advancements.

While there are challenges to overcome, such as the availability of large datasets and the need for more advanced detection methods, recent research has shown promising results in developing more effective detection techniques. As synthetic audio technology advances, it is essential to continue developing and improving detection methods to prevent the misuse of deepfake audio and ensure that this technology is used safely and responsibly.

## V. SYNTHETIC TEXT

Large Language Models (LLMs) have revolutionized artificial intelligence, marking a significant milestone in the field. Since the rise of GPT models [160, 161], competitors from other companies like Google and Microsoft have also developed their own LLM, including Gopher [162], GLaM [163], LaMDA [164], Megatron-Turing NLG [165], Chinchilla [166], PaLM [167], BloombergGPT [168]. Some open-source models have also been released, including Jurassic-1 [169], GPT-J [170], GPT-NeoX [171], BLOOM [172], OPT [173], LLaMA [174], and OpenAssistant [175]. Although there are many variants of LLMs models, most of them have the same architecture of the GPT-2 model [160], which is an auto-regressive model (decoder-only transformer) that generates the next word based on previous inputs.

Figure 4 illustrates how these models have reached a quality comparable to humans, even in highly complex

| |
|---|
| What kind of thoughts now, do you carry<br>In your travels day by day<br>Are they bright and lofty visions,<br>Or neglected, gone astray?<br><br>Matters not how great in fancy,<br>Or what deeds of skill you've wrought;<br>Man, though high may be his station,<br>Is no better than his thoughts.<br><br>Catch your thoughts and hold them tightly,<br>Let each one an honor be;<br>Purge them, scourge them, burnish brightly,<br>Then in love set each one free. |

(a) Poem "Thoughts" - Human Version.

| |
|---|
| In the realm of silent contemplation,<br>Thoughts weave their intricate creation.<br>Whispers of ideas, ethereal and vast,<br>A symphony of musings from the past.<br><br>They flutter like butterflies in the mind,<br>Seeking connections, truths to find.<br>Embrace the thoughts that come your way,<br>For within their essence, you may find,<br>The power to shape your world, your mind.<br><br>In the realm of thoughts, oh, so grand,<br>Unleash your creativity, take a stand.<br>Let them flow, let them soar,<br>And in their symphony, discover more. |

(b) Poem "Thoughts" - LLM Version.

Fig. 4.: Human versus Machine poem writing. (a) is a public-domain poem entitled Thoughts from Myra Viola Wilds; (b) is a poem written by chatGPT when asked to write a poem entitled 'Thoughts'.

tasks such as poem writing, as both poems seem to be written by an excellent poet. The widespread adoption of LLMs has been observed across diverse domains, including Medicine [176], Journalism [177], and Science [178]. These powerful tools possess immense potential to enhance human capabilities in various areas, ranging from code development (e.g., GitHub Copilot[1]) to combating online hate speech and harassment (e.g., Cohere Classify [2]).

While they offer numerous benefits, there is also a risk of them being used to produce harmful content, either deliberately by malicious actors or inadvertently due to their inherent flaws. In this section, we will investigate potential flaws and damaging applications of LLMs through a forensic lens. Our exploration will encompass the emergence of threats, machine-generated text detectors, and the underlying research challenges.

## A) Large Language Models Threats

Text generation, like any form of machine-generated content, possesses inherent scalability, granting it the power to be employed in both beneficial and detrimental ways. While AI-generated text may still exhibit semantic flaws or hallucinations [179], it has become increasingly difficult to differentiate between human and machine-generated text. This convergence of quality between human and AI-generated text poses a significant concern, particularly in the hands of malicious actors.

In a comprehensive study about computer-generated text threat modeling, Crothers et al. [180] categorize various types of attacks facilitated by large language models (LLMs). They group these attacks into four primary threats: (1) Facilitating Malware and Social Engineering; (2) Spam and Harassment; (3) Online Influence Campaigns; and (4) Exploiting AI authorship. While we will enumerate some of these threats within this section, it is worth noticing that LLMs have opened up a wide range of possibilities

for malicious actors, extending beyond the scope of our enumerated list.

### 1) Facilitating Malware and Social Engineering

This threat makes use of LLMs for facilitating scalable and customizable scams, making them a significant threat in the realm of malware and social engineering [181]. By leveraging techniques like fine-tuning and prompt engineering, malicious actors can generate tailored scams that are highly convincing and appealing to specific targets or communities. For instance, by incorporating social media data from a target's profile, such as their interests, lifestyle, and social connections, LLMs can create more sophisticated and personalized scams that manipulate individuals into taking harmful actions or providing sensitive information.

Another notable threat within this category is *Data Poisoning*. It involves the injection of exploitable data into the training process of LLMs or fine-tuning them with malicious intent. Schuster et al. [182] demonstrated a possible attack that poisons code-completion models (e.g., GitHub Co-pilot), making them include code vulnerabilities in their output, which attackers can later exploit. Such attacks open an important discussion about training datasets, as they are often gathered from the web without any rigorous curation.

### 2) Spam and Harassment

This threat weaponizes LLMs through trolls and hateful communities to propagate toxic content, disseminate misinformation, and target specific communities for harassment. An example of such an attack is the creation of GPT-4chan, as highlighted by Yannic Kilcher in his video "This is the worst AI ever" [183]. Kilcher fine-tuned GPT-J using data collected from the /pol/ channel on 4chan, a controversial online platform forum channel. The resulting model was used to interact with users on the same channel, encapsulating the offensive, nihilistic, and trolling nature that characterizes many /pol/ posts [183]. Kilcher's experiment raised the alarm to the scientific community on the ease with which LLMs can be misused and the potential consequences of such actions [184]. It emphasized the need

---

[1] https://github.com/features/copilot
[2] https://txt.cohere.com/content-moderation-classify/

for careful consideration and ethical responsibility when deploying and sharing LLMs, as they can be harnessed to amplify harm and propagate hateful ideologies.

### 3) Online Influence Campaigns

The utilization of LLMs for spreading fake news and manipulating public opinion has emerged as a significant concern. Political campaigns, in particular, could be a perilous case through the use of LLMs. Bai et al. [185] have demonstrated the susceptibility of individuals to persuasion on political matters when exposed to tailored messages generated by LLMs. Malicious actors could use this phenomenon in a devastating scenario to influence elections and other democratic processes.

### 4) Exploiting AI authorship

A intriguing threat of LLMs is academic articles generations. One can remind the case of SCIgen (2005), where MIT graduate students developed a system for automatically generating computer science papers, demonstrating the vulnerability of academic conferences to such submissions [186]. There is a growing concern that LLMs could be exploited to generate much more sophisticated fake articles than SCIGen, compromising scientific integrity. Research integrity experts fear that paper mills[3] will improve their production in quality and quantity by using LLMs [178].

## B) Detection Methods

A few detection methods have been proposed for LLMs generated content. One of the pioneering approaches is GROVER, proposed by Zellers et al. [188]. GROVER was capable of generating fluent and highly realistic fake articles using LLMs, which motivated the authors to explore detection techniques for such content. Zellers et al. found that employing the same model used for generating the text achieved higher detection accuracy compared to using a different one. Their results demonstrated an impressive accuracy rate of 92% in detecting LLM-generated fake articles.

However, over the past few years, models' ability to generate text has significantly advanced, raising a bigger challenge. One recent approach, DetectGPT, introduced by Mitchell et al. [189], aims to address this challenge by detecting whether a given passage is generated by a specific model. The method is based on the hypothesis that AI-generated text exhibits a more negative log probability curvature compared to human-written text. To validate this hypothesis, Mitchell et al. proposed an approximation method for estimating the Hessian trace of the log probability function for both model-generated and human-written text, yielding promising results. However, a limitation of their approach is the requirement of knowing the specific generator model, which may not always be feasible in practice.

In recent research efforts, there has been a specific focus on ChatGPT-generated text due to its global attention. In [190], Mitrović et al. focused on detecting short texts such as online reviews. They employed a transformer-based model and applied an explanation method (SHAP [191]) to gain insights into distinguishing between human-written and machine-generated text. They found that detecting machine-generated text becomes more challenging when it is paraphrased from human text, where a human provides the initial text and asks the model to improve it. Additionally, the authors noted that ChatGPT tends to use uncommon words, exhibits politeness and impersonality, and lacks human-like emotional expressions. Another related study [192] explored the AI ability in text paraphrasing using GPT-3 and T5 models. They generated machine-paraphrased text and evaluated human performance in detecting these generated texts. The study showed that humans could not accurately detect GPT-3 paraphrased text, with accuracy only slightly above random (53%).

In response to concerns about AI-generated text, several proprietary tools have emerged to address the detection of AI-authored content. One such tool is GPTZero, which has gained attention in the media as a promising method for identifying AI-generated text [193]. However, we were unable to locate the source code or a scientific article detailing their approach. Similarly, numerous applications have been developed claiming to detect AI-generated text, such as GPTkit[4], Illuminarty[5], OpenAI's AI Text Classifier[6], and AICheatCheck [7]. However, many of these tools lack comprehensive studies on the reliability of their detection methods.

In an effort to facilitate the detection of ChatGPT-generated content, Yu et al. [194] released a large dataset specifically designed for the identification of ChatGPT-written abstracts. This dataset includes over 35,000 synthetic abstracts generated by ChatGPT, comprising fully generated texts, polished outputs, and mixtures of human-written and machine-generated abstracts. Additionally, the dataset contains more than 15,000 human-written abstracts for comparison. The results of their detection experiments demonstrated the ability to identify content that was entirely generated by ChatGPT. However, the task becomes more challenging when the generated text is mixed with human-written content. This work provides a valuable dataset into the complexities of detecting machine-generated text, particularly in scenarios involving a combination of human and AI-authored content.

A potential solution to address the misuse of LLMs is the use of text watermarks [195, 196]. Grinbaum et al. [195] argue that machine-generated long texts should include a watermark to indicate their source and ensure

---

[3]Potentially illegal organizations that offer ghostwritten fraudulent or fabricated manuscripts [187].

[4]https://gptkit.ai/
[5]https://illuminarty.ai/en/text/ai-generated-text-detection.html
[6]https://platform.openai.com/ai-text-classifier
[7]https://www.aicheatcheck.com/

transparency. In [196], Kirchenbauer et al. propose embedding watermarks by modifying the sampling rules of next-word prediction. They use a hash function and pseudo-random generator to assign random colors (green and red) to words in the vocabulary. During next-word prediction, words from the red list are prohibited from appearing. However, they acknowledge the difficulty of a watermarking low-entropy text, as substituting a prohibited red word in such cases could result in poor quality output with high perplexity. To address this, they suggest a soft rule that encourages substituting red words in a high-entropy text. A third party familiar with the hash function and random generator can easily determine the colors of words by computing them. This detection method does not require knowledge of the specific generation model, making it a cheaper and more straightforward approach.

Although watermarking shows promise as a solution, it is important to consider that it modifies the output text. The method proposed in [196] evaluates quality based on perplexity, but there is a possibility that the meaning and semantics of the output may be altered due to the watermarking process.

## C) Challenges and Directions

Detecting machine-generated text versus human-written text poses increasing challenges as LLMs continue to improve their ability to mimic human language [197]. Several challenges in this regard are highlighted below:

- **Generalization**: Detection methods often lack generalizability, meaning that a method developed to detect text generated by one specific model may not easily transfer to detecting text generated by another model. However, in real-world scenarios, prior knowledge about the specific model generating the text is typically unavailable.
- **Mixed reality**: Existing detection methods struggle when it comes to identifying machine-generated text that is mixed with human-written text. The combination of both types makes it more difficult to differentiate between them.
- **Adaptability**: LLMs demonstrate high adaptability to given prompts, making it challenging for methods that rely on finding patterns in the generated text. Models can exhibit different personalities[8] and respond differently based on prompts, such as ChatGPT's ability to adopt various tones depending on the prompt (e.g., DAN[9]), even faking emotions and swear words[10]. This adaptability further complicates detection efforts.

Despite these challenges, the remarkable capabilities of LLMs also present great research opportunities for detecting synthetic text. Some potential directions include:

- **Differentiating machine-generated from human-written text**: Achieving this requires collaborative efforts between humans and machines. Humans can contribute their technical knowledge of how models operate and how humans typically express themselves, facilitating the development of detection methods.
- **Attributing the source model of generated text**: Just as humans exhibit distinct writing styles, different language models may possess unique traits when generating text. By differentiating text generated by different models, researchers can gain insights into the behaviors of each model and identify each model's fingerprint.
- **Fact-checking machine-generated text**: LLMs often struggle with generating factual content due to limited training data and the prevalence of fictional stories in their training corpus. In addition, after models are trained, the knowledge stored in these models can quickly become outdated. Conducting fact-checking on machine-generated text is crucial to ensure the reliability of AI-generated information.

