

# Deepfake Media Generation and Detection in the Generative AI Era: A Survey and Outlook

Florinel-Alin Croitoru, Andrei-Iulian Hiji, Vlad Hondru, Nicolae Cătălin Ristea, Paul Irofti, Marius Popescu, Cristian Rusu, Radu Tudor Ionescu, Fahad Shahbaz Khan, *Senior Member, IEEE*, and Mubarak Shah, *Fellow, IEEE*

**Abstract**—With the recent advancements in generative modeling, the realism of deepfake content has been increasing at a steady pace, even reaching the point where people often fail to detect manipulated media content online, thus being deceived into various kinds of scams. In this paper, we survey deepfake generation and detection techniques, including the most recent developments in the field, such as diffusion models and Neural Radiance Fields. Our literature review covers all deepfake media types, comprising image, video, audio and multimodal (audio-visual) content. We identify various kinds of deepfakes, according to the procedure used to alter or generate the fake content. We further construct a taxonomy of deepfake generation and detection methods, illustrating the important groups of methods and the domains where these methods are applied. Next, we gather datasets used for deepfake detection and provide updated rankings of the best performing deepfake detectors on the most popular datasets. In addition, we develop a novel multimodal benchmark to evaluate deepfake detectors on out-of-distribution content. The results indicate that state-of-the-art detectors fail to generalize to deepfake content generated by unseen deepfake generators. Finally, we propose future directions to obtain robust and powerful deepfake detectors. Our project page and new benchmark are available at <https://github.com/CroitoruAlin/biodeep>.

**Index Terms**—deepfake, deepfake generation, deepfake detection, deepfake benchmark.



## 1 INTRODUCTION

DEEPAKE media comprises image, video or audio files that are digitally altered or generated from scratch with AI tools in order to impersonate real or non-existent people. The recent groundbreaking progress of generative AI methods [1]–[6] has enabled the creation of realistic deepfake media with very little effort [7]–[18]. Unfortunately, the generated deepfake media can be used by scammers to spread misinformation on social media platforms to achieve large-scale political manipulation, and to deceive individuals or companies into financial frauds.

In an age where misinformation can quickly spread through social media platforms, deepfakes pose a critical threat to public trust and democracy, especially due to their growing online exploitation. A recent analysis of the fraud trends indicates that the number of fraud cases based on deepfakes registered a 10× increase in 2023, with respect to

2022<sup>1</sup>. Another recent study found that about 70% of people are unable to distinguish between a real and a deepfake voice<sup>2</sup>. The growing quality and quantity of deepfakes raise significant concerns, particularly regarding online fraud and manipulation. To prevent the spread of deepfake media, researchers have developed a broad range of unimodal [19]–[23] or multimodal [24]–[26] methods for deepfake detection. However, deepfake detectors trained on media generated with a certain set of AI tools typically fail on deepfakes generated with a distinct set of tools [20]–[22]. This has led to a relentless race to develop more powerful and robust deepfake detectors.

To this end, we conduct a comprehensive survey on the recent developments in deepfake media generation and detection. We first define a set of deepfake categories, which are determined based on the procedure used to generate the deepfake content. We identify both domain-agnostic and domain-specific deepfake types, and explain what kind of deepfake media belongs to each category. We next build a taxonomy of deepfake generation and detection methods, which creates a multi-perspective hierarchical categorization based on the considered media types, the employed architectures and the targeted tasks. As shown in Figure 1, we first divide contributions by task, into generation and detection. For each task, we identify the employed architectures. For deepfake generation, we find that the most popular architectures are Generative Adversarial Networks (GANs) [8], [14]–[16], [27], [28] and denoising diffusion models [11]–[13], [18], [29]–[31]. To detect deepfakes, the

- F.A. Croitoru, A.I. Hiji, V. Hondru, N.C. Ristea, P. Irofti, M. Popescu, C. Rusu and R.T. Ionescu are with the Department of Computer Science, University of Bucharest, Bucharest, Romania. F.A. Croitoru, A.I. Hiji, V. Hondru, and N.C. Ristea have contributed equally. R.T. Ionescu is the corresponding author.  
E-mail: [raducu.ionescu@gmail.com](mailto:raducu.ionescu@gmail.com)
- F.S. Khan is with Mohamed bin Zayed University of Artificial Intelligence (MBZUAI), UAE, and Linköping University, Sweden.
- M. Shah is with the Center for Research in Computer Vision (CRCV), Department of Computer Science, University of Central Florida, Orlando, FL, US.

Copyright 2024 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, including reprinting/republishing this material for advertising or promotional purposes, collecting new collected works for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works. This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

1. Sumsud Expert Roundtable: The Top KYC Trends Coming in 2024
2. Artificial Imposters—Cybercriminals Turn to AI Voice Cloning for a New Breed of Scam

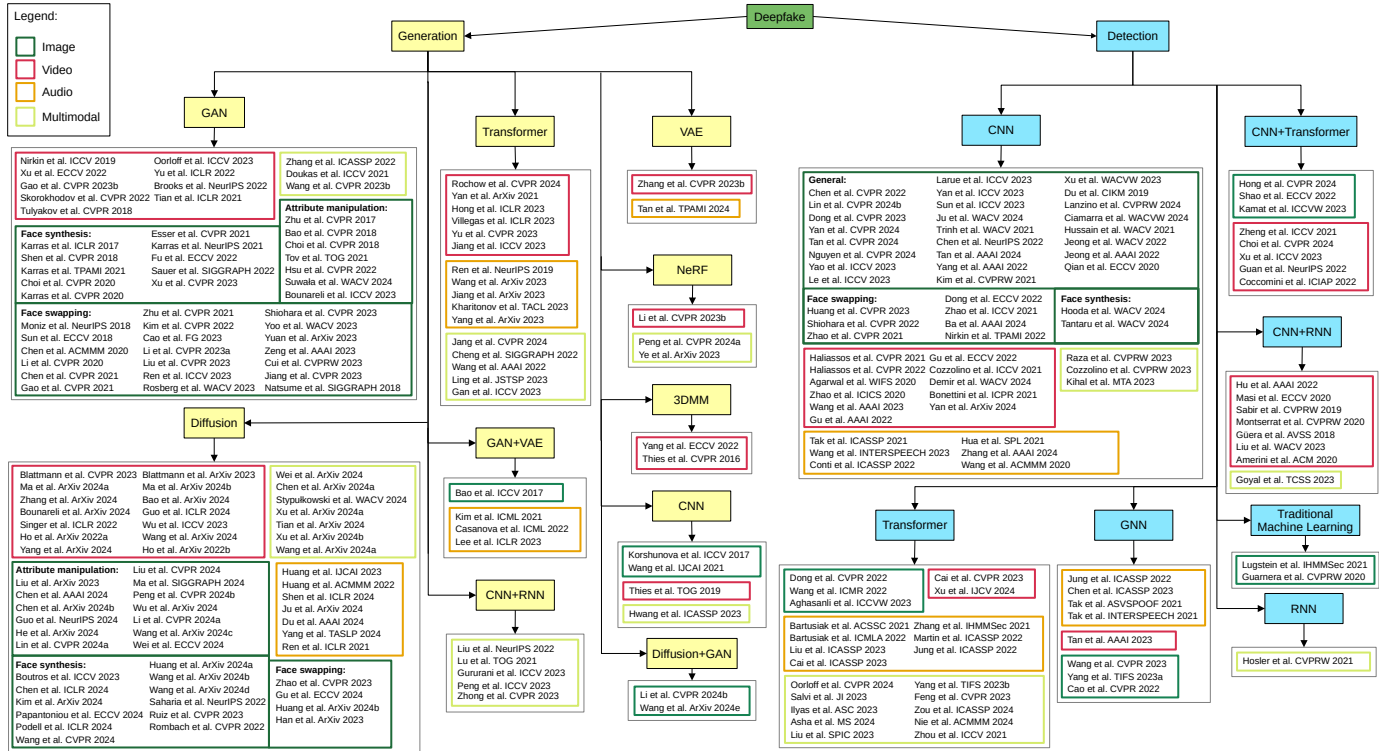


Fig. 1. A taxonomy of the state-of-the-art deepfake generation and detection methods. The methods are first divided according to the target task: generation versus detection. For each task, the methods are further divided into different kinds of architectures. For each architecture, we separate the methods based on the media types. Large groups are further divided according to the deepfake types presented in Section 3. References are clickable links to papers. Best viewed in color.

majority of methods are based on convolutional neural networks (CNNs) [19], [21], [24], [25], transformers [32]–[34], or hybrid architectures that combine CNNs either with transformers [35]–[37] or recurrent neural networks (RNNs) [38], [39]. For each type of architecture, we further divide the contributions with respect to the media types: image, video, audio or multimodal (audio-visual). Next, we present the main contributions in each category of articles included in the taxonomy. We further review existing datasets for deepfake detection in image, video and audio. We then aggregate the reported performance levels of deepfake detectors on the most popular datasets, thus facilitating a direct comparison of existing methods. In addition, we introduce a benchmark to test the generalization capacity of deepfake detectors to out-of-distribution content. Interestingly, we find that state-of-the-art deepfake detectors showcase poor generalization to realistic deepfakes generated by newer and more powerful generative models. Finally, we identify research gaps in current literature, proposing a series of future work directions that can lead to the development of better frameworks to detect deepfake media.

In summary, our contribution is fourfold:

- We conduct a comprehensive survey of deepfake generation and detection methods, comprising recent advancements in four domains: image, video, audio and multimodal.
- We construct a taxonomy of deepfake generation and detection methods, categorizing research articles according to tasks, architectures and media types.
- We collect and merge results reported on popular deepfake detection benchmarks, providing the

means to easily assess the current performance levels of deepfake detectors.