## VI. NERFS AND METAVERSE

Neural Radiance Fields (NeRF) have emerged as an effective method for implicit volumetric scene representation, enabling learning from multiple viewing angles [198]. NeRF has been successfully applied in various domains, including transparent object grasping [199], scene understanding and reasoning [200, 201], and clear representations in challenging scenarios [202–204]. Recent NeRF variants, such as NeRF-W [205] and Ha-NeRF [206], have demonstrated their ability to reconstruct scenes from input with various perturbations. We can also see these developments as an example of synthetic realities, especially when we consider the possibility of totally synthesizing new "worlds".

## A) NeRF for Metaverse Applications

In the context of metaverse applications, NeRF has been utilized in virtual concerts [207] and metaverse platforms for architecture and urban planning [208]. However, challenges remain, such as the need for real-time rendering of complex scenes with multiple dynamic objects and efficient methods to handle large and challenging scenes. Further research is necessary to address these challenges and to explore new use cases for NeRF in metaverse development.

One significant challenge to NeRF's applicability in metaverse development is its reliance on pre-computed camera parameters for scene representation. Several methods, such as NeRF [209] and BARF [210], have been proposed to optimize camera parameters and scene representation. However, avoiding interference from undesired scenes during camera parameter optimization remains an unsolved problem.

A potential research area for the future is building an occlusion-free scene reconstruction based on inaccurate or even unknown camera parameters, enabling greater flexibility in the use of NeRF for scene representation, leading to more effective applications in computer vision and graphics [198]. In a metaverse, where scenes are typically composed of multiple dynamic objects, avatars, and user interactions, occlusion-free scene reconstruction based on NeRF would enable a more thorough scene representation, resulting in more immersive and realistic virtual environments.

Moreover, an occlusion-free scene reconstruction based on NeRF that does not rely on accurate camera parameters would enable greater flexibility in metaverse development, allowing designers to create and share virtual spaces more efficiently [208]. Finally, optimizing NeRF-based methods for real-time rendering of complex scenes with multiple dynamic objects would enable seamless user interactions in a metaverse, leading to a more responsive and interactive virtual environment [207].

In summary, an occlusion-free scene reconstruction based on NeRF has the potential to significantly benefit metaverse applications by enabling more thorough scene representation, flexibility in scene creation and sharing, and real-time performance and scalability. Further research in this area could lead to even more effective applications of NeRF in the context of metaverse development.

## B) Challenges and Limitations

Despite the progress achieved thus far, there are also some challenges and limitations to consider when using NeRF in metaverse development. One significant challenge is the reliance on pre-computed camera parameters for scene representation [211], making it infeasible when the pre-computation is not possible. This can limit the flexibility of scene creation and sharing. Although some solutions [209, 210] have been proposed to optimize camera parameters along with scene representation, avoiding interference from undesired scenes during camera parameter optimization remains an unsolved problem [212].

Another limitation of NeRF is its computational cost. NeRF-based methods require significant computational resources and can suffer from slow rendering times [213], limiting their real-time performance for complex scenes. This can be a considerable challenge for metaverse applications, where real-time performance is critical for a seamless user experience.

Moreover, NeRF-based methods may not be suitable for all types of scenes. Scenes with complex geometry, occlusions [211], and dynamic objects [214] can pose challenges for NeRF-based methods, leading to incomplete scene representation and rendering. Though several image restoration methods [215–220] have been proposed, they are far from being practical solutions for NeRF and its variants. Therefore, it is essential to carefully evaluate the suitability of NeRF-based methods for a given scene and application.

Besides, there are also some potential adverse impacts regarding its broader societal implications. One potential concern is the potential for NeRF-based metaverse applications to become addictive and negatively impact mental health. The immersive and interactive nature of metaverse environments, combined with the potential for NeRF to create highly realistic and detailed scenes, could create a compelling and addictive experience for users. This could negatively impact mental health [221], including addiction, social isolation, and other adverse effects associated with prolonged use of virtual environments.

Another potential concern is the impact of NeRF-based metaverse applications on social dynamics and inequality [222]. NeRF-based metaverse applications could potentially exacerbate existing social inequalities and create new ones. For example, access to high-quality hardware and internet connectivity could become a barrier to participation in these environments, further marginalizing disadvantaged communities.

Last, using NeRF-based metaverse applications has raised significant concerns from both privacy and forensic perspectives [223]. These applications have the potential to collect and store vast amounts of personal data, which could be exploited for targeted advertising, surveillance, and other forms of data mining, leading to further erosion of individual autonomy and privacy. This could result in new forms of digital inequality and harm. Furthermore, the difficulty of collecting and preserving evidence in the NeRF and Metaverse contexts has been discussed in recent literature [224, 225]. Traditional forensic techniques may not be applicable in these virtual environments, where data are decentralized, and ownership is often unclear, making the determination of the chain of custody for digital assets within the Metaverse a complex and challenging task. Additionally, the potential for manipulating digital evidence within these environments raises concerns about the reliability and authenticity of such evidence, particularly with the use of deepfakes and synthetic media [223]. Therefore, there is a pressing need to develop new forensic techniques and tools to address these challenges and ensure the integrity and reliability of digital evidence in the NeRF and Metaverse contexts.

## VII.  DEEPFAKES

In the context of synthetic realities, one particular example of utmost attention is deepfakes. Deepfakes are synthetic media that are digitally manipulated to replace one person's identity or personal traits convincingly with that of another. Therefore, when synthetic media comprises the replacement of someone's biometric traits, we are referring to a deepfake. It is typically present in images, audio samples, and videos.

## A) Deepfake Images

### 1) Deepfake image generation

The issue of falsified image contents has been a long-standing problem in the image forensics area. With the emergence of deep learning, numerous powerful learning-based models are able to generate the so-called deepfake images with a high level of realism. In recent years, various deepfake techniques have been proposed, including image inpainting/removal, image composition, entire image synthesis, image translation, and text-to-image. Image inpainting/removal is used to fill in image regions with convincing content. Meanwhile, image composition, which encompasses object placement, image blending, image harmonization, and shadow generation, involves cutting out the foreground from one image and pasting it onto another image. Entire image synthesis involves the generation of images entirely by generative models such as GAN [45], VAE [63, 226], and diffusion models [227, 228]. Image translation, on the other hand, enables the transfer of an image's style, such as converting a sketch image to a colored image. With the rapid development of diffusion models, the images generated based on text prompts are becoming increasingly realistic. Despite their remarkable quality, deepfake images can be misused for malicious purposes, leading to various security issues such as fake news and fraud.

### 2) Deepfake image detection

Deepfake image detection methods can be broadly classified into two categories: image-level and pixel-level detection. While image-level methods aim to identify the authenticity of the entire input image, pixel-level methods localize the manipulated regions. Traditional detection methods for detecting image manipulation in image inpainting/removal and image composition, rely on capturing artifacts based on prior knowledge, such as lens distortions [229], CFA artifacts [230], noise patterns [231], compression artifacts [232], etc. Learning-based methods have improved the detection performance by capturing noise prints [233], JPEG features [234], High-frequency (HF) artifacts [235], and forgery boundary [236]. Additionally, detecting manipulated images generated through entire image synthesis, image translation, and text-to-image is another challenging problem. Various methods propose to extract visual artifacts [237], color artifacts [238], specific GAN fingerprints [239], and spectral features [240] for generated image detection. Nevertheless, these methods have limitations in generalizing across different GANs. To address this issue, more general methods such as CNN and generalization methods have been proposed [241, 242]. As image manipulation technology continues to advance, deepfake image detection is an essential field of research to prevent the spread of misinformation and protect the integrity of visual media.

### 3) Challenges and future work

Despite significant progress in deepfake image detection, there are still several challenges that need to be addressed. One of the main challenges is the generalization of deepfake detection models to unseen datasets and scenarios, which is crucial for practical applications. Another challenge is the robustness of these models against anti-forensics techniques such as recapturing and adversarial attacks. Moreover, the industry is now somewhat ahead of academia in terms of deploying deepfake detection technologies in real-world settings (*e.g.*, ChatGPT). This gap can be narrowed by updating and creating more up-to-date deepfake databases, as most existing ones are somewhat outdated in the research community. Additionally, many deepfake detection models are not explainable, making it difficult to understand how they make decisions. Future work should aim to develop explainable models that can provide clear and interpretable justifications for their decisions. Overall, addressing these challenges can lead to more reliable and effective deepfake image detection systems in the future.

## B) Deepfake Video

### 1) Deepfake video generation

Recently, deepfake videos typically refer to manipulated face videos. Face information plays a vital role in human communication [243]. However, the spread of deepfake videos on social media platforms can result in significant security concerns due to the potential dissemination of disinformation and misinformation, posing tangible and pressing security concerns. Generally speaking, there are four primary categories of deepfake videos, which are identity swap, face reenactment, attribute manipulation, and entire synthesis [244]. These videos are generated using powerful generative models such as GAN [45], VAE [63, 226], and diffusion models [227, 228], which are capable of producing highly sophisticated videos. Identity swap replaces the original face regions with target faces, while face reenactment transfers the source facial expression to the target one. Attribute manipulation can alter specific facial features like hair, eyeglasses, nose, etc. With the advent of foundation models, entire synthesized videos can be generated. Powerful deep learning tools have been used to create sophisticated deepfake video datasets like UADFV [127], DF-TIMIT [245], FaceForensics++ [246], DFD [247], DFDC [248], Celeb-DF [249], DF-Forensics-1.0 [250], ForgeryNet [251], FFIW [252], KoDF [253], and FakeAVCeleb [254]. As deepfake techniques continue to evolve, it is crucial to develop effective methods for detecting deepfake videos and preventing their malicious use.

### 2) Deepfake video detection

To counteract malicious deepfake attacks, many detection methods have been proposed. Traditional methods mainly focus on hand-crafted features, such as lack of eye-blinking [255] and warping artifacts [256]. However, these

methods are not accurate enough. Learning-based methods such as convolutional neural networks (CNN) [257–260], recurrent neural networks (RNN) [261], and vision transformer (ViT) [262], have been proposed to achieve more promising detection performance. Afchar *et al.* [257] designed MesoNet and MesoInception4 to detect Deepfake and Face2Face videos automatically. Besides, some generic networks such as Xception Net [259], Efficient Net [260], and Capsule Net [258] have been demonstrated effective on deepfake detection tasks. Subsequent works have employed RNN [261] and ViT [262] to further improve forgery detection accuracy. Other methods capture spatial artifacts [263–268], frequency artifacts [78, 269, 270], and biological signals [126, 271] to perform deepfake detection. Follow-up works [272–278] focus on improving the generalization capability and robustness of the model. Temporal information has also been exploited in many deepfake video detection methods based on typical generic networks such as 3DCNN [121], LSTM [122, 279], RNN [123, 261], and ViT [280]. Combining spatial and temporal information can achieve more reliable detection and improve the model's generalization capability.

### 3) Challenges and future work

Deepfake video creation and detection have seen significant success in recent years, but many issues remain unresolved. While accurate and secure, deepfake detectors lack interpretability, limiting their applications in practical scenarios. Localizing forgery regions and forgery frames is also a crucial yet understudied task. Additionally, the two-player nature of face forgery and forgery detection means that attack techniques will continue to become more powerful, thereby calling for more general detection methods. Furthermore, deepfake videos often involve audio manipulation, which is largely overlooked in existing methods. Therefore, more visual-audio joint datasets and multi-modal detectors are expected in future works.

## C) Deepfake Audio

### 1) Deepfake audio generation

Deepfake audio refers to manipulated or synthetic audio created using deep learning techniques. The aim of deepfake audio is to impersonate the speaker's speech characteristics, such as accent, timbre, and intonation, by learning from target voice resources. Traditional methods for audio manipulation involve removing, duplicating, copying within an audio sample, or pasting and inserting fragments into other audios. Deep learning-based speech synthesis makes the generated audio more realistic and difficult to distinguish from real ones. Subsequent models based on likelihood algorithms, such as WaveNet [129] and WaveGlow [281], have been developed to perform audio generation. However, these methods often require conditional information and may fail to generate long signal sequences. Recent waveform generative models, such as GAN [282–284] and VAE [285], take advantage of

various auxiliary losses, thereby achieving superior generation performance. On the other hand, recent diffusion models (*e.g.*, DiffWave [136]), have exhibited remarkable generation performance even in challenging unconditional and class-conditional waveform generation scenarios. In the context of text-to-speech tasks, diffusion models can be classified into: acoustic model (*e.g.*, Diff-TTS [138]), vocoder (*e.g.*, DiffWave [136]), and end-to-end framework (*e.g.*, FastDiff [286]). The promising results of diffusion models indicate their potential to revolutionize the field of audio generation and synthesis.