- We introduce a benchmark to test the out-of-domain generalization of deepfake detection models, showing that current detectors generally exhibit high performance drops on deepfakes generated by new and powerful generators.

## 2 RELATED SURVEYS

Several attempts have been made to survey deepfake detection and generation. In Table 1, we gather related surveys and illustrate the tasks, domains, methods and other aspects covered by the gathered surveys. Some surveys only cover the generation part [43], [45], while others are particularly focused on detection [40], [41], [44], [50]. Many surveys consider only one input media type, e.g. video [40], [42], [43], [45] or audio [44], [50]. There are a few surveys [41], [46], [47] that cover all media types (image, video, audio and multimodal), but only Masood *et al.* [46] and Patel *et al.* [47] address both detection and generation tasks. Although the surveys of Masood *et al.* [46] and Patel *et al.* [47] are comprehensive, they do not cover the most recent developments, such as diffusion models and vision transformers.

In summary, we find that existing surveys are either outdated or limited in terms of coverage, including only specific tasks (generation or detection) or media types (image, audio or video). In contrast, we conduct an *extensive survey of current literature*, covering both generation and detection, as well as all deepfake media types. Moreover, we create a multi-level taxonomy to ease the navigation through the current deepfake literature, providing direct links to the



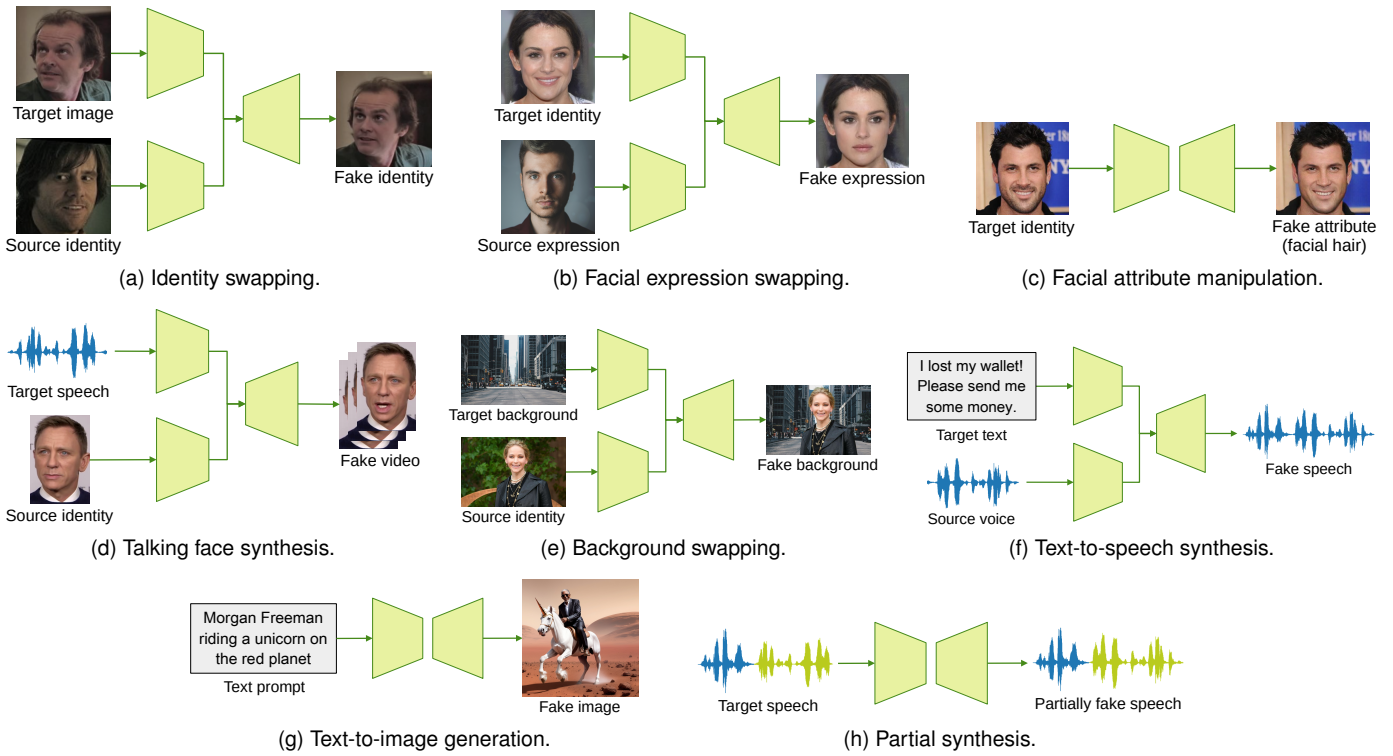


Fig. 2. Deepfake types according to the general procedure used to synthesize the fake content. For deepfake types that apply to multiple domains, we provide the illustration for only one domain. Best viewed in color.

TABLE 1

Comparing our survey with related surveys in terms of the covered tasks, domains, methods and other aspects. There are at least three factors that differentiate our survey from each of the other surveys.

Survey	Task		Domain				Method					Taxonomy	Datasets	New Benchmark	
	Generation	Detection	Image	Video	Audio	Multimodal	GANs	Diffusion	CNNs	RNNs	Transformers				Others
Das <i>et al.</i> [40]		✓		✓				✓	✓						
Heidari <i>et al.</i> [41]		✓	✓	✓	✓	✓		✓	✓			✓			
Kaur <i>et al.</i> [42]	✓	✓		✓			✓	✓	✓	✓	✓	✓	✓	✓	✓
Lei <i>et al.</i> [43]	✓			✓			✓	✓				✓		✓	
Li <i>et al.</i> [44]		✓			✓			✓	✓	✓	✓	✓	✓	✓	✓
Li <i>et al.</i> [45]	✓			✓			✓	✓	✓	✓	✓	✓	✓	✓	✓
Masood <i>et al.</i> [46]	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Patel <i>et al.</i> [47]	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Pei <i>et al.</i> [48]	✓	✓	✓	✓			✓	✓	✓	✓	✓	✓	✓	✓	✓
Seow <i>et al.</i> [49]	✓	✓	✓	✓			✓	✓	✓	✓	✓	✓	✓	✓	✓
Yi <i>et al.</i> [50]		✓			✓			✓	✓	✓	✓	✓	✓	✓	✓
Zhang [51]	✓	✓	✓	✓	✓		✓	✓	✓	✓	✓	✓	✓	✓	✓
Ours	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

referenced papers. To our knowledge, our survey is the first to propose a novel benchmark to test the generalization capacity of deepfake detectors to out-of-distribution data.

### 3 DEEPAKE TYPES

To date, a number of alternative procedures have been employed to generate deepfakes. In order to simplify the task of producing realistic deepfakes, one commonly used procedure is to only alter a certain aspect of a media file,

e.g. modifying the identity of a person in an existing video, or changing the emotion of a speech, while preserving the speech content and the speaker’s identity. By employing recent and powerful generative models [3], [52], generating deepfake content from scratch has also become prevalent. We further categorize the deepfake content according to the procedure employed to obtain the respective deepfake type. We illustrate the identified categories in Figure 2 and present them in detail below. Interestingly, we identify deepfake categories that are domain-agnostic (which have been applied to all media types) and domain-specific (which have only been applied to a certain media type).

**Identity swapping.** Deepfakes based on identity swapping imply replacing the identity information of a target person with that of a source person, while preserving identity-agnostic attributes, such as facial expressions. In the visual domain, this kind of deepfake is often referred to as *face swapping*, while in the audio domain, it is known as *voice conversion* or *voice swapping*. Voice conversion seeks to change the timbre and prosody of a speaker with those of another speaker, while preserving the content of the speech.

**Expression/emotion swapping.** In contrast to identity swapping, expression or emotion swapping involves altering the facial expression/emotion without changing the identity information. In the image domain, this task is known as *facial expression swapping*. In the video domain, the task is also known as *face reenactment*, and it implies altering the facial movement, which might often require facial motion capturing technology. In the audio domain, *emotion swapping* is the task of changing the emotion of speech, while retaining the speech content and the speaker’s identity.

**Facial attribute manipulation.** Producing deepfake image

and video content via facial attribute manipulation implies changing certain semantic attributes of a target face, while maintaining the identity information. Some of the attributes that are usually altered are age, gender, skin color and hair.

**Talking face synthesis.** Talking face synthesis is perhaps the most complex procedure to obtain deepfake audio-video content, but also the most flexible procedure. The task seeks to generate an audio-video file of a talking face, where the source character is engaged in the act of speech. The generated content is conditioned on some target information, provided in the form of text, audio, video or even multimodal content. The facial expressions, head movements, lip movements, speech emotions, and spoken content in the synthesized audio-video content are consistent with those of the source character.

**Background swapping.** Deepfakes based on background swapping are generated by changing the background scene with a different one. In the visual domain, this involves segmenting the source person and blending this person in a new scene. In the audio domain, the background sound of the original recording is replaced with another background sound using audio editing technologies, while preserving speech content and the speaker's identity.

**Text-to-speech synthesis.** Deepfake audio can be created with the help of a machine learning model for speech synthesis, starting from a piece of text. Current text-to-speech (TTS) synthesis models can produce natural speech and emulate the voice of a source identity.

**Text-to-image/video generation.** With the recent development of diffusion models, such as Stable Diffusion [3] and GLIDE [53], a new type of deepfake has emerged. Deepfake content can be easily generated by simply prompting a text-conditional diffusion model. All the necessary details, including the name of the source person, the facial attributes, the body pose and the actions can be specified through the prompt. This approach can be used to generate both images and videos. In text-to-speech synthesis, the input text is literally pronounced by the system, while in text-to-image/video generation, the prompt is rather interpreted by the system as a set of instructions.

**Partial synthesis.** As the name implies, partial synthesis involves changing only a part of an existing media file to create a deepfake. In the video domain, this kind of deepfake can be obtained by changing a subset of frames. In the audio domain, partial synthesis seeks to change only a subset of words in an utterance. The changes are performed in such a way that the target identity is maintained for the entire duration of the video or audio clip.

## 4 DEEPAKE GENERATION

In Figure 1, we first divide deepfake research by task, into deepfake generation methods and deepfake detection methods. We organize our literature review according to this split, discussing generative methods in the current section. We find that the most prominent approaches for deepfake generation are based on GANs or diffusion models. In some domains, such as video and audio, transformer-based methods are also very popular. Less frequently encountered methods are based on Variational Autoencoders (VAEs), Neural Radiance Fields (NeRF), 3D Morphable Models (3DMMs) and CNNs. Some models are only applied to

specific media types, *e.g.* NeRF and 3DMMs are typically applied in the video domain. A number of studies use hybrid models, combining GANs with VAEs on the one hand, or CNNs with RNNs on the other. We further structure our presentation according to the media type. For each media type, we divide the surveyed studies according to the underlying architectures. Moreover, we provide tutorials for the most important generative frameworks in the supplementary.

### 4.1 Image

#### 4.1.1 GAN-based methods

**Face synthesis.** Creating realistic faces is essential for deepfake generation, and GANs are widely employed to achieve this [5], [6], [54]–[60]. Shen *et al.* [54] introduce a third model into the traditional adversarial framework, tasked with determining whether the generated images retain the identity from a reference image. This approach enables the method to perform conditional generation. Conditional generation is also the focus of Xu *et al.* [60]. Their approach synthesizes high-quality 3D heads with control over the camera poses and other facial attributes. StyleGAN [56]–[58] improves the quality of the synthesized images by changing the generator architecture, and leveraging a mapping network to map the usual Gaussian vector to an intermediary latent space. The generative network adopts these new latents at different scales in the architecture via AdaIN [61] layers. Fu *et al.* [59] demonstrate that StyleGAN is also effective for generating full body images. Sauer *et al.* [5] extend the StyleGAN model, presenting a method that leverages Projected GAN training [62], progressive growing and classifier guidance [63], unlocking image synthesis at a resolution of  $1024 \times 1024$ . Different from StyleGAN, Esser *et al.* [6] utilize GANs to learn a perceptually rich codebook, representing images as a sequence of codebook entries. This approach enables the use of transformers for training.

**Face swapping.** One of the most widely-used methods for generating deepfakes is face swapping. This technique involves replacing the face in a target image with that of another individual, sourced from a different image. The key challenge lies in seamlessly integrating the source face, while maintaining non-identity-specific attributes, such as facial expressions and lighting conditions. Thanks to their well-known capacity of generating realistic images, GANs [7]–[10], [27], [55], [64]–[69] are widely adopted in face swapping frameworks.

In GAN-based face swapping pipelines, the generator is usually conditioned on identity information from the source image and the attributes extracted from the target image [9], [10], [65], [66], [68]–[73]. The work of Bao *et al.* [65] is one of the earliest contributions in this direction. The authors employ a face recognition model to extract an identity embedding from the source image. The attributes of the target image are extracted by a neural network trained to minimize the Euclidean distance between the target and generated images, applying a lower weight when the identities in the target and source images differ. The concatenated representations are processed by a generator that is trained in an adversarial setting. Li *et al.* [69] advance the previous framework by introducing a multi-level attribute encoder that is trained in a self-supervised manner,

providing a more detailed representation of the target image than the approach of Bao *et al.* [65]. Similarly, Chen *et al.* [66] propose the ID Injection Module, which integrates identity information through Adaptive Instance Normalization (AdaIN) [61] layers into the target image features. More recent works [9], [10], [68], [70], [71], [74] increase the quality and quantity of the conditional identity information. Kim *et al.* [9] enforce smoothness to the identity encoding space through contrastive learning. Cao *et al.* [74] and Cui *et al.* [72] explore the effectiveness of the transformer architecture for identity embedding in face swapping. Shiohara *et al.* [10] improve the embeddings by introducing BlendFace, a face encoder model trained with synthetic face images featuring swapped attributes. Rosberg *et al.* [68] leverage the feature maps provided by multiple layers of the face encoder to better represent the identity. Zeng *et al.* [70] show that a masked autoencoder (MAE) [75], pre-trained on a large-scale face dataset, is an effective encoder for face swapping. Yoo *et al.* [71] present the Triplet Adaptive Normalization block to integrate the pose and identity features within their generator.

Slightly distinct from earlier studies, another line of research [7], [8], [27], [67], [76]–[78] explores the manipulation of identity and attribute features within the latent space of the generator. Natsume *et al.* [27] create the latent space of the generator by merging the outputs of two encoders. One of the encoders is responsible for identity information, and the other for attributes. Two encoders are also utilized by Ren *et al.* [67] to separately learn embeddings for facial non-identity and non-facial attributes. This separation eliminates the need for skip connections, preventing identity leakage from the target image. Gao *et al.* [76] perform the disentanglement between identity and attributes through their novel Informative Identity Bottleneck layers that are included in a frozen face recognition model. Zhu *et al.* [77] train a model to perform GAN inversion and obtain the latent code for a given image. Subsequently, they employ another model to integrate the attributes from the target image into the latent code of the source face. The resulting latent code is fed into StyleGAN2 [57] to generate the swapped image. Similarly, Li *et al.* [7] employ learnable GAN inversion, but in their case, they leverage the latent space of a 3D GAN [79] to synthesize multi-view swapped images. The latent space of StyleGAN is further exploited for face swapping by Liu *et al.* [8], where the GAN inversion is extended at region level through the use of facial semantic masks. Jiang *et al.* [78] introduce identity preserving semantic bases (StyleIPSB) for StyleGAN [56]. Their approach identifies direction vectors in the latent space that modify attributes of generated images, while maintaining the identity of the subject.

GANs are also used in face swapping for the purpose of fixing the swapped image and making it more realistic [80]–[82]. Specifically, Moniz *et al.* [80] and Sun *et al.* [81] employ GANs to perform the blending of the source face in the target image, leveraging inpainting pipelines or frameworks such as CycleGAN [83]. Chen *et al.* [82] take a step further and present a method to correct deepfake images perturbed with adversarial attacks.

**Face editing.** Altering facial attributes such as age, gender, hair color or pose can be used to create counterfeit content. GANs support this kind of applications, either via image-

to-image translation between different domains [83]–[86] or via latent code manipulation [15], [87], [88]. CycleGAN [83] was first proposed for image-to-image translation between two domains. The method comprises two generators, one for each of the two domains. The primary contribution of Zhu *et al.* [83] is the introduction of the cycle-consistency loss, which ensures that the pipeline can reconstruct the original image, after translating it from one domain to the other and back. The main limitation of CycleGAN is its ability to handle only two domains. StarGAN [84], [85] addresses this limitation and supports image translation from multiple domains. The model allows this feature by including an additional condition as input, along with the conditional image. Hsu *et al.* [86] employ a dual-generator approach for image-to-image translation. The first generator produces a landmark image matching the pose of a reference image, and the second uses this landmark to recreate a source identity in the specified pose. Different from these approaches, other works [15], [87], [88] harness the latent space of StyleGAN. Tov *et al.* [87] study the latent space of StyleGAN and design an encoder for inversion, which is suitable for image editing. Similarly, Suwała *et al.* [88] design a plugin for the latent space of StyleGAN. This plugin disentangles the latent codes into attribute and non-attribute features, allowing attribute manipulation for facial editing. Bounareli *et al.* [15] project an identity image in the latent space of StyleGAN, and then, they harness pose and appearance encoders to create offsets, allowing the generator to change the pose of the identity latent.

#### 4.1.2 Diffusion-based methods

**Text-to-image.** Diffusion models are effectively applied in text-to-image generation [2]–[4], [89], utilizing large language models to encode textual descriptions that guide image creation. This capability allows users to generate counterfeit content featuring public figures simply by including their names in the text description used as input for generation. One of the most popular methods for text-to-image generation is Stable Diffusion [3], which leverages the latent space of a vector quantized (VQ) GAN [6] to perform the diffusion processes. SDXL [4] scales up the Stable Diffusion architecture, improving the quality and text fidelity of the generated images.

**Personalized generation.** Although text-to-image diffusion models allow deepfake content generation of public figures, some results do not accurately replicate the identity of the person. Thus, these models might have limited application in deepfake generation. However, there is another direction of research [11]–[13], [29], [30], [90]–[110] focused on generating images that contain a specific identity or concept depicted in an image or a set of images given as input. Such methods are more likely to be employed in deepfake generation. We can distinguish these contributions into two main approaches, those that perform test-time fine-tuning [90], [94], [97], [100], [102], [111]–[113] and those that leverage large-scale datasets and learn how to incorporate the additional images offline [11]–[13], [29], [30], [91]–[93], [95], [96], [98], [99], [102]–[110].

Test-time fine-tuning approaches use different components to integrate and learn the new identity. Some works introduce a new text token for the identity and learn to

embed it [97], [100], [111]. Other approaches [94] either use low-rank adaptation (LoRA) [114] or directly fine-tune the weights of the denoising network [90], [102], [112], [113]. Overall, test-time fine-tuning methods yield impressive results in terms of identity preservation, but their main disadvantage is the expensive optimization, which significantly increases the generation time. To this end, many works address the efficiency issue, to some extent. For example, Chen *et al.* [113] try to incorporate the knowledge of multiple subject-specific models into a single model. Ruiz *et al.* [112] leverage a HyperNetwork architecture to predict the network weights from a face image. Subsequently, they use these weights as a starting point for test-time fine-tuning, reaching faster convergence than previous work [90]. Despite these advancements, test-time fine-tuning methods still suffer from high generation times.