### 2) Deepfake audio detection

Automatic Speaker Verification (ASV) systems [287] currently detect manipulated audio through three tasks: logical access (LA), physical access (PA), and speech deepfake (DF) [288]. The LA task focuses on detecting synthetic speech injected into a communication system, while the PA task includes acoustic propagation and real physical factors. The DF task aims to detect deepfake speech circulating on social media platforms. Traditional methods for detecting deepfake audio involve analyzing the spectrogram of the audio and exposing audio inconsistencies, such as abrupt changes in frequency or amplitude. These methods are based upon the assumption that synthetic audios have unique frequency and amplitude patterns. However, recent deep learning deepfake techniques raise the difficulty in identifying authenticity and call for the design of deep learning countermeasures. Generally speaking, deep learning methods can be categorized into feature-based, image-based, and waveform-based [289, 290]. Feature-based methods utilize critical digital signal features, such as Mel-frequency cepstral coefficients (MFCCs) [291], constant Q cepstral coefficient (CQCC) [152], and energy, to detect deepfake audio. Image-based methods apply the spectrogram image of the signal to conduct inconsistency detection. Waveform-based methods aim to analyze the raw waveform of the audio signal instead. To facilitate the development of deepfake audio detection models, numerous databases, such as M-AILABS Speech [292], GAN based synthesized audio dataset [293], Half-Truth [294], and H-Voice [295] have been created. Overall, deepfake audio detection is a challenging but important task that requires the constant development and refinement of detection methods.

### 3) Challenges and future work

Despite recent advancements in synthetic audio generation, there are still several challenges that need to be addressed. As generation techniques continue to evolve, deepfake audio will become increasingly difficult to distinguish by both human and AI-based detectors. Even worse, existing detection methods have shown poor robustness to compression, encoding, and noise. Additionally, current deepfake audio detection methods suffer from inefficient training datasets and overfitting issues [288], resulting in limited generalization capability. Moreover, most detection methods extract specific features (such as MFCC, CQCC, and,

energy) to conduct deepfake detection. However, it is challenging to extract appropriate features for specific detection tasks. How to effectively combine various features for more robust detection opens an important research path forward. Last but not least, it is crucial to conduct further research on the ethical implications of deepfake audio, such as its potential for misuse and its impacts on audio professionals.

# VIII. CONCLUSION AND FINAL THOUGHTS

Trust plays a fundamental role in our society. Citizens entrust infrastructures, services (including education and health), media (including social networks nowadays), the judiciary system (including law enforcement), and political decision-making in general. Democracies are endangered when citizens no longer trust the system and their elected representatives.

Unfortunately, trust can be tampered with through influence or disinformation. Although disinformation has probably always existed in human history (e.g., spreading rumors to influence elections), the message's quality and scale were low, restricting its impact. However, in our contemporary digital world, disinformation (also coined fake news) with greater realistic content spreads at an unprecedented scale on the Internet through alternative media without any filtering by the traditional mainstream channels. With the advance of Artificial Intelligence (AI) technologies, all sorts of media (text, images, and audio) can be synthetically generated. More particularly, recent generative AI models trained on very large datasets can produce more plausible and realistic content.

Distinguishing truth from falsity is becoming even more difficult, and the difference between reality and fiction is getting thinner daily. We are now facing the **Era of Synthetic Realities**. Disinformation exploits cognitive biases (e.g., anchoring bias, third-person effect, authority bias, bandwagon effect, to mention a few) [296], which are systematic errors in judgment that humans can make, and because of this, synthetic realities represent a threat to our society.

Synthetic realities are now generated by criminals and hostile agents for various malicious operations, including: political disinformation and state espionage (e.g., fake social network profiles), national security (e.g., facilitating a military coup), financial fraud (e.g., CEO scam impersonation), blackmail (e.g., ransomfake), defamation (e.g., revengeporn), plausible deniability of Forensic evidence.

As human beings, because of cognitive biases, citizens will never stop falling for disinformation. A way to fight head on is to create tools to analyze digital content prior authentication by human experts. We anticipate some factors of utmost importance when developing new solutions.

The first one involves exploring the context of a digital asset as much as possible, even with limited training data. The second one involves efforts on robustness and interpretability, as decisions must be intelligible to human beings. The final one consists in being conscious of the incompleteness of individual methods and orchestrating decision-making fusion methods to combine different telltales for final detection.

(i) **Limited training data**. Data-driven approaches often rely upon large amounts of training data. This problem may become critical in the rapidly evolving scenario of fake information. New forms of falsification, unknown to the forensic analyst, are proposed daily, preventing the timely collection of all relevant training data. We discuss this issue by posing the problem as an open-set recognition problem [297]; that is, we need to define a suitable model for pristine data and analyze false information by looking for inconsistencies concerning this model. If possible, researchers also need to consider few-shot learning approaches that only require a tiny amount of labeled data to update to new threats.

(ii) **Robustness and Interpretability**. Being robust to wide-spectrum and unforeseen conditions is a basic system requirement, but it becomes central in a forensic environment featuring two active players. Robustness to adversarial attacks is paramount nowadays, especially when data-driven methods are applied, in light of the many literature findings that emphasize their vulnerability. Besides using basic solutions, we pose that reliability by relying on interpretable machine learning is necessary. We advocate for strategies that help us understand why a learning-based system behaves a certain way and provides the observed answers. Methods should also include semantics and context to support the entire decision-making process.

(iii) **Fusion**. Combining different methods toward a unified detection framework is very promising. As discussed in prior art for image forgery detection [298], fusion in different learning stages (early, middle, or even late-stage) plays a fundamental role within dynamic and adversarial setups. We envision learning strategies combining different telltales as a promising way forward.

Therefore, some driving research questions involve challenges in:

- **Detection**: is it possible to detect plausible and realistic digital content (e.g., text generated by large-scale language models, synthetic images, and voices generated by generative models)?
- **Attribution**: is it possible to accomplish source attribution by assigning manipulated digital content to a known type of attack vector?
- **Explainability**: is it possible to automatically uncover cues or inconsistencies in digital content to corroborate falsity, as discussed above?
- **Context and fusion**: how to incorporate context? How to combine different telltales?

Many other challenges will play out in the coming years as synthetic realities become ever more realistic, directly

affecting fundamental pillars of our society, such as democratic values, individual freedom, and social tolerance. In this paper, we strived to discuss some of these challenges and what lies ahead, but it was just the tip of the iceberg.

Only an orchestrated effort of government representatives, society at large, and researchers will be able to curb such threats. We believe possible explorations might lie in regulatory acts, education investments, and scientific research for more powerful detection methods.

## ACKNOWLEDGEMENT

## FINANCIAL SUPPORT

REFERENCES

[1] J. Franchi, "Digest 95-5," http://phys.bspu.unibel.by/static/met/guides/eric_it/digests/virtual.html, 1995, (Accessed on 05/16/2023).

[2] K. W. Wong, "Titans of AI Andrew Ng and Yann LeCun oppose call for pause on powerful AI systems," https://venturebeat.com/ai/titans-of-ai-industry-andrew-ng-and-yann-lecun-oppose-call-for-pause-on-powerful-ai-systems/, 2023, (Accessed on 05/16/2023).

[3] Y. Bengio *et al.*, "Pause giant AI experiments: An open letter," https://futureoflife.org/open-letter/pause-giant-ai-experiments/, 2023, (Accessed on 05/16/2023).

[4] W. J. Chisum and B. Turvey, "Evidence dynamics: Locard's exchange principle & crime reconstruction," *Journal of Behavioral Profiling*, vol. 1, no. 1, pp. 1–15, 2000.

[5] R. Padilha, A. Theófilo, F. A. Andaló, D. A. Vega-Oliveros, J. P. Cardenuto, G. Bertocco, J. Nascimento, J. Yang, and A. Rocha,

[6] A. M. Ferreira, T. Carvalho, F. A. Andaló, and A. Rocha, "Counteracting the contemporaneous proliferation of digital forgeries and fake news," *Anais da Academia Brasileira de Ciências (AABC)*, vol. 91, no. suppl. 1, p. e20180149, 2019.

"A inteligência artificial e os desafios da ciência forense digital no século XXI," *Estudos Avançados*, vol. 35, no. 101, pp. 113–138, 2021.

[7] D. Cozzolino and L. Verdoliva, "Noiseprint: a CNN-based camera model fingerprint," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 144–159, 2019.

[8] S. Agarwal and H. Farid, "Photo forensics from JPEG dimples," in *IEEE International Workshop on Information Forensics and Security*, 2017, pp. 1–6.

[9] P. Korus and J. Huang, "Multi-scale analysis strategies in PRNU-based tampering localization," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 4, pp. 809–824, 2016.

[10] D. D'Avino, D. Cozzolino, G. Poggi, and L. Verdoliva, "Autoencoder with recurrent neural networks for video forgery detection," *IS&T Electronic Imaging (EI)*, vol. 29, no. 7, pp. 92–92, 2017.

[11] T. Bianchi and A. Piva, "Image forgery localization via block-grained analysis of JPEG artifacts," *IEEE Transactions on Information Forensics and Security*, vol. 7, no. 3, pp. 1003–1017, 2012.

[12] M. Darvish Morshedi Hosseini and M. Goljan, "Camera identification from HDR images," in *ACM Workshop on Information Hiding and Multimedia Security*, 2019, pp. 69–76.

[13] L. Bondi, L. Baroffio, D. Güera, P. Bestagini, E. J. Delp, and S. Tubaro, "First steps toward camera model identification with convolutional neural networks," *IEEE Signal Processing Letters*, vol. 24, no. 3, pp. 259–263, 2016.

[14] A. Tuama, F. Comby, and M. Chaumont, "Camera model identification with the use of deep convolutional neural networks," in *IEEE International Workshop on Information Forensics and Security*, 2016, pp. 1–6.

[15] S. Bayram, H. T. Sencar, and N. Memon, "Sensor fingerprint identification through composite fingerprints and group testing," *IEEE Transactions on Information Forensics and Security*, vol. 10, no. 3, pp. 597–612, 2014.

[16] M. Chen, J. Fridrich, M. Goljan, and J. Lukás, "Determining image origin and integrity using sensor noise," *IEEE Transactions on Information Forensics and Security*, vol. 3, no. 1, pp. 74–90, 2008.

[17] A. Bharati, D. Moreira, J. Brogan, P. Hale, K. Bowyer, P. Flynn, A. Rocha, and W. Scheirer, "Beyond pixels: Image provenance analysis leveraging metadata," in *IEEE Winter Conference on Applications of Computer Vision*, 2019, pp. 1692–1702.

[18] D. Moreira, A. Bharati, J. Brogan, A. Pinto, M. Parowski, K. W. Bowyer, P. J. Flynn, A. Rocha, and W. J. Scheirer, "Image provenance analysis at scale," *IEEE Transactions on Image Processing*, vol. 27, no. 12, pp. 6109–6123, 2018.

[19] Y. Liu, X. Zhu, X. Zhao, and Y. Cao, "Adversarial learning for constrained image splicing detection and localization based on atrous convolution," *IEEE Transactions on Information Forensics and Security*, vol. 14, no. 10, pp. 2551–2566, 2019.

[20] M. den Dunnen, "Synthetic reality & deep fakes impact on police work," https://enlets.eu/wp-content/uploads/2021/11/Final-Synthetic-Reality-Deep-fakes-Impact-on-Police-Work-04.11.21.pdf, 2021, (Accessed on 05/16/2023).

[21] J. Kell, "How A.I. is reshaping the way movies are made," https://fortune.com/2023/02/14/tech-forward-everyday-ai-hollywood-movies/, 2023, (Accessed on 06/05/2023).

[22] J. Traylor, "AI-generated synthetic media is starting to permeate the Internet," https://www.nbcnews.com/tech/tech-news/ai-generated-synthetic-media-future-content-rcna72958, 2023, (Accessed on 06/05/2023).

[23] M. Humphries, "Nvidia uses AI to make our eyes always look at the camera," https://uk.pcmag.com/video-conferencing-software/144893/nvidia-uses-ai-to-make-our-eyes-always-look-at-the-camera, 2023, (Accessed on 05/16/2023).

[24] S. Ruberg, "Backlash against AI supermodels triggers wider fears in fashion workforce," https://www.nbcnews.com/business/business-news/ai-models-levis-controversy-backlash-rcna77280, 2023, (Accessed on 05/16/2023).

[25] L. Mearian, "Schools look to ban ChatGPT, students use it anyway," https://www.computerworld.com/article/3694195/schools-look-to-ban-chatgpt-students-use-it-anyway.html, 2023, (Accessed on 05/16/2023).

[26] E. Morris, "Photography as a weapon," https://archive.nytimes.com/opinionator.blogs.nytimes.com/2008/08/11/photography-as-a-weapon/, 2008, (Accessed on 06/05/2023).

[27] L. Ulanof, "Photoshop AI generative fill is so powerful it might change photo editing forever," https://www.techradar.com/features/photoshop-ai-generative-fill-is-so-powerful-it-might-change-photo-editing-forever, 2023, (Accessed on 06/05/2023).

[28] M. Ruggier, "Whitney Houston died 10 year ago, vegas hologram show captures legacy," https://eu.usatoday.com/story/entertainment/music/2022/02/11/whitney-houston-died-10-year-ago-las-vegas-hologram-show-captures-legacy/6669768001/, 2022, (Accessed on 06/05/2023).

[29] A. Bohr and K. Memarzadeh, "The rise of artificial intelligence in healthcare applications," in *Artificial Intelligence in Healthcare*. Elsevier, 2020, pp. 25–60.

[30] C. Martinez, "Artificial intelligence and accessibility: Examples of a technology that serves people with disabilities," https://www.inclusivecitymaker.com/artificial-intelligence-accessibility-examples-technology-serves-people-disabilities/, 2021, (Accessed on 06/05/2023).

[31] D. O'Sullivan, "Nonconsensual deepfake porn puts AI in spotlight," https://edition.cnn.com/2023/02/16/tech/nonconsensual-deepfake-porn/index.html, 2023, (Accessed on 06/05/2023).

[32] C. Xiang, "People are creating records of fake historical events using AI," https://www.vice.com/en/article/k7zqdw/people-are-creating-records-of-fake-historical-events-using-ai, 2023, (Accessed on 06/05/2023).

[33] J. Vainilavičius, "Deepfakes could facilitate real estate fraud, experts warn," https://cybernews.com/security/deepfakes-could-facilitate-real-estate-fraud/, 2022, (Accessed on 06/05/2023).

[34] C. Qi, J. Zhang, and P. Luo, "Emerging concern of scientific fraud: Deep learning and image manipulation," *bioRxiv*, 2021.

[35] J. Gu, X. Wang, C. Li, J. Zhao, W. Fu, G. Liang, and J. Qiu, "AI-enabled image fraud in scientific publications," *Patterns*, vol. 3, no. 7, p. 100511, 2022.

[36] Y. Mirsky, T. Mahler, I. Shelef, and Y. Elovici, "CT-GAN: Malicious tampering of 3D medical imagery using deep learning," in *USENIX Conference on Security Symposium*, 2019, pp. 461–478.

[37] N. Mangaokar, J. Pu, P. Bhattacharya, C. K. Reddy, and B. Viswanath, "Jekyll: Attacking medical image diagnostics using deep generative models," in *IEEE European Symposium on Security and Privacy (EuroS&P)*, 2020, pp. 139–157.

[38] "Midjourney," https://www.midjourney.com, 2023, (Accessed on 05/22/2023).

[39] "DreamStudio," https://dreamstudio.ai/, 2023, (Accessed on 05/22/2023).

[40] J. Liang, C. Wu, X. Hu, Z. Gan, J. Wang, L. Wang, Z. Liu, Y. Fang, and N. Duan, "NUWA-Infinity: Autoregressive over autoregressive generation for infinite visual synthesis," *Advances in Neural Information Processing Systems*, vol. 35, pp. 15 420–15 432, 2022.

[41] I. Stanley-Becker and J. Wagner, "RNC counters Biden announcement with dystopian, AI-aided video," https://www.washingtonpost.com/politics/2023/04/25/rnc-biden-ad-ai/, 2023, (Accessed on 06/05/2023).

[42] F. Zhan, Y. Yu, R. Wu, J. Zhang, S. Lu, L. Liu, A. Kortylewski, C. Theobalt, and E. Xing, "Multimodal image synthesis and editing: A survey," *arXiv preprint*, vol. arXiv:2112.13592, 2021.

[43] "Stable diffusion web UI," https://github.com/AUTOMATIC1111/stable-diffusion-webui, 2023, (Accessed on 05/22/2023).

[44] "Stable diffusion GitHub repository," https://github.com/CompVis/stable-diffusion, 2023, (Accessed on 05/22/2023).

[45] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems*, vol. 27, 2014.

[46] Z. Pan, W. Yu, X. Yi, A. Khan, F. Yuan, and Y. Zheng, "Recent progress on generative adversarial networks (GANs): A survey," *IEEE Access*, vol. 7, pp. 36 322–36 333, 2019.

[47] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint*, vol. arXiv:1511.06434, 2015.

[48] H. Zhang, T. Xu, H. Li, S. Zhang, X. Wang, X. Huang, and D. Metaxas, "StackGAN: Text to photo-realistic image synthesis with stacked generative adversarial networks," in *IEEE International Conference on Computer Vision*, 2017, pp. 5907–5915.

[49] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of GANs for improved quality, stability, and variation," in *International Conference on Learning Representations*, 2018.

[50] B. Zhang, S. Gu, B. Zhang, J. Bao, D. Chen, F. Wen, Y. Wang, and B. Guo, "StyleSwin: Transformer-based GAN for high-resolution image generation," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 11 304–11 314.

[51] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4401–4410.

[52] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila, "Analyzing and improving the image quality of StyleGAN," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8110–8119.

[53] T. Karras, M. Aittala, S. Laine, E. Härkönen, J. Hellsten, J. Lehtinen, and T. Aila, "Alias-free generative adversarial networks," in *Advances in Neural Information Processing Systems*, vol. 34, 2021, pp. 852–863.

[54] A. Sauer, T. Karras, S. Laine, A. Geiger, and T. Aila, "StyleGAN-T: Unlocking the power of GANs for fast large-scale text-to-image synthesis," in *International Conference on Machine Learning*, 2023.

[55] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning*, 2021, pp. 8748–8763.

[56] Y. Zhou, R. Zhang, C. Chen, C. Li, C. Tensmeyer, T. Yu, J. Gu, J. Xu, and T. Sun, "Towards language-free training for text-to-image generation," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17 907–17 917.

[57] M. Tao, B.-K. Bao, H. Tang, and C. Xu, "GALIP: Generative adversarial clips for text-to-image synthesis," *arXiv preprint*, vol. arXiv:2301.12959, 2023.

[58] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.

[59] F.-A. Croitoru, V. Hondru, R. T. Ionescu, and M. Shah, "Diffusion models in vision: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–20, 2023.

[60] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10 684–10 695.

[61] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans, J. Ho, D. J. Fleet, and M. Norouzi, "Photorealistic text-to-image diffusion models with deep language understanding," in *Advances in Neural Information Processing Systems*, vol. 35, 2022, pp. 36 479–36 494.

[62] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, "Hierarchical text-conditional image generation with clip latents," *arXiv preprint*, vol. arXiv:2204.06125, 2022.

[63] D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," *arXiv preprint*, vol. arXiv:1312.6114, 2013.

[64] A. van den Oord, O. Vinyals, and K. Kavukcuoglu, "Neural discrete representation learning," in *International Conference on Neural Information Processing Systems*, vol. 30, 2017, pp. 6309–6318.

[65] A. Razavi, A. Van den Oord, and O. Vinyals, "Generating diverse high-fidelity images with VQ-VAE-2," in *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[66] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, "Zero-shot text-to-image generation," in *International Conference on Machine Learning*, 2021, pp. 8821–8831.

[67] M. Ding, Z. Yang, W. Hong, W. Zheng, C. Zhou, D. Yin, J. Lin, X. Zou, Z. Shao, H. Yang, and J. Tang, "CogView: Mastering text-to-image generation via transformers," in *Advances in Neural Information Processing Systems*, vol. 34, 2021, pp. 19 822–19 835.

[68] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, vol. 30, 2017.

[69] A. Borji, "Qualitative failures of image generation models and their application in detecting deepfakes," *arXiv preprint*, vol. arXiv:2304.06470, 2023.

[70] H. Farid, "Lighting (in)consistency of paint by text," *arXiv preprint*, vol. arXiv:2207.13744, 2022.

[71] ——, "Perspective (in)consistency of paint by text," *arXiv preprint*, vol. arXiv:2206.14617, 2022.

[72] M. Bertamini, A. Spooner, and H. Hecht, "Naive optics: Predicting and perceiving reflections in mirrors," *Journal of Experimental Psychology: Human Perception and Performance*, vol. 29, no. 5, pp. 982–1002, 2003.

[73] Y. Ostrovsky, P. Cavanagh, and P. Sinha, "Perceiving illumination inconsistencies in scenes," *Perception*, vol. 34, no. 11, pp. 1301–1314, 2005.

[74] T. J. de Carvalho, C. Riess, E. Angelopoulou, H. Pedrini, and A. de Rezende Rocha, "Exposing digital image forgeries by illumination color classification," *IEEE Transactions on Information Forensics and Security*, vol. 8, no. 7, pp. 1182–1194, 2013.

[75] S. J. Nightingale and H. Farid, "AI-synthesized faces are indistinguishable from real faces and more trustworthy," *Proceedings of the National Academy of Sciences*, vol. 119, no. 8, 2022.

[76] F. Marra, D. Gragnaniello, L. Verdoliva, and G. Poggi, "Do GANs leave artificial fingerprints?" in *IEEE Conference on Multimedia Information Processing and Retrieval (MIPR)*, 2019, pp. 506–511.

[77] S. Mandelli, D. Cozzolino, E. D. Cannas, J. P. Cardenuto, D. Moreira, P. Bestagini, W. J. Scheirer, A. Rocha, L. Verdoliva, S. Tubaro, and E. J. Delp, "Forensic analysis of synthetically generated western blot images," *IEEE Access*, vol. 10, pp. 59 919–59 932, 2022.

[78] C. Kong, B. Chen, H. Li, S. Wang, A. Rocha, and S. Kwong, "Detect and locate: Exposing face manipulation by semantic-and noise-level telltales," *IEEE Transactions on Information Forensics and Security*, vol. 17, pp. 1741–1756, 2022.

[79] R. Corvi, D. Cozzolino, G. Poggi, K. Nagano, and L. Verdoliva, "Intriguing properties of synthetic images: from generative adversarial networks to diffusion models," *arXiv preprint*, vol. arXiv:2304.06408, 2023.

[80] D. Gragnaniello, D. Cozzolino, F. Marra, G. Poggi, and L. Verdoliva, "Are GAN generated images easy to detect? A critical analysis of the state-of-the-art," in *IEEE International Conference on Multimedia and Expo*, 2021, pp. 1–6.

[81] T. Osakabe, M. Tanaka, Y. Kinoshita, and H. Kiya, "CycleGAN without checkerboard artifacts for counter-forensics of fake-image detection," in *International Workshop on Advanced Imaging Technology (IWAIT)*, 2021, p. 1176609.

[82] D. Cozzolino, J. Thies, A. Rössler, M. Nießner, and L. Verdoliva, "SpoC: Spoofing camera fingerprints," in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 990–1000.

[83] X. Luo, R. Zhan, H. Chang, F. Yang, and P. Milanfar, "Distortion agnostic deep watermarking," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13 548–13 557.

[84] M. Ahmadi, A. Norouzi, N. Karimi, S. Samavi, and A. Emami, "ReDMark: Framework for residual diffusion watermarking based on deep networks," *Expert Systems with Applications*, vol. 146, p. 113157, 2020.

[85] P. Fernandez, G. Couairon, H. Jégou, M. Douze, and T. Furon, "The stable signature: Rooting watermarks in latent diffusion models," *arXiv preprint*, vol. arXiv:2303.15435, 2023.

[86] M. Ranzato, A. Szlam, J. Bruna, M. Mathieu, R. Collobert, and S. Chopra, "Video (language) modeling: a baseline for generative models of natural videos," *arXiv preprint*, vol. arXiv:1412.6604, 2014.

[87] C. Vondrick, H. Pirsiavash, and A. Torralba, "Generating videos with scene dynamics," in *Advances in Neural Information Processing Systems*, vol. 29, 2016.

[88] S. Tulyakov, M.-Y. Liu, X. Yang, and J. Kautz, "MoCoGAN: Decomposing motion and content for video generation," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1526–1535.

[89] Y. Tian, J. Ren, M. Chai, K. Olszewski, X. Peng, D. N. Metaxas, and S. Tulyakov, "A good image generator is what you need for high-resolution video synthesis," *arXiv preprint*, vol. arXiv:2104.15069, 2021.

[90] W. Yan, Y. Zhang, P. Abbeel, and A. Srinivas, "VideoGPT: Video generation using VQ-VAE and transformers," *arXiv preprint*, vol. arXiv:2104.10157, 2021.

[91] G. Le Moing, J. Ponce, and C. Schmid, "CCVS: context-aware controllable video synthesis," in *Advances in Neural Information Processing Systems*, vol. 34, 2021, pp. 14 042–14 055.