In contrast to test-time fine-tuning methods, the approaches that harness offline training [11]–[13], [29], [30], [91]–[93], [95], [96], [98], [99], [102]–[109] are faster in terms of generation time, but their primary issue is identity preservation. Therefore, solving the latter problem constitutes the priority of these works. Zhao *et al.* [12] propose an identity preservation loss for which they construct a better estimation of the original image given the predicted noise, at training time. The same idea is studied by Liu *et al.* [13], who improve the estimation even further. Peng *et al.* [92] employ an identity loss, but only for certain noise levels. Other methods [11], inspired by the GAN literature, employ the ArcFace model as identity embedding extractor for better identity representations. Similarly, Li *et al.* [109] improve representations by stacking multiple embeddings when multiple images are available. Lastly, Wang *et al.* [30] decouple the generation of background and identity-related content by training two separate denoising networks, out of which only one knows to generate images of a given person. **Tools.** The most popular tools for personalized generation are LoRA-based variants of Stable Diffusion [3] and SDXL [4], that are specialized on particular public personalities, *e.g.* Elon Musk<sup>3</sup> or Alan Turing<sup>4</sup>. Different from these options, another powerful tool is Midjourney<sup>5</sup>. In contrast to the LoRA-based methods, Midjourney is not popular for personalized generation, but for text-to-image synthesis. However, given the quality of its generative results, Midjourney is a popular tool for creating counterfeit images.

#### 4.1.3 Other methods

Unlike previous methods centered around generative models, some approaches rely on alternative techniques. Bitouk *et al.* [115] identify the closest match in terms of lighting and pose from a large set of face images, and perform face replacement using key point alignment. Korshunova *et al.* [116] use a multi-resolution CNN in the VGG feature space, aligning target and source images to minimize cosine distance between corresponding patches. Wang *et al.* [117] propose an encoder-decoder architecture with a 3D identity extractor and a Semantic Facial Fusion module to enhance resolution and preserve identity. In contrast, other works

combine generative methods to produce higher-quality images. Bao *et al.* [118] introduce CVAE-GAN, a method which combines VAEs with GANs. The generator and the encoder are trained with an adversarial objective, but also with a mean feature matching objective and a pixel-wise reconstruction loss, respectively. Li *et al.* [119] merge diffusion models and GANs by representing identity in Stable Diffusion through the latent space of StyleGAN, integrating latent codes into the U-Net via cross-attention layers.

## 4.2 Video

### 4.2.1 GAN-based methods

The early works for generating deepfake videos employ conventional GAN models for face swapping and reenactment, which are applied frame by frame [14], [120]. Nevertheless, these methods are usually part of more complex frameworks which have zero-shot capabilities, either involving more steps [14] or enhanced architectures [120]. Gao *et al.* [121] introduce a face reenactment GAN, focusing on generating videos of talking heads. The facial landmarks, expressions and head poses are extracted from both source and target frames to fit a face 3DMM and obtain predefined keypoints.

To depart from the conventional paradigm and improve the video generation using GANs, subsequent works leverage the latent space of StyleGAN2 [57]. A consistent number of methods divide the latent space in which they operate into two: one for content and one for motion [28], [122], [123]. While some utilize an RNN for sampling the motion trajectory [122], [123] and employ two discriminators, one for individual frames and one for the video sequence, Skorokhodov *et al.* [28] compute the motion embeddings with 1D convolutional layers and use only one video discriminator. Oorloff *et al.* [16] take a different approach by encoding both source and target frames, fusing their latent representations, then generating the output frame, while also utilizing multiple latent spaces [124], [125]. Yu *et al.* [126] treat videos as continuous-time signals and, with the help of an Implicit Neural Representation [127], they map an input signal (pixel coordinates and time) to RGB values in order to generate the corresponding video. Brook *et al.* [128] propose to generate multiple consecutive frames in low-resolution, and then increase their resolution with a super-resolution network. This approach ensures that training long video sequences is feasible.

### 4.2.2 Diffusion-based methods

Latent diffusion models [3] use a cross-attention mechanism that facilitates conditioning diffusion models for image generation. However, the main challenge in generating deepfake videos with diffusion models is employing a conditioning mechanism, while achieving temporal cohesion. The studies of Ho *et al.* [129] and Blattman *et al.* [130] represent the stepping stones in adopting diffusion models for video generation. Their methods extend diffusion models in several ways. The architectural changes applied on the U-Net mainly consist of replacing 2D convolutions with 3D convolutions, and appending additional self-attention layers for temporal attention. In a subsequent work, Blattmann *et al.* [131] demonstrate the benefits of using a large curated

3. <https://civitai.com/models/603798/elon-musk-sdxl>

4. <https://civitai.com/models/796450/alan-turing-mathematical-flux>

5. <https://www.midjourney.com/>

dataset for training a video generator. Wu *et al.* [132] introduce a one-shot method for editing a video given a text prompt. Inspired by Ho *et al.* [129], a text-to-image diffusion model is extended to an additional dimension (time) by Wu *et al.* [132], where the added self-attention layers operate on the current frame and the previous two frames.

Newer diffusion-based video generation methods [133]–[135] depart from the U-Net architecture and adopt a transformer-based one, namely ViT, which provides an innate mechanism for both spatial and temporal attention. This allows longer videos to be generated. Additionally, Guo *et al.* [134] introduce a plug-and-play module that can be integrated into a text-to-image diffusion model to induce the ability to generate videos. Inspired by this module, Wang *et al.* [136] present a method for text-to-video generation composed of several stages, in which a ControlNet is applied to improve guidance. Different from previous studies employing Stable Diffusion as the base model, Singer *et al.* [137] ground their work on DALLE-2 [138], while Ho *et al.* [139] utilize Imagen [89]. Nevertheless, similar architectural changes are implemented, where the network is extended to support the temporal dimension.

An important line of research is represented by portrait animation, in which a video is generated from a source frame and various conditional inputs. Most works in this area [140]–[143] aim to apply a sequence of facial expressions over the image. Two different approaches are used to condition the diffusion model. One is based on intermediate representations of the facial expressions, such as facial keypoints [140], [141], [143], and the other is based on directly encoding the frames containing the target facial movements [142].

Currently, the ability of the video generation methods based on diffusion modeling is not satisfactory, often requiring quality enhancements at a later stage in the pipeline. For example, super-resolution models are sometimes employed to increase video resolution [136], [137], [139], while the frame rate is usually increased through frame interpolation [131], [136], [137], [141].

#### 4.2.3 Other methods

Transformers represent the most popular architectural choice for video generation. For instance, Rochow *et al.* [144] leverage cross-attention blocks to guide the generation (using encoded facial keypoints and expressions), while Villegas *et al.* [145] apply the attention mechanism on frames to generate longer and coherent videos.

An alternative choice for video generation consists of employing some VQ autoencoder, either variational [146] or standard [147]. Similar to diffusion models, the generation process is carried out in the latent space of the autoencoder, which is lower dimensional. Within this vector space, a transformer is used to generate video tokens [147]–[150]. Unlike other related approaches, Jiang *et al.* [149] carefully design the latent space such that it is decomposed into an appearance and a pose representation, respectively.

A few methods harness the 3D space for face reenactment. In this context, warping is often employed, which involves computing a flow field between the source frame and the driving frame, then applying it on the former frame. In the same context, NeRF models [151] are used to

generate novel views of 3D face models. For example, Thies *et al.* [152] obtain a coarse 3D representation from the source frame using a traditional graphics pipeline, and then feed it to a neural network to obtain a neural texture, a high-dimensional embedding space, from which a Deferred Neural Renderer (based on U-Net) generates the target image. Thies *et al.* [153] and Yang *et al.* [154] synthesize faces by applying a deformation transfer between two 3DMM-based intermediate representations of the source and driving video frames, the mouth being further refined through warping. Zhang *et al.* [155] also apply warping based on dense landmarks, while Li *et al.* [156] combine warping with NeRF. Finally, to increase the performance, some works adopt pre-training strategies that involve masking the input and then reconstruct the signal [145], [147].

### 4.3 Audio

#### 4.3.1 GAN/VAE-based methods

A number of text-to-speech models employ popular generative frameworks, such as GANs and VAEs, either alone or in combination with more recent developments in the field. Kim *et al.* [157] propose an end-to-end TTS framework that augments variational inference with normalizing flows and uses an adversarial training procedure to enhance the representation potential. The method of Casanova *et al.* [158] constructs on the previous model, introducing new procedures, such as the concatenation of language embeddings with the ones of the input characters to allow training in a multilingual fashion.

In [159], the authors introduce new modules to develop NaturalSpeech, another VAE-based TTS. A differentiable durator is used to improve the duration prediction, a memory mechanism simplifies the waveform reconstruction, and a bidirectional prior/posterior module improves the prior from text, while simplifying the posterior from speech.

Lee *et al.* [160] present a GAN-based vocoder that improves the generator by introducing anti-aliased feature representation and periodic non-linearities, delivering state-of-the-art results and robustness for out-of-distribution scenarios, such as novel languages and speakers.

#### 4.3.2 Transformer-based methods

A few recent methods employ transformers to obtain competitive generation performance. FastSpeech [161] introduces a transformer-based model that speeds up speech synthesis by parallelizing Mel-spectrogram generation through a feed-forward architecture. A length regulator is used to match the length of the hidden states with the length of the Mel-spectrograms, and a duration predictor provides the duration for the phonemes. Jiang *et al.* [162] reuse the length regulator and the duration predictor from FastSpeech, adding separate modules for content, timbre and prosody modeling. A global timbre encoder is used to extract a global timbre vector, while a latent code language model fits the prosody distribution.

Wang *et al.* [163] proposed VALL-E, a framework that uses intermediate representations consisting of audio codec codes instead of Mel-spectrograms. A pre-trained neural codec model generates acoustic codes that are used alongside corresponding phoneme sequences during training, allowing the neural language model to extract both speaker



information and content. SPEAR-TTS [164] removes the necessity to supply the transcripts of audio prompts by decoupling the generation of semantic tokens and the acoustic tokens.