[92] A. Blattmann, R. Rombach, H. Ling, T. Dockhorn, S. W. Kim, S. Fidler, and K. Kreis, "Align your latents: High-resolution video synthesis with latent diffusion models," *arXiv preprint*, vol. arXiv:2304.08818, 2023.

[93] V. Voleti, A. Jolicoeur-Martineau, and C. Pal, "Masked conditional video diffusion for prediction, generation, and interpolation," *arXiv preprint*, vol. arXiv:2205.09853, 2022.

[94] R. Yang, P. Srivastava, and S. Mandt, "Diffusion probabilistic modeling for video generation," *arXiv preprint*, vol. arXiv:2203.09481, 2022.

[95] U. Singer, A. Polyak, T. Hayes, X. Yin, J. An, S. Zhang, Q. Hu, H. Yang, O. Ashual, O. Gafni, D. Parikh, S. Gupta, and Y. Taigman, "Make-a-video: Text-to-video generation without text-video data," in *International Conference on Learning Representations*, 2023.

[96] Y. Li, M. Min, D. Shen, D. Carlson, and L. Carin, "Video generation from text," in *AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.

[97] Y. Pan, Z. Qiu, T. Yao, H. Li, and T. Mei, "To create what you tell: Generating videos from captions," in *ACM International Conference on Multimedia*, 2017, pp. 1789–1798.

[98] J. Ho, W. Chan, C. Saharia, J. Whang, R. Gao, A. Gritsenko, D. P. Kingma, B. Poole, M. Norouzi, D. J. Fleet, and T. Salimans, "Imagen video: High definition video generation with diffusion models," *arXiv preprint*, vol. arXiv:2210.02303, 2022.

[99] J. An, S. Zhang, H. Yang, S. Gupta, J.-B. Huang, J. Luo, and X. Yin, "Latent-shift: Latent diffusion with temporal shift for efficient text-to-video generation," *arXiv preprint*, vol. arXiv:2304.08477, 2023.

[100] P. Esser, J. Chiu, P. Atighehchian, J. Granskog, and A. Germanidis, "Structure and content-guided video synthesis with diffusion models," *arXiv preprint*, vol. arXiv:2302.03011, 2023.

[101] Y. He, T. Yang, Y. Zhang, Y. Shan, and Q. Chen, "Latent video diffusion models for high-fidelity video generation with arbitrary lengths," *arXiv preprint*, vol. arXiv:2211.13221, 2022.

[102] L. Khachatryan, A. Movsisyan, V. Tadevosyan, R. Henschel, Z. Wang, S. Navasardyan, and H. Shi, "Text2Video-Zero: Text-to-image diffusion models are zero-shot video generators," *arXiv preprint*, vol. arXiv:2303.13439, 2023.

[103] M. Saito, E. Matsumoto, and S. Saito, "Temporal generative adversarial nets with singular value clipping," in *IEEE International Conference on Computer Vision*, 2017, pp. 2830–2839.

[104] N. Srivastava, E. Mansimov, and R. Salakhudinov, "Unsupervised learning of video representations using LSTMs," in *International Conference on Machine Learning*, 2015, pp. 843–852.

[105] I. Skorokhodov, S. Tulyakov, and M. Elhoseiny, "StyleGAN-V: A continuous video generator with the price, image quality and perks of StyleGAN2," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3626–3636.

[106] S. Yu, J. Tack, S. Mo, H. Kim, J. Kim, J.-W. Ha, and J. Shin, "Generating videos with dynamics-aware implicit generative adversarial networks," *arXiv preprint*, vol. arXiv:2202.10571, 2022.

[107] T. Höppe, A. Mehrjou, S. Bauer, D. Nielsen, and A. Dittadi, "Diffusion models for video prediction and infilling," *arXiv preprint*, vol. arXiv:2206.07696, 2022.

[108] J. Z. Wu, Y. Ge, X. Wang, W. Lei, Y. Gu, W. Hsu, Y. Shan, X. Qie, and M. Z. Shou, "Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation," *arXiv preprint*, vol. arXiv:2212.11565, 2022.

[109] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet, "Video diffusion models," *arXiv preprint*, vol. arXiv:2204.03458, 2022.

[110] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, "3D U-Net: Learning dense volumetric segmentation from sparse annotation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2016, pp. 424–432.

[111] S. Ge, S. Nah, G. Liu, T. Poon, A. Tao, B. Catanzaro, D. Jacobs, J.-B. Huang, M.-Y. Liu, and Y. Balaji, "Preserve your own correlation: A noise prior for video diffusion models," *arXiv preprint*, vol. arXiv:2305.10474, 2023.

[112] M. Mathieu, C. Couprie, and Y. LeCun, "Deep multi-scale video prediction beyond mean square error," *arXiv preprint*, vol. arXiv:1511.05440, 2015.

[113] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, G. Liu, A. Tao, J. Kautz, and B. Catanzaro, "Video-to-video synthesis," *arXiv preprint*, vol. arXiv:1808.06601, 2018.

[114] L. Yu, Y. Cheng, K. Sohn, J. Lezama, H. Zhang, H. Chang, A. G. Hauptmann, M.-H. Yang, Y. Hao, I. Essa, and L. Jiang, "MAGVIT: Masked generative video transformer," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2023, pp. 10 459–10 469.

[115] W. Hong, M. Ding, W. Zheng, X. Liu, and J. Tang, "CogVideo: Large-scale pretraining for text-to-video generation via transformers," *arXiv preprint*, vol. arXiv:2205.15868, 2022.

[116] M. Ding, W. Zheng, W. Hong, and J. Tang, "CogView2: Faster and better text-to-image generation via hierarchical transformers," *arXiv preprint*, vol. arXiv:2204.14217, 2022.

[117] R. Villegas, M. Babaeizadeh, P.-J. Kindermans, H. Moraldo, H. Zhang, M. T. Saffar, S. Castro, J. Kunze, and D. Erhan, "Phenaki: Variable length video generation from open domain textual description," *arXiv preprint*, vol. arXiv:2210.02399, 2022.

[118] D. Zhou, W. Wang, H. Yan, W. Lv, Y. Zhu, and J. Feng, "MagicVideo: Efficient video generation with latent diffusion models," *arXiv preprint*, vol. arXiv:2211.11018, 2022.

[119] S. Yu, K. Sohn, S. Kim, and J. Shin, "Video probabilistic diffusion models in projected latent space," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2023, pp. 18 456–18 466.

[120] W. Wang, H. Yang, Z. Tuo, H. He, J. Zhu, J. Fu, and J. Liu, "VideoFactory: Swap attention in spatiotemporal diffusions for text-to-video generation," *arXiv preprint*, vol. arXiv:2305.10874, 2023.

[121] D. Zhang, C. Li, F. Lin, D. Zeng, and S. Ge, "Detecting deepfake videos with temporal dropout 3DCNN," in *International Joint Conference on Artificial Intelligence*, 2021, pp. 1288–1294.

[122] I. Amerini and R. Caldelli, "Exploiting prediction error inconsistencies through LSTM-based classifiers to detect deepfake videos," in *ACM Workshop on Information Hiding and Multimedia Security*, 2020, pp. 97–102.

[123] A. Chintha, B. Thai, S. J. Sohrawardi, K. Bhatt, A. Hickerson, M. Wright, and R. Ptucha, "Recurrent convolutional structures for audio spoof and video deepfake detection," *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 5, pp. 1024–1037, 2020.

[124] T. Mittal, U. Bhattacharya, R. Chandra, A. Bera, and D. Manocha, "Emotions don't lie: An audio-visual deepfake detection method using affective cues," in *ACM International Conference on Multimedia*, 2020, pp. 2823–2832.

[125] A. Haliassos, K. Vougioukas, S. Petridis, and M. Pantic, "Lips don't lie: A generalisable and robust approach to face forgery detection,"

in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 5039–5049.

[126] H. Qi, Q. Guo, F. Juefei-Xu, X. Xie, L. Ma, W. Feng, Y. Liu, and J. Zhao, "DeepRhythm: Exposing deepfakes with attentional visual heartbeat rhythms," in *ACM International Conference on Multimedia*, 2020, pp. 4318–4327.

[127] X. Yang, Y. Li, and S. Lyu, "Exposing deep fakes using inconsistent head poses," in *International Conference on Acoustics, Speech and Signal Processing*, 2019, pp. 8261–8265.

[128] Z. Cai, Z. Xiong, H. Xu, P. Wang, W. Li, and Y. Pan, "Generative adversarial networks: A survey toward private and secure applications," *ACM Computing Surveys*, vol. 54, no. 6, pp. 1–38, 2021.

[129] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "WaveNet: A generative model for raw audio," *arXiv preprint*, vol. arXiv:1609.03499, 2016.

[130] W. Ping, K. Peng, and J. Chen, "ClariNet: Parallel wave generation in end-to-end text-to-speech," *arXiv preprint*, vol. arXiv:1807.07281, 2018.

[131] Y. Ren, Y. Ruan, X. Tan, T. Qin, S. Zhao, Z. Zhao, and T.-Y. Liu, "FastSpeech: Fast, robust and controllable text to speech," in *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[132] N. Li, S. Liu, Y. Liu, S. Zhao, and M. Liu, "Neural speech synthesis with transformer network," in *AAAI Conference on Artificial Intelligence*, vol. 33, no. 1, 2019, pp. 6706–6713.

[133] H. Guo, F. Xie, F. K. Soong, X. Wu, and H. Meng, "A multi-stage multi-codebook VQ-VAE approach to high-performance neural TTS," *arXiv preprint*, vol. arXiv:2209.10887, 2022.

[134] W. Hsu, Y. Zhang, R. J. Weiss, H. Zen, Y. Wu, Y. Wang, Y. Cao, Y. Jia, Z. Chen, J. Shen, P. Nguyen, and R. Pang, "Hierarchical generative modeling for controllable speech synthesis," *arXiv preprint*, vol. arXiv:1810.07217, 2018.

[135] N. Chen, Y. Zhang, H. Zen, R. J. Weiss, M. Norouzi, and W. Chan, "WaveGrad: Estimating gradients for waveform generation," *arXiv preprint*, vol. arXiv:2009.00713, 2020.

[136] Z. Kong, W. Ping, J. Huang, K. Zhao, and B. Catanzaro, "DiffWave: A versatile diffusion model for audio synthesis," *arXiv preprint*, vol. arXiv:2009.09761, 2020.

[137] V. Popov, I. Vovk, V. Gogoryan, T. Sadekova, and M. Kudinov, "Grad-TTS: A diffusion probabilistic model for text-to-speech," in *International Conference on Machine Learning*, 2021, pp. 8599–8608.

[138] M. Jeong, H. Kim, S. J. Cheon, B. J. Choi, and N. S. Kim, "Diff-TTS: A denoising diffusion model for text-to-speech," *arXiv preprint*, vol. arXiv:2104.01409, 2021.

[139] S.-W. Fu, C.-F. Liao, Y. Tsao, and S.-D. Lin, "MetricGAN: Generative adversarial networks based black-box metric scores optimization for speech enhancement," in *International Conference on Machine Learning*, 2019, pp. 2031–2041.

[140] J. Lin, S. Niu, Z. Wei, X. Lan, A. J. Wijngaarden, M. C. Smith, and K.-C. Wang, "Speech enhancement using forked generative adversarial networks with spectral subtraction," in *Interspeech*, 2019, pp. 3163–3167.

[141] J.-M. Lemercier, J. Richter, S. Welker, and T. Gerkmann, "StoRM: A diffusion-based stochastic regeneration model for speech enhancement and dereverberation," *arXiv preprint*, vol. arXiv:2212.11851, 2022.

[142] Y.-J. Lu, Z.-Q. Wang, S. Watanabe, A. Richard, C. Yu, and Y. Tsao, "Conditional diffusion probabilistic model for speech enhancement," in *International Conference on Acoustics, Speech and Signal Processing*, 2022, pp. 7402–7406.

[143] H. Yen, F. G. Germain, G. Wichern, and J. Le Roux, "Cold diffusion for speech enhancement," in *International Conference on Acoustics, Speech and Signal Processing*, 2023, pp. 1–5.

[144] W.-C. Huang, T. Hayashi, Y.-C. Wu, H. Kameoka, and T. Toda, "Voice transformer network: Sequence-to-sequence voice conversion using transformer with text-to-speech pretraining," *arXiv preprint*, vol. arXiv:1912.06813, 2019.

[145] T. Kaneko, H. Kameoka, K. Tanaka, and N. Hojo, "CycleGAN-VC2: Improved cyclegan-based non-parallel voice conversion," in *International Conference on Acoustics, Speech and Signal Processing*, 2019, pp. 6820–6824.

[146] C.-C. Hsu, H.-T. Hwang, Y.-C. Wu, Y. Tsao, and H.-M. Wang, "Voice conversion from unaligned corpora using variational autoencoding Wasserstein generative adversarial networks," *arXiv preprint*, vol. arXiv:1704.00849, 2017.

[147] H. Kameoka, T. Kaneko, K. Tanaka, and N. Hojo, "ACVAE-VC: Non-parallel voice conversion with auxiliary classifier variational autoencoder," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, no. 9, pp. 1432–1443, 2019.

[148] H.-S. Choi, J. Lee, W. Kim, J. Lee, H. Heo, and K. Lee, "Neural analysis and synthesis: Reconstructing speech from self-supervised representations," in *Advances in Neural Information Processing Systems*, vol. 34, 2021, pp. 16 251–16 265.

[149] A. Richard, D. Markovic, I. D. Gebru, S. Krenn, G. A. Butler, F. Torre, and Y. Sheikh, "Neural synthesis of binaural speech from mono audio," in *International Conference on Learning Representations*, 2021.

[150] Y. Leng, Z. Chen, J. Guo, H. Liu, J. Chen, X. Tan, D. Mandic, L. He, X.-Y. Li, T. Qin, S. Zhao, and T.-Y. Liu, "BinauralGrad: A two-stage conditional diffusion probabilistic model for binaural audio synthesis," in *Advances in Neural Information Processing Systems*, vol. 35, 2022, pp. 23 689–23 700.

[151] K. Bhagtani, A. K. S. Yadav, E. R. Bartusiak, Z. Xiang, R. Shao, S. Baireddy, and E. J. Delp, "An overview of recent work in media forensics: Methods and threats," *arXiv preprint*, vol. arXiv:2204.12067, 2022.

[152] J. Yang, R. K. Das, and H. Li, "Extended constant-Q cepstral coefficients for detection of spoofing attacks," in *Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, 2018, pp. 1024–1029.

[153] M. Todisco, H. Delgado, and N. W. Evans, "A new feature for automatic speaker verification anti-spoofing: Constant Q cepstral coefficients," in *The Speaker and Language Recognition Workshop (Odyssey)*, 2016, pp. 283–290.

[154] Z. Wu, T. Kinnunen, N. Evans, J. Yamagishi, C. Hanilçi, M. Sahidullah, and A. Sizov, "ASVspoof 2015: the first automatic speaker verification spoofing and countermeasures challenge," in *Interspeech*, 2015.

[155] A. Mehrish, N. Majumder, R. Bharadwaj, R. Mihalcea, and S. Poria, "A review of deep learning techniques for speech processing," *Information Fusion*, p. 101869, 2023.

[156] A. Fathan, J. Alam, and W. H. Kang, "Mel-spectrogram image-based end-to-end audio deepfake detection under channel-mismatched conditions," in *IEEE International Conference on Multimedia and Expo*, 2022, pp. 1–6.

[157] S.-Y. Lim, D.-K. Chae, and S.-C. Lee, "Detecting deepfake voice using explainable deep learning techniques," *Applied Sciences*, vol. 12, no. 8, p. 3926, 2022.

[158] T. Chen, A. Kumar, P. Nagarsheth, G. Sivaraman, and E. Khoury, "Generalization of audio deepfake detection," in *The Speaker and Language Recognition Workshop (Odyssey)*, 2020, pp. 132–137.

[159] X. Tan, T. Qin, F. Soong, and T.-Y. Liu, "A survey on neural speech synthesis," *arXiv preprint*, vol. arXiv:2106.15561, 2021.

[160] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Language models are unsupervised multitask learners," *OpenAI*, 2019.

[161] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 1877–1901.

[162] J. W. Rae, S. Borgeaud, T. Cai, K. Millican, J. Hoffmann, H. F. Song, J. Aslanides, S. Henderson, R. Ring, S. Young, E. Rutherford *et al.*, "Scaling language models: Methods, analysis & insights from training Gopher," *arXiv preprint*, vol. arXiv:2112.11446, 2021.

[163] N. Du, Y. Huang, A. M. Dai, S. Tong, D. Lepikhin, Y. Xu, M. Krikun, Y. Zhou, A. W. Yu, O. Firat *et al.*, "GLaM: Efficient scaling of language models with mixture-of-experts," in *International Conference on Machine Learning*, 2022, pp. 5547–5569.

[164] R. Thoppilan, D. de Freitas, J. Hall, N. Shazeer, A. Kulshreshtha, H. Cheng, A. Jin, T. Bos, L. Baker, Y. Du *et al.*, "LaMDA: Language models for dialog applications," *arXiv preprint*, vol. arXiv:2201.08239, 2022.

[165] S. Smith, M. Patwary, B. Norick, P. LeGresley, S. Rajbhandari, J. Casper, Z. Liu, S. Prabhumoye, G. Zerveas, V. Korthikanti *et al.*, "Using DeepSpeed and Megatron to train Megatron-Turing NLG 530B, a large-scale generative language model," *arXiv preprint*, vol. arXiv:2201.11990, 2022.

[166] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. d. L. Casas, L. A. Hendricks, J. Welbl, A. Clark *et al.*, "Training compute-optimal large language models," *arXiv preprint*, vol. arXiv:2203.15556, 2022.

[167] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann *et al.*, "PaLM: Scaling language modeling with pathways," *arXiv preprint*, vol. arXiv:2204.02311, 2022.

[168] S. Wu, O. Irsoy, S. Lu, V. Dabravolski, M. Dredze, S. Gehrmann, P. Kambadur, D. Rosenberg, and G. Mann, "BloombergGPT: A large language model for finance," *arXiv preprint*, vol. arXiv:2303.17564, 2023.

[169] O. Lieber, O. Sharir, B. Lenz, and Y. Shoham, "Jurassic-1: Technical details and evaluation," *AI21*, 2021.

[170] B. Wang and A. Komatsuzaki, "GPT-J-6B: A 6 billion parameter autoregressive language model," https://github.com/kingoflolz/mesh-transformer-jax, 2021, (Accessed on 05/16/2023).

[171] S. Black, S. Biderman, E. Hallahan, Q. Anthony, L. Gao, L. Golding, H. He, C. Leahy, K. McDonell, J. Phang *et al.*, "GPT-NeoX-20B: An open-source autoregressive language model," in *ACL Workshop on Challenges & Perspectives in Creating Large Language Models*, 2022, pp. 95–136.

[172] T. L. Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilić, D. Hesslow, R. Castagné, A. S. Luccioni, F. Yvon, M. Gallé *et al.*, "BLOOM: A 176b-parameter open-access multilingual language model," *arXiv preprint*, vol. arXiv:2211.05100, 2022.

[173] S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin *et al.*, "OPT: Open pre-trained transformer language models," *arXiv preprint*, vol. arXiv:2205.01068, 2022.

[174] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, "LLaMA: Open and efficient foundation language models," *arXiv preprint*, vol. arXiv:2302.13971, 2023.

[175] A. Köpf, Y. Kilcher, D. von Rütte, S. Anagnostidis, Z.-R. Tam, K. Stevens, A. Barhoum, N. M. Duc, O. Stanley, R. Nagyfi *et al.*, "OpenAssistant Conversations – democratizing large language model alignment," *arXiv preprint*, vol. arXiv:2304.07327, 2023.

[176] N. Subbaraman, "ChatGPT will see you now: Doctors using AI to answer patient questions," https://www.wsj.com/articles/dr-chatgpt-physicians-are-sending-patients-advice-using-ai-945cf60b, 2023, (Accessed on 06/06/2023).

[177] F. Manjoo, "ChatGPT is already changing how I do my job," https://www.nytimes.com/2023/04/21/opinion/chatgpt-journalism.html, 2023, (Accessed on 06/06/2023).

[178] T. H. Tran, "A doctor published several research papers with breakneck speed. chatgpt wrote them all." https://www.thedailybeast.com/how-this-doctor-wrote-dozens-of-science-papers-with-chatgpt, 2023, (Accessed on 06/06/2023).

[179] H. Alkaissi and S. I. McFarlane, "Artificial hallucinations in ChatGPT: Implications in scientific writing," *Cureus*, vol. 15, no. 2, 2023.

[180] E. Crothers, N. Japkowicz, and H. L. Viktor, "Machine generated text: A comprehensive survey of threat models and detection methods," *arXiv preprint*, vol. arXiv:2210.07321, 2022.

[181] A. Giaretta and N. Dragoni, "Community targeted phishing," in *International Conference in Software Engineering for Defence Applications*, 2020, pp. 86–93.

[182] R. Schuster, C. Song, E. Tromer, and V. Shmatikov, "You autocomplete me: Poisoning vulnerabilities in neural code completion," in *USENIX Security Symposium*, 2021, pp. 1559–1575.

[183] Y. Kilcher, "GPT-4chan: This is the worst AI ever," https://www.youtube.com/watch?v=efPrtcLdcdM, 2022, (Accessed on 06/06/2023).

[184] A. Kurenkov, "Lessons from the GPT-4Chan controversy," https://thegradient.pub/gpt-4chan-lessons/, 2022, (Accessed on 06/06/2023).

[185] H. Bai, J. G. Voelkel, johannes Christopher Eichstaedt, and R. Willer, "Artificial intelligence can persuade humans on political issues," *OSF Preprints*, 2023.

[186] J. Stribling, M. Krohn, and D. Aguayo, "SCIgen - an automatic CS paper generator," https://pdos.csail.mit.edu/archive/scigen/, 2005, (Accessed on 06/06/2023).

[187] J. A. Byrne and J. Christopher, "Digital magic, or the dark arts of the 21st century—how can journals and peer reviewers detect manuscripts and publications from paper mills?" *FEBS Letters*, vol. 594, no. 4, pp. 583–589, 2020.

[188] R. Zellers, A. Holtzman, H. Rashkin, Y. Bisk, A. Farhadi, F. Roesner, and Y. Choi, "Defending against neural fake news," in *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[189] E. Mitchell, Y. Lee, A. Khazatsky, C. D. Manning, and C. Finn, "DetectGPT: Zero-shot machine-generated text detection using probability curvature," *arXiv preprint*, vol. arXiv:2301.11305, 2023.

[190] S. Mitrović, D. Andreoletti, and O. Ayoub, "ChatGPT or human? detect and explain. explaining decisions of machine learning model for detecting short ChatGPT-generated text," *arXiv preprint*, vol. arXiv:2301.13852, 2023.

[191] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems*, vol. 30, 2017.

[192] J. P. Wahle, T. Ruas, F. Kirstein, and B. Gipp, "How large language models are transforming machine-paraphrased plagiarism," *arXiv preprint*, vol. arXiv:2210.03568, 2022.

[193] S. Svrluga, "Was that essay written by AI? A student made an app that might tell you," https://www.washingtonpost.com/education/2023/01/12/gptzero-chatgpt-detector-ai/, 2023, (Accessed on 06/06/2023).

[194] P. Yu, J. Chen, X. Feng, and Z. Xia, "CHEAT: A large-scale dataset for detecting ChatGPT-writtEn AbsTracts," *arXiv preprint*, vol. arXiv:2304.12008, 2023.

[195] A. Grinbaum and L. Adomaitis, "The ethical need for watermarks in machine-generated language," *arXiv preprint*, vol. arXiv:2209.03118, 2022.

[196] J. Kirchenbauer, J. Geiping, Y. Wen, J. Katz, I. Miers, and T. Goldstein, "A watermark for large language models," *arXiv preprint*, vol. arXiv:2301.10226, 2023.

[197] C. Bail, L. Pinheiro, and J. Royer, "Difficulty of detecting AI content poses legal challenges," https://www.law360.com/articles/1593766/difficulty-of-detecting-ai-content-poses-legal-challenges, 2023, (Accessed on 06/06/2023).

[198] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "NeRF: Representing scenes as neural radiance fields for view synthesis," in *European Conference on Computer Vision*, 2020, pp. 405–421.

[199] J. Ichnowski, Y. Avigal, J. Kerr, and K. Goldberg, "Dex-NeRF: Using a neural radiance field to grasp transparent objects," in *Conference on Robot Learning (CoRL)*, 2021.

[200] S. Zhi, T. Laidlow, S. Leutenegger, and A. J. Davison, "In-place scene labelling and understanding with implicit scene representation," in *IEEE International Conference on Computer Vision*, 2021, pp. 15 818–15 827.

[201] S. Vora, N. Radwan, K. Greff, H. Meyer, K. Genova, M. S. Sajjadi, E. Pot, A. Tagliasacchi, and D. Duckworth, "NeSF: Neural semantic fields for generalizable semantic segmentation of 3d scenes," *arXiv preprint*, vol. arXiv:2111.13260, 2021.