Yang *et al.* [165] introduce UniAudio, a hierarchical transformer framework that learns both inter-frame and intra-frame correlations separately, reducing the computational complexity. It supports the generation of multiple types of audio by employing LLM-style next token prediction and tokenization via a universal neural codec.

#### 4.3.3 Diffusion-based methods

Following the success of diffusion models in vision [1], several generation methods adopted the diffusion modeling framework to generate deepfake audio. Huang *et al.* [166] present FastDiff-TTS, a conditional diffusion model that follows the architectural design proposed by Ren *et al.* [167]. The authors employ time-aware location variable convolutions for long-term dependency modeling and a noise schedule predictor for sampling acceleration. ProDiff [168] is another framework with an architecture inspired by Ren *et al.* [167], which uses a denoising model with a parametrization that directly predicts the clean data, halving the number of diffusion steps via knowledge distillation.

The audio encoder/decoder, the phoneme encoder and the duration and pitch predictors proposed by Tan *et al.* [159] are reused in NaturalSpeech 2 [169], alongside a diffusion model that learns to predict latent representations conditioned on the input text. To promote zero-shot generation, a speech prompting mechanism helps the diffusion model and the duration and pitch predictors to follow prosody, style and speaker identity from the supplied audio prompt. The encoder/decoder and the duration predictor are further used in NaturalSpeech 3 [170], where, in contrast to previous studies [159], [169], each of the following speech attributes are independently generated by a novel factorized diffusion model: duration, content, prosody and acoustic details.

Du *et al.* [171] introduce UniCATS, a framework capable of performing speech editing tasks, where speech is synthesized by taking into account both preceding and following contexts. UniCATS can achieve this with a contextual VQ-diffusion-based acoustic model and a contextual vocoder.

Yang *et al.* [172] aim to solve learning problems specific to expressive TTS, a novel task that tries to control the speaking style of the synthesized speech. The proposed framework uses RoBERTa [173] to extract the style representation from a natural language prompt.

## 4.4 Multimodal

### 4.4.1 Transformer-based methods

Recent advancements in talking face generation focus on improving the synchronization of facial movements with speech, while maintaining natural motion and emotional consistency [174]–[177]. These approaches address challenges such as lip-sync accuracy [175]–[177], motion stability [174], [177] and speaker-specific styles [174], [175], aiming for realistic human-video synthesis. Jang *et al.* [174] introduce a system that combines talking face generation with TTS, addressing the challenge of generating natural head poses and maintaining consistent speech patterns even with

varying facial motions. Their approach leverages a motion sampler and a conditioning method to ensure fluidity in both aspects. Building on the idea of synchronizing audio and visual elements, Cheng *et al.* [175] propose VideoReTalking, a method designed to edit real-world talking head videos for perfect lip-sync and emotional consistency. In a similar fashion, Wang *et al.* [176] develop a one-shot talking face generation framework. They introduce an audio-visual correlation transformer, which improves lip-sync accuracy by mapping audio to dense motion fields through phoneme and keypoint representations. Addressing another challenge in the field, Ling *et al.* [177] tackle the issue of lip motion jitter in speech-driven talking face generation. Their solution, StableFace, identifies key problems such as noise in the 3D face representation and mismatches between training and inference stages. To address emotion-agnostic talking head generation, Gan *et al.* [178] propose emotional adaptation for audio-driven talking-head. The method enhances emotion-agnostic talking-head models by adding three lightweight adaptations: deep emotional prompts, an emotional deformation network, and an emotional adaptation module.

### 4.4.2 Diffusion-based methods

Several frameworks are designed to generate high-quality audio-driven portrait animations, aiming to achieve realism and synchronization [18], [31], [179]–[183]. AniPortrait [179] transforms audio into photorealistic animations by extracting 3D facial meshes and head poses, allowing for flexible facial motion editing. Building on this, V-Express [180] focuses on precisely synchronizing lip movements with audio, while maintaining control over facial identity and background through progressive training techniques. EchoMimic [181] offers another solution by integrating audio and facial landmarks using a denoising U-Net architecture, which stabilizes and enhances the natural flow of portrait videos. Similarly, Stypułkowski *et al.* [31] propose an autoregressive diffusion model to achieve realistic talking heads with smooth, expressive movements and accurate lip-sync. Xu *et al.* [18] also use a diffusion-based framework to improve lip-sync accuracy and motion diversity, employing a hierarchical audio-driven visual synthesis module. In a similar direction, VASA [182] produces talking faces, capturing synchronized lip movements and dynamic expressions using a diffusion-based model in a latent facial space, enabling real-time interactions with high realism. Distinctly, EMO [183] generates expressive talking head videos without relying on 3D models, excelling in natural transitions and seamless identity preservation.

### 4.4.3 Other methods

Recent advancements in talking head generation leverage different models, such as NeRF [184], [185], GANs [186]–[188], RNNs [189]–[191], VAEs [192] or CNNs [193], to address the challenges of synchronization, realism, and efficiency. SyncTalk [184] is a NeRF-based approach which focuses on speech-driven video generation. SyncTalk enhances synchronization between lip movements, facial expressions and head poses by using a face-sync controller for precise lip-sync, a head-sync stabilizer for natural head movements, and a portrait-sync generator to integrate the generated head with the torso. Similarly, GeneFace++ [185] builds on NeRF technology to produce real-time talking face

videos with arbitrary speech audio. By improving audio-lip synchronization using pitch contour analysis and incorporating a fast motion-to-video renderer, GeneFace++ offers a robust and efficient solution.

In the context of GANs, Text2Video [186] presents an approach for synthesizing videos directly from text, reducing reliance on audio-driven models. Using a phoneme-pose dictionary and a GAN-based architecture, the method achieves high-quality video synthesis with just one minute of training data. HeadGAN [187] is developed for head reenactment and editing from a single reference image. It integrates 3DMMs for real-time reenactment at approximately 20 FPS. Furthermore, it incorporates audio features for enhanced mouth movement accuracy. Wang *et al.* [188] introduce TalkLip, a speech-to-lip generation model that enhances lip-speech intelligibility by incorporating a lip-reading expert to penalize incorrect outputs.

For gesture generation, RNN-based frameworks such as hierarchical audio-to-gesture (HA2G) [189] introduce ways to generate co-speech gestures. HA2G extracts multi-level audio features using a hierarchical audio learner. Another RNN-based model [190] offers a real-time pipeline to generate personalized photorealistic talking-head animations. This model operates at over 30 FPS and follows a three-stage process: extracting deep audio features, predicting facial dynamics and head motions with an auto-regressive model, and rendering high-fidelity faces through image-to-image translation. Some RNN-based methods [194], [195] rely on audio-visual cues for realistic face synthesis. Peng *et al.* [194] present a speech-driven 3D face animation model that separates speech content and emotion using an Emotion Disentangling Encoder, while Zhong *et al.* [195] introduce a two-stage framework for audio-driven person-generic talking face video generation. Both approaches apply RNNs on top of CNN features. The CNN-based DisCoHead [193] offers an unsupervised approach to disentangling head motion from facial expressions. By applying geometric transformations to isolate head motion and using speech audio for facial expressions, DisCoHead efficiently generates realistic talking heads.

## 5 DEEPPAKE DETECTION

The second part of our taxonomy illustrated in Figure 1 comprises deepfake detection methods. The taxonomy clearly indicates that most deepfake detectors are based on CNN architectures. However, with the recent advent of vision and audio transformers, a large body of work on deepfake detection is now based on multi-head attention. To boost detection performance, a considerable number of studies employ hybrid models, combining CNNs with transformers or RNNs, respectively. Less prevalent architectures in deepfake detection are graph and recurrent neural networks. We organize our subsequent presentation of deepfake detection methods according to the input domain. The reviewed articles are further separated according to the employed architectures.

### 5.1 Image

#### 5.1.1 CNN-based methods

Convolutional nets are the most prevalent type of architecture for deepfake detection [19]–[22], [196]–[223]. The detec-

tion task is commonly formalized as a binary classification, where a CNN backbone is used as feature extractor. The most frequently chosen backbones in this line of work are EfficientNet [224], XceptionNet [225] and ResNet [226].

The main direction of research for deepfake detection focuses on developing methods that generalize well across different types of manipulations [20], [21], [196]–[206], [208], [223], [227]. To improve generalization, Chen *et al.* [20] propose an adversarial training pipeline that dynamically identifies which type of forgeries are most challenging for the deepfake detector, and uses them to improve the overall performance. Similarly, Yan *et al.* [196] increase the diversity and complexity of forgeries, but different from Chen *et al.* [20], they achieve this through latent space manipulations. Other studies [21], [198], [200], [202], [203], [206], [208] try to identify the common artifacts or features for different types of forgeries. Thus, some methods [21], [198], [206], [208] base their solution on local artifact detection and less on the overall identity. Yan *et al.* [200] propose an explicit disentanglement approach using multi-task learning to analyze image information, allowing the detection of features that are shared across various types of forgeries. Some studies [202], [203], [223], [228], [229] indicate that the generalization capabilities of deepfake detectors can also be improved by leveraging artifacts in the frequency domain. All these advancements in the generalization of deepfake detectors are validated by Yao *et al.* [199], who demonstrate that detectors modeling low-order interactions exhibit superior generalization capabilities. In contrast to previous works, some deepfake detection methods [19], [207], [214], [216] are specialized in identifying specific forgery techniques. Huang *et al.* [207] and Shiohara *et al.* [19] focus on face swapping. Huang *et al.* [207] argue that a face-swapped image contains information about the identity present in the target image. Their method builds a face recognition model to detect this identity from a face-swapped image, and leverages the difference between its embedding and the embedding of the source identity to detect deepfakes. Shiohara *et al.* [19] improve face-swapping detection by generating more challenging examples for a face-swapping detector. In this regard, they use face-swapping pipelines where the target and source images are of the same identity and closely resemble each other in terms of face position and other non-identity attributes. In contrast, Hooda *et al.* [214] and Tantar *et al.* [216] tackle the detection of forgeries in images generated with diffusion models. Hooda *et al.* [214] detect deepfakes using an ensemble in which the models are using disjoint parts of the input features, aligning with the aforementioned finding of Yao *et al.* [199].