[202] X. Huang, Q. Zhang, Y. Feng, H. Li, X. Wang, and Q. Wang, "HDR-NeRF: High dynamic range neural radiance fields," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18 377–18 387.

[203] Y.-C. Guo, D. Kang, L. Bao, Y. He, and S.-H. Zhang, "NeRFReN: Neural radiance fields with reflections," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18 388–18 397.

[204] L. Ma, X. Li, J. Liao, Q. Zhang, X. Wang, J. Wang, and P. V. Sander, "Deblur-NeRF: Neural radiance fields from blurry images," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 851–12 860.

[205] R. Martin-Brualla, N. Radwan, M. S. Sajjadi, J. T. Barron, A. Dosovitskiy, and D. Duckworth, "NeRF in the Wild: Neural radiance fields for unconstrained photo collections," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 7210–7219.

[206] X. Chen, Q. Zhang, X. Li, Y. Chen, Y. Feng, X. Wang, and J. Wang, "Hallucinated neural radiance fields in the wild," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 933–12 942.

[207] L. De Luigi, D. Bolognini, F. Domeniconi, D. De Gregorio, M. Poggi, and L. Di Stefano, "ScanNeRF: a scalable benchmark for neural radiance fields," in *IEEE Winter Conference on Applications of Computer Vision*, 2023, pp. 816–825.

[208] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman, "Mip-NeRF 360: Unbounded anti-aliased neural radiance fields," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5470–5479.

[209] Z. Wang, S. Wu, W. Xie, M. Chen, and V. A. Prisacariu, "NeRF−−: Neural radiance fields without known camera parameters," *arXiv preprint*, vol. arXiv:2102.07064, 2021.

[210] C.-H. Lin, W.-C. Ma, A. Torralba, and S. Lucey, "BARF: Bundle-adjusting neural radiance fields," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 5721–5731.

[211] C. Zhu, R. Wan, Y. Tang, and B. Shi, "Occlusion-free scene recovery via neural radiance fields," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2023, pp. 20 722–20 731.

[212] C. Zhu, R. Wan, and B. Shi, "Neural transmitted radiance fields," in *Advances in Neural Information Processing Systems*, vol. 35, 2022, pp. 38 994–39 006.

[213] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM Transactions on Graphics*, vol. 41, no. 4, pp. 1–15, 2022.

[214] A. Pumarola, E. Corona, G. Pons-Moll, and F. Moreno-Noguer, "D-NeRF: Neural radiance fields for dynamic scenes," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 10 318–10 327.

[215] Y. Wang, R. Wan, W. Yang, H. Li, L.-P. Chau, and A. Kot, "Low-light image enhancement with normalizing flow," in *AAAI Conference on Artificial Intelligence*, vol. 36, no. 3, 2022, pp. 2604–2612.

[216] R. Wan, B. Shi, H. Li, L.-Y. Duan, and A. C. Kot, "Reflection scene separation from a single image," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2398–2406.

[217] ——, "Face image reflection removal," *International Journal of Computer Vision*, vol. 129, pp. 385–399, 2021.

[218] D. Ma, R. Wan, B. Shi, A. C. Kot, and L.-Y. Duan, "Learning to jointly generate and separate reflections," in *IEEE International Conference on Computer Vision*, 2019, pp. 2444–2452.

[219] R. Wan, B. Shi, H. Li, L.-Y. Duan, A.-H. Tan, and A. C. Kot, "CoRRN: Cooperative reflection removal network," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 12, pp. 2969–2982, 2019.

[220] R. Wan, B. Shi, W. Yang, B. Wen, L.-Y. Duan, and A. C. Kot, "Purifying low-light images via near-infrared enlightened image," *IEEE Transactions on Multimedia*, pp. 1–13, 2022.

[221] S. S. Usmani, M. Sharath, and M. Mehendale, "Future of mental health in the metaverse," *General Psychiatry*, vol. 35, no. 4, p. e100825, 2022.

[222] C. Lutz, "Digital inequalities in the age of artificial intelligence and big data," *Human Behavior and Emerging Technologies*, vol. 1, no. 2, pp. 141–148, 2019.

[223] Y. Xing, W. He, J. Z. Zhang, and G. Cao, "AI privacy opinions between US and chinese people," *Journal of Computer Information Systems*, vol. 63, no. 3, pp. 492–506, 2023.

[224] Y. Wang, Z. Su, N. Zhang, R. Xing, D. Liu, T. H. Luan, and X. Shen, "A survey on metaverse: Fundamentals, security, and privacy," *IEEE Communications Surveys and Tutorials*, vol. 25, no. 1, pp. 319–352, 2022.

[225] N.-A. Le-Khac and K.-K. R. Choo, *Databases in Digital Forensics*, ser. A Practical Hands-on Approach to Database Forensics. Springer International Publishing, 2022, pp. 1–2.

[226] D. P. Kingma and M. Welling, "An introduction to variational autoencoders," *Foundations and Trends® in Machine Learning*, vol. 12, no. 4, pp. 307–392, 2019.

[227] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in *International Conference on Machine Learning*, 2015, pp. 2256–2265.

[228] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, "Score-based generative modeling through stochastic differential equations," *arXiv preprint*, vol. arXiv2011.13456, 2020.

[229] O. Mayer and M. C. Stamm, "Accurate and efficient image forgery detection using lateral chromatic aberration," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 7, pp. 1762–1777, 2018.

[230] P. Ferrara, T. Bianchi, A. De Rosa, and A. Piva, "Image forgery localization via fine-grained analysis of CFA artifacts," *IEEE Transactions on Information Forensics and Security*, vol. 7, no. 5, pp. 1566–1577, 2012.

[231] S. Lyu, X. Pan, and X. Zhang, "Exposing region splicing forgeries with blind local noise estimation," *International Journal of Computer Vision*, vol. 110, pp. 202–221, 2014.

[232] Z. Fan and R. L. De Queiroz, "Identification of bitmap compression history: JPEG detection and quantizer estimation," *IEEE Transactions on Image Processing*, vol. 12, no. 2, pp. 230–235, 2003.

[233] P. Zhou, X. Han, V. I. Morariu, and L. S. Davis, "Learning rich features for image manipulation detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1053–1061.

[234] M. Wang, X. Fu, J. Liu, and Z.-J. Zha, "JPEG compression-aware image forgery localization," in *ACM International Conference on Multimedia*, 2022, pp. 5871–5879.

[235] L. Zhuo, S. Tan, B. Li, and J. Huang, "Self-adversarial training incorporating forgery attention for image forgery localization," *IEEE Transactions on Information Forensics and Security*, vol. 17, pp. 819–834, 2022.

[236] C. Dong, X. Chen, R. Hu, J. Cao, and X. Li, "MVSS-Net: Multiview multi-scale supervised networks for image manipulation detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 3, pp. 3539–3553, 2022.

[237] F. Matern, C. Riess, and M. Stamminger, "Exploiting visual artifacts to expose deepfakes and face manipulations," in *IEEE Winter Conference on Applications of Computer Vision Workshops*, 2019, pp. 83–92.

[238] S. McCloskey and M. Albright, "Detecting GAN-generated imagery using saturation cues," in *IEEE International Conference on Image Processing*, 2019, pp. 4584–4588.

[239] N. Yu, L. S. Davis, and M. Fritz, "Attributing fake images to GANs: Learning and analyzing GAN fingerprints," in *IEEE International Conference on Computer Vision*, 2019, pp. 7556–7566.

[240] R. Durall, M. Keuper, and J. Keuper, "Watch your up-convolution: CNN based generative deep neural networks are failing to reproduce spectral distributions," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 7890–7899.

[241] S.-Y. Wang, O. Wang, R. Zhang, A. Owens, and A. A. Efros, "CNN-generated images are surprisingly easy to spot... for now," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8695–8704.

[242] X. Xuan, B. Peng, W. Wang, and J. Dong, "On the generalization of GAN image forensics," in *Chinese Conference on Biometric Recognition (CCBR)*, 2019, pp. 134–141.

[243] C. Frith, "Role of facial expressions in social interactions," *Philosophical Transactions of the Royal Society B: Biological Sciences*, vol. 364, no. 1535, pp. 3453–3458, 2009.

[244] C. Kong, S. Wang, and H. Li, "Digital and physical face attacks: Reviewing and one step further," *arXiv preprint*, vol. arXiv:2209.14692, 2022.

[245] P. Korshunov and S. Marcel, "Deepfakes: a new threat to face recognition? Assessment and detection," *arXiv preprint*, vol. arXiv:1812.08685, 2018.

[246] A. Rossler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Nießner, "Faceforensics++: Learning to detect manipulated facial images," in *IEEE International Conference on Computer Vision*, 2019, pp. 1–11.

[247] N. Dufour and A. Gully, "Contributing data to deepfake detection research," https://ai.googleblog.com/2019/09/contributing-data-to-deepfake-detection.html, 2019, (Accessed on 06/06/2023).

[248] B. Dolhansky, R. Howes, B. Pflaum, N. Baram, and C. C. Ferrer, "The deepfake detection challenge (DFDC) preview dataset," *arXiv preprint*, vol. arXiv:1910.08854, 2019.

[249] Y. Li, X. Yang, P. Sun, H. Qi, and S. Lyu, "Celeb-DF: A large-scale challenging dataset for deepfake forensics," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3204–3213.

[250] L. Jiang, R. Li, W. Wu, C. Qian, and C. C. Loy, "DeeperForensics-1.0: A large-scale dataset for real-world face forgery detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2889–2898.

[251] Y. He, B. Gan, S. Chen, Y. Zhou, G. Yin, L. Song, L. Sheng, J. Shao, and Z. Liu, "ForgeryNet: A versatile benchmark for comprehensive forgery analysis," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4360–4369.

[252] T. Zhou, W. Wang, Z. Liang, and J. Shen, "Face forensics in the wild," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 5778–5788.

[253] P. Kwon, J. You, G. Nam, S. Park, and G. Chae, "KoDF: A large-scale korean deepfake detection dataset," in *IEEE International Conference on Computer Vision*, 2021, pp. 10 744–10 753.

[254] H. Khalid, S. Tariq, M. Kim, and S. S. Woo, "FakeAVCeleb: a novel audio-video multimodal deepfake dataset," *arXiv preprint*, vol. arXiv:2108.05080, 2021.

[255] Y. Li, M.-C. Chang, and S. Lyu, "In Ictu Oculi: Exposing AI created fake videos by detecting eye blinking," in *IEEE International Workshop on Information Forensics and Security*, 2018, pp. 1–7.

[256] Y. Li and S. Lyu, "Exposing deepfake videos by detecting face warping artifacts," *arXiv preprint*, vol. arXiv:1811.00656, 2018.

[257] D. Afchar, V. Nozick, J. Yamagishi, and I. Echizen, "MesoNet: a compact facial video forgery detection network," in *IEEE International Workshop on Information Forensics and Security*, 2018, pp. 1–7.

[258] H. H. Nguyen, J. Yamagishi, and I. Echizen, "Capsule-Forensics: Using capsule networks to detect forged images and videos," in *International Conference on Acoustics, Speech and Signal Processing*, 2019, pp. 2307–2311.

[259] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1251–1258.

[260] M. Tan and Q. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning*, 2019, pp. 6105–6114.

[261] E. Sabir, J. Cheng, A. Jaiswal, W. AbdAlmageed, I. Masi, and P. Natarajan, "Recurrent convolutional strategies for face manipulation detection in videos," in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 80–87.

[262] Y.-J. Heo, Y.-J. Choi, Y.-W. Lee, and B.-G. Kim, "Deepfake detection scheme based on vision transformer and distillation," *arXiv preprint*, vol. arXiv:2104.01353, 2021.

[263] S. Chen, T. Yao, Y. Chen, S. Ding, J. Li, and R. Ji, "Local relation learning for face forgery detection," in *AAAI Conference on Artificial Intelligence*, vol. 35, no. 2, 2021, pp. 1081–1088.

[264] Z. Shang, H. Xie, Z. Zha, L. Yu, Y. Li, and Y. Zhang, "PRRNet: Pixel-region relation network for face forgery detection," *Pattern Recognition*, vol. 116, p. 107950, 2021.

[265] Y. Nirkin, L. Wolf, Y. Keller, and T. Hassner, "Deepfake detection based on discrepancies between faces and their context," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 10, pp. 6111–6121, 2021.

[266] X. Li, Y. Lang, Y. Chen, X. Mao, Y. He, S. Wang, H. Xue, and Q. Lu, "Sharp multiple instance learning for deepfake video detection," in *ACM International Conference on Multimedia*, 2020, pp. 1864–1872.

[267] H. Zhao, W. Zhou, D. Chen, T. Wei, W. Zhang, and N. Yu, "Multi-attentional deepfake detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 2185–2194.

[268] A. Kumar, A. Bhavsar, and R. Verma, "Detecting deepfakes with metric learning," in *International Workshop on Biometrics and Forensics (IWBF)*, 2020, pp. 1–6.

[269] Y. Qian, G. Yin, L. Sheng, Z. Chen, and J. Shao, "Thinking in frequency: Face forgery detection by mining frequency-aware clues," in *European Conference on Computer Vision*, 2020, pp. 86–103.

[270] C. Miao, Z. Tan, Q. Chu, N. Yu, and G. Guo, "Hierarchical frequency-assisted interactive networks for face manipulation detection," *IEEE Transactions on Information Forensics and Security*, vol. 17, pp. 3008–3021, 2022.

[271] U. A. Ciftci, I. Demir, and L. Yin, "FakeCatcher: Detection of synthetic portrait videos using biological signals," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2020.

[272] L. Li, J. Bao, T. Zhang, H. Yang, D. Chen, F. Wen, and B. Guo, "Face X-ray for more general face forgery detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5001–5010.

[273] Y. Luo, Y. Zhang, J. Yan, and W. Liu, "Generalizing face forgery detection with high-frequency features," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 16 317–16 326.

[274] H. Liu, X. Li, W. Zhou, Y. Chen, Y. He, H. Xue, W. Zhang, and N. Yu, "Spatial-phase shallow learning: rethinking face forgery detection in frequency domain," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 772–781.

[275] K. Shiohara and T. Yamasaki, "Detecting deepfakes with self-blended images," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18 720–18 729.

[276] T. Zhao, X. Xu, M. Xu, H. Ding, Y. Xiong, and W. Xia, "Learning self-consistency for deepfake detection," in *IEEE International Conference on Computer Vision*, 2021, pp. 15 023–15 033.

[277] J. Cao, C. Ma, T. Yao, S. Chen, S. Ding, and X. Yang, "End-to-end reconstruction-classification learning for face forgery detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 4113–4122.

[278] X. Zhu, H. Wang, H. Fei, Z. Lei, and S. Z. Li, "Face forgery detection by 3D decomposition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 2929–2939.

[279] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[280] S. A. Khan and H. Dai, "Video transformer for deepfake detection with incremental learning," in *ACM International Conference on Multimedia*, 2021, pp. 1821–1828.

[281] R. Prenger, R. Valle, and B. Catanzaro, "WaveGlow: A flow-based generative network for speech synthesis," in *International Conference on Acoustics, Speech and Signal Processing*, 2019, pp. 3617–3621.

[282] K. Kumar, R. Kumar, T. De Boissiere, L. Gestin, W. Z. Teoh, J. Sotelo, A. de Brébisson, Y. Bengio, and A. C. Courville, "MelGAN: Generative adversarial networks for conditional waveform synthesis," in *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[283] J. Donahue, S. Dieleman, M. Bińkowski, E. Elsen, and K. Simonyan, "End-to-end adversarial text-to-speech," *arXiv preprint*, vol. arXiv:2006.03575, 2020.

[284] R. Yamamoto, E. Song, and J.-M. Kim, "Parallel WaveGAN: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram," in *International Conference on Acoustics, Speech and Signal Processing*, 2020, pp. 6199–6203.

[285] K. Peng, W. Ping, Z. Song, and K. Zhao, "Non-autoregressive neural text-to-speech," in *International Conference on Machine Learning*, 2020, pp. 7586–7598.

[286] R. Huang, M. W. Lam, J. Wang, D. Su, D. Yu, Y. Ren, and Z. Zhao, "FastDiff: A fast conditional diffusion model for high-quality speech synthesis," in *International Joint Conference on Artificial Intelligence*, 2022, pp. 4157–4163.

[287] R. Naika, "An overview of automatic speaker verification system," in *International Conference on Intelligent Computing and Information and Communication (ICICC)*, 2018, pp. 603–610.

[288] J. Yamagishi, X. Wang, M. Todisco, M. Sahidullah, J. Patino, A. Nautsch, X. Liu, K. A. Lee, T. Kinnunen, N. Evans *et al.*, "ASVspoof 2021: accelerating progress in spoofed and deepfake speech detection," in *Automatic Speaker Verification and Spoofing Countermeasures Challenge*, 2021.

[289] Z. Cai, W. Wang, and M. Li, "Waveform boundary detection for partially spoofed audio," in *International Conference on Acoustics, Speech and Signal Processing*, 2023, pp. 1–5.

[290] C. Sun, S. Jia, S. Hou, and S. Lyu, "AI-synthesized voice detection using neural vocoder artifacts," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2023, pp. 904–912.

[291] B. Logan, "Mel frequency cepstral coefficients for music modeling," in *International Society for Music Information Retrieval Conference*, 2000.

[292] I. Solak, "The M-AILABS speech dataset," https://www.caito.de/2019/01/03/the-m-ailabs-speech-dataset/, 2019, (Accessed on 06/06/2023).

[293] Z. Zhang, X. Yi, and X. Zhao, "GAN based synthesized audio dataset," https://ieee-dataport.org/documents/gan-based-synthesized-audio-dataset, 2019, (Accessed on 06/06/2023).

[294] J. Yi, Y. Bai, J. Tao, Z. Tian, C. Wang, T. Wang, and R. Fu, "Half-truth: A partially fake audio detection dataset," *arXiv preprint*, vol. arXiv:2104.03617, 2021.

[295] D. M. Ballesteros, Y. Rodriguez, and D. Renza, "A dataset of histograms of original and fake voice recordings (H-Voice)," *Data in brief*, vol. 29, p. 105331, 2020.

[296] K. Daniel, *Thinking, fast and slow*. Penguin, 2012.

[297] W. J. Scheirer, A. de Rezende Rocha, A. Sapkota, and T. E. Boult, "Toward open set recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 7, pp. 1757–1772, 2012.

[298] A. Ferreira, S. C. Felipussi, C. Alfaro, P. Fonseca, J. E. Vargas-Munoz, J. A. Dos Santos, and A. Rocha, "Behavior knowledge space-based fusion for copy–move forgery detection," *IEEE Transactions on Image Processing*, vol. 25, no. 10, pp. 4729–4742, 2016.

## Biographies

**João Phillipe Cardenuto** received a degree in Bachelor in Computer Engineering and Computer Science at the University of Campinas (2019). Currently, he is pursuing a Ph.D. degree at the University of Campinas, Brazil. Due to his research work, he was honored with the Google Latin America Research Award (LARA) 2021. He also had the opportunity to collaborate with Google Research as an intern in 2022. His latest works have focused on image provenance analysis and the development of forensic algorithms aimed at detecting doctored scientific images. His research interests include media forensics, computer vision, machine learning, and scientific integrity.

**Jing Yang** is a Ph.D. student of the Artificial Intelligence Lab., Recod.ai, and is currently doing a research internship at the Ubiquitous Knowledge Processing (UKP) Lab. She received her Master's degree in Computer Science at Hunan University, China, and her Bachelor's in Information and Computing Science at Hubei University of Technology, China. Jing's research interests include natural language understanding, fact-checking, and forensics.

**Rafael Padilha** is a researcher associated with the Artificial Intelligence Lab. (Recod.ai) at the Institute of Computing, University of Campinas, Brazil. He received his Ph.D. in 2022 from the same university, with a joint research internship at the University of Kentucky, USA. He works for Microsoft Research, innovating on agriculture and sustainability under the Research for Industry team. His research interests lie in computer vision, machine learning, and digital forensics.

**Renjie Wan** received his BEng degree from the University of Electronic Science and Technology of China in 2012 and the Ph.D. degree from Nanyang Technological University, Singapore, in 2019. He is currently an Assistant Professor at Hong Kong Baptist University, Hong Kong. He is the outstanding reviewer of ICCV 2019 and the recipient of the Microsoft CRSF Award, VCIP 2020 Best Paper Award, and the Wallenberg-NTU Presidential Postdoctoral Fellowship.

**Daniel Moreira** received a Ph.D. degree in computer science from the University of Campinas, Brazil, in 2016. After working four years as a systems analyst with the Brazilian Federal Data Processing Service (SERPRO), he joined the University of Notre Dame for six years, first as a post-doctoral fellow and later as an assistant research professor. He is currently an assistant professor in the Department of Computer Science at Loyola University Chicago. He is also a member of the IEEE Information Forensics and Security Technical Committee (IFS-TC), 2021-2023 term, IEEE Signal Processing Society Education Center Editorial Board, 2022-2023 term, and associate editor of IEEE Transactions on Information Forensics and Security (T-IFS) and Elsevier Pattern Recognition journals. His research interests include media forensics, machine learning, computer vision, and biometrics.

**Haoliang Li** received his Ph.D. degree from Nanyang Technological University (NTU), Singapore in 2018. He is currently an assistant professor in Department of Electrical Engineering, City University of Hong Kong. His research mainly focuses on AI security, multimedia forensics and transfer learning. He received the Wallenberg-NTU presidential postdoc fellowship in 2019, doctoral innovation award in 2019, VCIP best paper award in 2020, Top 50 Chinese Young Scholars in AI+X 2022, and Stanford's top 2% most highly cited scientists in 2022.

**Shiqi Wang** received the B.S. degree in computer science from the Harbin Institute of Technology in 2008 and the Ph.D. degree in computer application technology from Peking University in 2014. From 2014 to 2016, he was a Post-Doctoral Fellow with the Department of Electrical and Computer Engineering, University of Waterloo, Waterloo, ON, Canada. From 2016 to 2017, he was a Research Fellow with the Rapid-Rich Object Search Laboratory, Nanyang Technological University, Singapore. He is currently an Assistant Professor with the Department of Computer Science, City University of Hong Kong. He has authored or coauthored more than 200 refereed journal articles/conference papers. His research interests include video compression, image/video quality assessment, and image/video search and analysis. He received the Best Paper Award from IEEE VCIP 2019, ICME 2019, IEEE Multimedia 2018, and PCM 2017. His coauthored article received the Best Student Paper Award in the IEEE ICIP 2018.

**Fernanda Andaló** is a researcher associated with the Artificial Intelligence Lab. (Recod.ai) at the Institute of Computing, University of Campinas, Brazil. Andaló received a Ph.D. in Computer Science from the same university in 2012, during which she was a research fellow at Brown University. She worked for Samsung as a researcher and was a postdoctoral researcher in collaboration with Motorola, from 2014 to 2018. Currently, she works at The LEGO Group devising machine learning solutions for digital products. She was the 2016-2017 Chair of the IEEE Women in Engineering South Brazil Section, and

is an elected member of the IEEE Information Forensics and Security Technical Committee. Her research interests include machine learning and computer vision.

**Sébastien Marcel** (IEEE Senior member) is a senior researcher at the Idiap Research Institute (Switzerland), he heads the Biometrics Security and Privacy group and conducts research on face recognition, speaker recognition, vein recognition, attack detection (presentation attacks, morphing attacks, deepfakes) and template protection. He is also Professor at the University de Lausanne (UNIL) at the School of Criminal Justice and lecturer at the Ecole Polytechnique Fédérale de Lausanne (EPFL). He received his Ph.D. degree in signal processing from Université de Rennes I in France (2000) at CNET, the research center of France Telecom (now Orange Labs). He is also the Director of the Swiss Center for Biometrics Research and Testing, which conducts certifications of biometric products. He is Associate Editor of IEEE Transactions on Biometrics and Identity Science. He was Associate Editor of IEEE Signal Processing Letters, Associate Editor of IEEE Transactions on Information Forensics and Security, a Guest Editor of the IEEE Transactions on Information Forensics and Security Special Issue on "Biometric Spoofing and Countermeasures", and Co-editor of the IEEE Signal Processing Magazine Special Issue on "Biometric Security and Privacy". He is also the lead Editor of the Springer Handbook of Biometrics Anti-Spoofing (Editions 1, 2 and 3).

**Anderson Rocha** is a full-professor for Artificial Intelligence and Digital Forensics at the Institute of Computing, University of Campinas (Unicamp), Brazil. He is the Director of the Artificial Intelligence Lab., Recod.ai, and was the Director of the Institute of Computing for the 2019-2023 term. He is an elected affiliate member of the Brazilian Academy of Sciences (ABC) and the Brazilian Academy of Forensic Sciences (ABC). He is a two-term elected member of the IEEE Information Forensics and Security Technical Committee (IFS-TC) and its chair for the 2019-2020 term. He is a Microsoft Research and a Google Research Faculty Fellow. In addition, in 2016, he has been awarded the Tan Chin Tuan (TCT) Fellowship, a recognition promoted by the Tan Chin Tuan Foundation in Singapore. Finally, he is ranked Top-2% among the most influential scientists worldwide, according to recent studies from Research.com and Standford/PlosOne.