Other studies explore deepfake detection methods by examining their fairness [22], [215], level of explainability [209], [217], or level of vulnerability to adversarial attacks [230]. Ju *et al.* [215] propose the first approach to tackle fairness in deepfake detectors. Their method uses groups of people specified by the user and ensures that the loss of these groups is similar to each other. The method of Ju *et al.* [215] performs well when tested on the same type of forgery as in the training set. Lin *et al.* [22] extend the work of Ju *et al.* [215], aiming to improve generalization across different forgery types. More specifically, their work introduces a disentanglement loss to separate the demo-

graphic and forgery specific features. These features are then combined and used in a fairness loss function, which aims to ensure equal importance across different demographic groups. In terms of explainability, Dong *et al.* [209] try to identify the visual concepts that are relevant for deepfake detectors. Their findings indicate that the features specific to source and target images are, in general, ignored. The focus of the detection models is on visual artifacts. Trinh *et al.* [217] reinforce this finding by showing that, in addition to visual artifacts, temporal artifacts also serve as evidence for detectors. Hussain *et al.* [230] show that, despite the progress of deepfake detectors, these models are susceptible to adversarial attacks and future works need to address this drawback.

### 5.1.2 GCN-based methods

As stated before, Yao *et al.* [199] demonstrate that deepfake detectors with strong generalization capabilities tend to model low-order interactions. Due to their ability to capture such relationships, some studies employ graph convolutional networks (GCNs) for deepfake detection [231]–[233]. Yang *et al.* [231] design their graphs with vertices representing features of facial regions and edges capturing the correlations between these regions. They also introduce a masking strategy that removes edges based on their weights. These graphs are then given as input to a GCN, which extracts features for a binary classifier. Wang *et al.* [232] propose a similar method, but along with the spatial domain features, they also include frequency domain information.

### 5.1.3 Transformer-based methods

State-of-the-art results in a broad range of computer vision tasks are achieved by transformer architectures [234]. As a result, these architectures are often adopted for deepfake detection [32], [33], [235]. Aghasanli *et al.* [235] propose a direct application of transformers for deepfake detection, where the model is used as feature extractor for a binary classifier. The method of Dong *et al.* [33] is similar to that of Huang *et al.* [207], as both approaches harness discrepancies between the explicit identity depicted in the image and the one given by the outer region of the face. However, in contrast to Huang *et al.* [207], Dong *et al.* [33] use the same model for both identities and differentiate between the two regions with two additional tokens. Wang *et al.* [32] use a multi-scale transformer which operates on patches of different sizes to extract features for deepfake detection.

### 5.1.4 Hybrid architectures

A natural strategy for improving deepfake detection is to combine the aforementioned architectures in a joint pipeline. The primary research focus is on combining transformer and CNN architectures [37], [236], [237], though the combination of GANs and CNNs [238] is also explored. For instance, Jeong *et al.* [238] train a GAN to generate perturbation maps, which are added to both real and fake images to minimize differences at the frequency level. They argue that using this augmented data to train a CNN-based classifier prevents overfitting to method-specific frequency artifacts, thereby improving the generalization.

Kamat *et al.* [237] focus on exploring different techniques of combining CNN-based and transformer-based feature extractors for deepfake detection. Shao *et al.* [37] and Hong *et*

*al.* [236] harness the CNN and transformer combination for a slightly different problem. Their goal is to determine the sequence of facial manipulations used to create a fake image, because deepfakes are frequently created with several manipulation steps. Both methods [37], [236] rely on a CNN backbone for feature extraction followed by a transformer that returns the sequence of manipulated regions.

### 5.1.5 Traditional machine learning methods

The earliest approaches to deepfake detection examine the effectiveness of classical machine learning algorithms [239], [240]. Guarnera *et al.* [240] conjecture that transposed convolutional layers within GANs generate local pixel correlations. Leveraging this, they design an algorithm to extract local features from images, which are then passed to machine learning algorithms such as SVM and k-NN for detection. Lugstein *et al.* [239] also employ an SVM, but focus on extracting features from the photo response non-uniformity (PRNU) signal.

## 5.2 Video

### 5.2.1 CNN-based methods

Most deepfake video detection methods are based on plain convolutional models. In the majority of works, 3D convolutions are applied to extract spatio-temporal features from a whole video sequence. Nevertheless, 2D convolutions are also used to extract salient spatial features from individual frames, and then, the results are aggregated to make the final prediction. Agarwal *et al.* [241] extract facial features (both static and temporal), and then compare them with a reference set, for each biometric source data, to identify a similar data point, whose label is used for prediction. Similarly, Cozzolino *et al.* [242] compare the extracted biometric features from the input video to those of a pristine video.

Some works propose to capture temporal inconsistencies in the video. This is either achieved from successive frames, with some specialized sequential convolutional blocks [243], or through a hierarchical framework from both local (frame) and global (video snippet) perspectives that can differentiate between real and fake videos [244], [245]. Based on this objective, some papers focus only on inconsistencies of specific aspects of the face. Haliassos *et al.* [246], [247] study mouth movements and propose to learn spatio-temporal representations of mouth motion via two-stage frameworks. Demir *et al.* [248] magnify the motion of the face and then classify the videos, while also identifying the source generation method.

### 5.2.2 Hybrid CNN and RNN architectures

To obtain a prediction based on the temporal dimension, many detection methods employ a recurrent network to aggregate the latent features extracted by CNNs. Multiple variants of recurrent architectures are used, such as simple RNNs [38], [249], gated recurrent units (GRUs) [39], [250], [251] and Long Short-Term Memory (LSTM) networks [252], [253]. Unlike the rest, Masi *et al.* [253] use two branches, each with a different specialization in extracting features, one in the frequency domain and one in the color domain. The resulting representations are aggregated in a bi-directional LSTM, which is optimized via an improved loss function. Montserrat *et al.* [251] and Liu *et al.* [39] extract face crops

and employ ArcFace [254] for a better representation of the backbone features.

### 5.2.3 Other hybrid architectures

Aside from combining CNNs and RNNs, some attempts try to fuse other types of neural networks. An important category is represented by the integration of attention [255]–[257]. Bonettini *et al.* [255] integrate an attention mechanism into each network in an ensemble of CNNs. Wang *et al.* [256] extract noise features from the face crop, as well as a background crop, and feed them into a multi-head attention module. Furthermore, these works adopt the contrastive learning paradigm in their deepfake video detectors. Different from previous methods, Coccomini *et al.* [258] combine various ViTs with an EfficientNet [224], the latter being used for feature extraction.

Due to its demonstrated strength in many tasks, the transformer architecture [259] is often employed to capture temporal incoherence. For instance, in a number of studies, the transformer is used together with 3D CNNs to extract temporal features [35], [36], [260]. Choi *et al.* [36] propose a complex framework that utilizes latent features from a pre-trained StyleGAN model [261], which are further encoded with GRUs. Cai *et al.* [34] employ the masked autoencoder pre-training framework to learn facial representations by guiding the masking strategy to focus on hiding face information. The encoder is further fine-tuned on deepfake detection.

Another interesting direction is to formulate the problem as a graph classification task. Tan *et al.* [262] extract embeddings associated with the actions of different facial elements, systematically arrange them in a graph, and then employ a GCN to classify the video. Xu *et al.* [263] propose a novel strategy: to randomly sample frames from a video and combine them into a single image, called thumbnail. Then, the thumbnail is processed by a Swin Transformer [264] to obtain feature embeddings. Finally, these are fed into a GCN to capture any inconsistency and thus identify fake videos.

## 5.3 Audio

### 5.3.1 CNN-based methods

The ability of CNNs to extract local features allows them to achieve competitive results in spoofed audio detection. Tak *et al.* [265] bring small modifications, such as fixed sinc filters, to the RawNet2 architecture and use it for spoofed speech detection. The same base architecture is further improved by Wang *et al.* [266] with orthogonal convolutions and temporal convolution networks (TCNs) to enhance the discrimination capability. In contrast, Conti *et al.* [267] introduce a new pipeline architecture that uses a Speech Emotion Recognition (SER) system as the feature extractor, and a Random Forest as the final classifier. Emotion features are extracted from an intermediate layer of a 3D-Convolutional Recurrent Neural Network. Hua *et al.* [268] introduce the Time-domain Synthetic Speech Detection Net (TSSDNet), an end-to-end framework that considers Inception-style convolutions and ResNet-like skip connections. The resulting architectures obtain state-of-the-art results on the ASVspoof 2019 dataset.

### 5.3.2 GNN-based methods

Some models use the ability of GNNs to model relationships between entities in order to enhance spoofed speech detection. Graph attention networks (GATs) are used in [269]–[271] to detect artifacts from both temporal and spectral domains. Tak *et al.* [269] use two separate GATs for relationship modeling between neighboring temporal segments and different sub-bands, respectively, fusing the scores for the final prediction. In a different study [270], the authors use a third GAT to integrate information from temporal and spectral sub-graphs, while Jung *et al.* [271] propose to combine the two sub-graphs into a single heterogeneous graph via a heterogeneous attention mechanism. Chen *et al.* [272] use a GCN to model the relationships from a graph constructed from patches of a spectrogram, outperforming competing models.

### 5.3.3 Transformer-based methods

A growing number of methods use transformers for the synthesized speech detection task. Bartusiak *et al.* [273] employ a compact convolutional transformer (CCT) to extract feature maps with a convolutional block from supplied spectrograms and further analyze them, after concatenation, using an attention mechanism. The CCT is extended in [274] to produce a compact attribution transformer (CAT) for the speech synthesizer attribution task, which aims to identify the tool/method that was used to synthesize the speech input. The proposed method also uses spectrograms as input, further producing a probability distribution over the set of known synthesizers.

Martín-Doñas *et al.* [275] present a model trained in a self-supervised manner, employing representations from different transformer layers of a pre-trained wav2vec 2.0 model to detect spoofed speech. The intermediate representations from these layers are used to construct a vector for each time step. wav2vec 2.0 is also used by Cai *et al.* [276], who address partially fake audio detection. They identify fake audio segments by discovering the discontinuity between them. Features are extracted with the aforementioned model, and frame embeddings are obtained by a ResNet-1D module, while transformer-based encoders capture the context of the frames with respect to the sequence.

Zhang *et al.* [277] aggregate a transformer architecture and a residual network, where the ability of the transformer to model long-term dependencies allows it to find correlations between audio frames. Different data augmentation techniques are used to increase the size of the training set and the final score is computed with a logistic regression meta-classifier that takes the scores from multiple models. Rawformer [278] aims to improve AASIST [271] by replacing the GAT with a transformer encoder. A positional aggregator augments the feature maps obtained by the RawNet2 feature extractor with positional information.

### 5.3.4 Other methods

A method based on monitoring the behavior of neurons from a speaker recognition (SR) model is designed by Wang *et al.* [279]. The activated neurons from convolutional and fully-connected layers are used as feature vectors in the training process of a shallow network that classifies the input speech as genuine or synthesized. Zhang *et al.* [280]

introduce Radian Weight Modification (RWM), a continual learning method that adjusts the direction of the gradient based on the means of the intra-class cosine distances of the samples from the current batch.

## 5.4 Multimodal

### 5.4.1 CNN-based methods

Recent advances in deepfake and multimedia manipulation detection focus on combining audio-visual elements to improve model robustness and accuracy [25], [281], [282]. In Multimodaltrace [281], a ResNet-based framework blends audio and visual features, both independently and jointly, to enhance deepfake detection, offering insights into model focus areas using integrated gradient analysis. Similarly, Cozzolino *et al.* [25] propose a Person-of-Interest detector which leverages unique identity markers through contrastive learning with ResNet-50, excelling at detecting inconsistencies in low-quality videos. Kihal *et al.* [282] introduce VTA-CNN-RF, a deep multimodal spam detection system, achieving over 98% precision in text, image, and audio spam classification using CNNs and Random Forests.

### 5.4.2 Transformer-based methods

Recent deepfake detection frameworks based on attention mechanisms exploit both audio and visual cues to effectively identify manipulated content [24], [283]–[291]. Zhou *et al.* [24] propose a joint audio-visual detection method that leverages the synchronization between modalities, significantly boosting detection accuracy by late-fusing joint predictions with inter-attention mechanisms. Building on this approach, Oorloff *et al.* [283] develop a two-stage audio-visual feature fusion method, using contrastive learning and autoencoders in the initial phase to capture audio-visual correspondences, followed by fine-tuning of transformer-based encoders for precise deepfake classification. Additional frameworks that employ audio-visual cues have been proposed. For example, Salvi *et al.* [284] introduce a framework analyzing audio-visual feature discrepancies over time, uniquely trained on separate monomodal datasets to identify unseen deepfakes. Similarly, AVFakeNet [285] is a unified model with dense Swin Transformer modules, which aptly handles variations in facial poses, lighting, and demographic diversity. Asha *et al.* [286] propose an ensemble-based D-Fence model, utilizing cross-modal attention and self-attenuated neural networks to emphasize correlations between visual and audio elements for improved detection accuracy.

For both intra and inter modality deepfake detection, Liu *et al.* [287] introduce the Forgery Clues Magnification Transformer (FCMT), which amplifies both intra-modal and cross-modal forgery cues through a distribution difference-based inconsistency computing module. Feng *et al.* [288] tackle audio-visual inconsistencies through an anomaly detection method that trains autoregressive transformers to flag low-probability sequences, using a joint ResNet-18 and VGG-M encoder. Zou *et al.* [289] advance cross-modality and within-modality regularization by aligning distinct audio and visual signals through multimodal transformers, while Nie *et al.* [290] introduce FRADE, which relies on adaptive forgery-aware injection and audio-distilled cross-modal interaction to effectively bridge the audio-visual domain gap.

TABLE 2

Datasets that are commonly used in deepfake detection literature, separated by domain. AV stands for audio-video (multimodal).

Dataset	#Real	#Fake	Resolution/frequency	
Image	DFFD [292]	58,703	240,336	250×250 - 1024×1024
	FakeSpotter [293]	6,000	6,000	224×224
	ForgeryNet [294]	1,438,201	1,457,861	240×240 - 1080×1080
	DiffusionFace [295]	30,000	600,000	256×256
Video	FaceForensics++ [296]	1,000	4,000	512×512
	DeeperForensics [297]	48,475	11,000	1920×1080
	Celeb-DF [298]	590	5,639	256×256
	WildDeepfake [299]	3,805	3,509	varying
	DeepFake-TIMIT [300]	0	620	64×64/128×128
	UADFV [301]	98	98	294×500
GenVideo [302]	1,224,511	1,089,671	224×224 - 1280×2048	
Audio	WaveFake [303]	0	117,985	16 kHz
	ASVspooF 2019-LA [304]	12,483	108,978	16 kHz
	ASVspooF 2021-LA [305]	16,492	148,148	16 kHz
	ASVspooF 2021-DF [305]	20,637	572,616	16 kHz
	In-the-Wild [306]	19,963	11,816	16 kHz
	ADD 2022 [307]	91,464	358,082	16 kHz
	ADD 2023 [308]	243,194	273,874	16 kHz
	FoR [309]	111,000	87,285	16 kHz
	MLAAD [310]	0	154,000	22 kHz
AV	FakeAVCeleb [311]	500	19,500	224×224
	LAV-DF [312]	36,431	99,873	224×224
	DFDC [313]	23,654	104,500	1920×1080/1080×1920

Moreover, Yang *et al.* [314] introduce AVoid-DF, a model based on a temporal-spatial encoder and a multimodal joint decoder. AVoid-DF captures inter-modal and intra-modal disharmony, achieving good performance across various forgery techniques.

### 5.4.3 Other methods

Hosler *et al.* [315] introduce a method for detecting deepfakes by analyzing emotional consistency in human faces and voices using LSTM networks. By predicting emotions from audio and video features, the approach identifies unnatural emotional patterns to flag deepfake media.

## 6 DATASETS AND REPORTED RESULTS

In Table 2, we present the most frequently used datasets for deepfake detection, along with the number of real and fake samples, as well as the resolution (for visual datasets) or the bit rate (for audio datasets). We next describe the main steps that are usually employed to build deepfake datasets. The first step of creating a dataset for deepfake detection is collecting the real data. Except for Dolhansky *et al.* [313], who create the original data by recording movies of paid actors, the basic procedure is to scrape the Internet for videos, especially YouTube. Even image datasets use frames extracted from videos. After acquiring real data, various deepfake methods are applied to generate the fake samples. Due to their excellent trade-off between performance and speed, GANs are adopted for the creation of most datasets, *e.g.* Celeb-DF [298], DeepFake-TIMIT [300], DFFD [292], FakeSpotter [293] and FakeAVCeleb [311]. VAEs represent the method of choice only for a few datasets, *e.g.* DeeperForensics [297] and ASVspooF 2019-LA [304]. Diffusion models are recent and powerful generative methods, yet they require more computation. Hence, only a couple of recent datasets employ them to create the fake samples, *e.g.* GenVideo [302] and DiffusionFace [295]. Several datasets, especially the audio ones, are generated by using more than one



TABLE 3  
Results of top scoring image deepfake detection methods on DFFD [292], DiffusionFace [295], ForgeryNet [294].

Dataset	Method	Accuracy	AUC
DFFD [292]	BNext-M [221]	99.18%	0.9994
	VGG-16 [292]	-	0.9967
	XceptionNet [292]	-	0.9964
DiffusionFace [295]	GramNet [227]	62.60%	0.7250
	GFF [228]	61.10%	0.7250
	RECCE [233]	64.40%	0.7130
	F3Net [229]	59.80%	0.6960
ForgeryNet [294]	SNRFilters-Xception [240]	81.09%	0.9052
	GramNet [227]	80.89%	0.9020
	F3Net [229]	80.86%	0.9015
	XceptionNet [296]	80.78%	0.9012

TABLE 4  
Results of top scoring video deepfake detection methods on FaceForensics++ [296], DFDC [313] and Celeb-DF [298] datasets.

Dataset	Method	Accuracy	AUC
FaceForensics++ [296]	TALL++ [263]	98.65%	0.9987
	LipForensics [246]	98.90%	0.9970
	FADE [262]	92.89%	0.9952
	M2TR [32]	97.93%	0.9951
	App.+Beh. [241]	98.90%	0.9900
	RealForensics [247]	-	0.9900
DFDC [313]	Efficient ViT [258]	-	0.9510
	TALL++ [263]	-	0.9068
	CNN Ensemble [255]	-	0.8782
	RealForensics [247]	-	0.7590
	LipForensics [246]	-	0.7350
Celeb-DF [298]	App.+Beh. [241]	98.50%	0.9900
	FInfer [250]	90.47%	0.9330
	RealForensics [247]	-	0.8690
	LipForensics [246]	-	0.8240

method, *e.g.* GenVideo [302], ForgeryNet [294], FaceForensics++ [296], DFDC [313], LAV-DF [312], WaveFake [303], ASVspoo 2019-LA [304], ASVspoo 2021-LA/DF [305], FoR [309] and MLAAD [310]. A few visual datasets [296], [302] rely on online tools to create deepfakes, the most popular tool being FaceSwap<sup>6</sup>.

Depending on the input modality, different metrics are commonly reported. For the visual modalities, the most frequent metric is the area under the curve (AUC). Given that deepfake detection is a binary classification task, the AUC score can illustrate the ability of the model to differentiate between real and fake samples. The AUC is obtained by plotting the True Positive Rate against the False Positive Rate for multiple thresholds, and then computing the area under the resulting curve. Accuracy is an alternative metric that can be used to assess the overall performance of a deepfake detection model. Nevertheless, deepfake detection datasets are usually imbalanced, making accuracy a less preferred option. For the audio modality, models are regularly evaluated via the equal error rate (EER). Its popularity is given by the robustness to class imbalance, while equally assessing false positives and false negatives. EER is computed by finding the intersection of the False Acceptance Rate and the False Rejection Rate.

## 6.1 Results on Popular Benchmarks

**Image.** Table 3 includes top accuracy and AUC scores on three commonly-used datasets of deepfake images, namely DFFD [292], DiffusionFace [295] and ForgeryNet [294]. The

TABLE 5  
Results of top scoring audio deepfake detection methods on the ASVspoo 2019-LA [304] and ASVspoo 2021-LA [305] datasets.

Dataset	Method	EER (%)
ASVspoo 2019-LA [304]	GCN [272]	0.58
	Rawformer [278]	0.59
	AASIST [271]	0.83
	RawGAT [270]	1.06
	TO-RawNet [266]	1.58
	TSSDNet [268]	1.64
ASVspoo 2021-LA [305]	wav2vec2+AASIST [316]	0.82
	wav2vec2+MLP [275]	3.54
	TO-RawNet [266]	3.70
	Rawformer [278]	4.53
	AASIST [271]	9.15

TABLE 6  
Results of top scoring multimodal deepfake detection methods on the FakeAVCeleb [311] and DFDC [313] datasets.

Dataset	Method	Accuracy	AUC
FakeAVCeleb [311]	FRADE [290]	98.60%	0.9980
	AVFF [283]	98.60%	0.9910
	MIS-AVoiDD [317]	96.20%	0.9730
	PVASS-MDD [318]	95.70%	0.9730
	SSVF [288]	94.20%	0.9450
	MRDF [289]	94.05%	0.9243
DFDC [313]	FRADE [290]	97.20%	0.9900
	PVASS-MDD [318]	96.30%	0.9890
	AVoiD-DF [314]	91.40%	0.9480
	AVA-CL [291]	84.20%	0.8864
	AVFakeNet [285]	82.80%	0.8620

results suggest that the oldest dataset, DFFD, has become saturated due to advancements in recent deepfake detectors. In contrast, the most recent dataset, DiffusionFace, featuring faces generated by diffusion models, poses a significantly greater challenge for state-of-the-art detectors. This highlights the need for future developments in deepfake detection to effectively differentiate between genuine faces and those synthesized by diffusion models.

**Video.** Table 4 provides the performance levels of a few of the most effective methods for video deepfake detection on three distinctive datasets: FaceForensics++ [296], DeepFake Detection Challenge (DFDC) [313] and Celeb-DF [298]. The main metric in this area is the AUC, but we also report the accuracy, whenever it is available. All the included models are outstanding, most of them almost achieving flawless performance. This is not only true for the newest methods that employ the most recent trends (such as transformers), but also for the preceding ones, that solely utilize CNNs. Nevertheless, on the more difficult datasets (DFDC and Celeb-DF), it can be observed that ViT-based architectures are superior.

**Audio.** In Table 5, we report the Equal Error Rate (EER) values of top audio deepfake detection methods on ASVspoo 2019-LA [304] and ASVspoo 2021-LA [305], two of the most popular audio deepfake detection datasets. GNN-based methods achieve the lowest EER values on both datasets, demonstrating their effectiveness in the synthesized speech detection task.

**Multimodal.** In Table 6, we present the performance levels of top deepfake detection methods on two popular multimodal datasets, namely FakeAVCeleb [311] and DFDC [313]. For each method, we report the performance in terms of accuracy and AUC. For the FakeAVCeleb dataset, we include results from six state-of-the-art detection methods,

6. <https://github.com/deepfakes/faceswap>



Fig. 3. Randomly sampled frames captured from the fake videos included in BioDeepAV exhibit a high level of realism. Best viewed in color.

with FRADE [290] and AVFF [283] achieving the highest accuracy rates, both at 98.60%. On the DFDC dataset, we include five methods, with FRADE and PVASS-MDD demonstrating top performance, attaining accuracy rates of 97.20% and 96.30%, respectively. The reported results hint towards important advancements in multimodal methods, with recent methods nearly saturating the benchmarks.

## 6.2 Proposed Benchmark

We create a new dataset, called BioDeepAV<sup>7</sup>, to assess the out-of-domain generalization capabilities of deepfake detection models. Our primary focus is on generating videos featuring talking faces, but we also include audio-video examples with audio-only manipulations. Figure 3 depicts a few frames from various deepfake videos in BioDeepAV.

**Generated Data.** We generate over 1,600 deepfake videos using four recent methods specialized in talking-face synthesis [181], [184], [319], [320]. These approaches base their solutions on the recent development of diffusion models [181], [319], NeRFs [184] and Gaussian Splatting [320]. We use three face image sources to sample target identities. First, we create 300 synthetic faces using RealVisXLv5<sup>8</sup> and supplement these with faces from the LAION-Face [321] and HDTF [322] datasets. Three of the methods [181], [184], [320] also require head motion information as a conditioning signal, which we obtain from the HDTF [322] dataset. In addition to face images and motion cues, all methods also require an audio file to condition their talking-face generation. For this, we use the samples from a dataset of English dialects [323], the audio from the HDTF [322] dataset, and over 700 deepfake audio samples created by us. To generate these synthetic audio samples, we employ StyleTTS [324], SSR-Speech [325] and YourTTS [158], which support both text-to-speech synthesis [158], [324], [325] and voice conversion [158]. We source text prompts for text-to-speech synthesis from the LibriTTS dataset [326], and use the speakers from this dataset for voice conversion, with target audio sourced from the dataset of English dialects [323]. **Real Data.** We sample real videos for our experiments from two datasets, HDTF [322] and TalkingHead-1KH [327]. We include all available videos from HDTF, while sampling an additional 2,000 videos from TalkingHead-1KH.

7. Available at: <https://github.com/CroituruAlin/biodeep>

8. <https://civitai.com/models/139562/realvisxl-v50>

TABLE 7

Results (in terms of AUC) of four state-of-the-art deepfake detectors on the original test sets versus BioDeepAV. UCF [200], RECCE [233], TALL [263], F3Net [229], StA [257] and XceptionNet [296] are originally tested on FaceForensics++ [296], while MRDF is originally tested on FakeAVCeleb [311].

Method	Venue	Original Test	BioDeepAV
StA [257]	ArXiv 2024	0.9420	0.6195
XceptionNet [296]	ICCV 2019	0.9637	0.5677
F3Net [229]	ECCV 2020	0.9449	0.5010
RECCE [233]	CVPR 2022	0.9422	0.5001
TALL [263]	ICCV 2023	0.9987	0.4935
UCF [200]	ICCV 2023	0.9527	0.4882
MRDF [289]	ICASSP 2024	0.9243	0.5852

**Experiments.** We run the experiments using the DeepfakeBench benchmark [328], which implements state-of-the-art deepfake detectors. For our analysis, we choose three image-based detectors, namely UCF [200], RECCE [233] and a model based on XceptionNet [296], one detector applied on the frequency domain, namely F3Net [229], two video-based detectors, namely TALL [263] and StA [257], and one audio-visual detector, namely MRDF [289]. MRDF is not implemented in the DeepFakeBench benchmark, so we employ the official implementation in our experiments. UCF [200], RECCE [233], TALL [263], F3Net [229], StA [257] and XceptionNet [296] are trained on FaceForensics++ [296], while MRDF [289] is trained on FakeAVCeleb [311]. In Table 7, we report the video AUC of these detectors on BioDeepAV and their original test sets, respectively. The considered methods always surpass the 90% threshold when tested in-domain, yet all methods register drastic performance drops (higher than 30%) on BioDeepAV. The findings clearly demonstrate that current detectors struggle to identify the authenticity of talking faces generated by the novel (unseen) generative models included in BioDeepAV, highlighting the need for further research in this area.

## 7 CONCLUSIONS AND FUTURE DIRECTIONS

In this paper, we reviewed deepfake generation and detection methods, constructing a comprehensive taxonomy of methods across image, video, audio and multimodal domains. After discussing the methods included in our taxonomy, we turned our attention to datasets used for deepfake detection, with a particular focus on the results reported by top performing models. Moreover, we evaluated some of the best methods on our novel benchmark, BioDeepAV, aiming to assess the generalization capacity of current deepfake detectors to out-of-distribution data. The results show that the distribution gap can greatly affect state-of-the-art deepfake detectors, pinpointing the need for more robust models.

**Future directions.** Based on the observed gaps in deepfake literature, there are several directions which we recommend exploring in future work. The most important future direction is the development of deepfake detectors that can generalize across multiple generative tools. Our new benchmark, BioDeepAV, will come in handy to test the generalization capacity of deepfake detection models in the future. Another area that is not sufficiently explored is the development of explainable deepfake detectors. Knowing when and why deepfake detectors fail is an important aspect for making deepfake detectors more user friendly, but this has

often been disregarded. While current detectors mostly rely on deep neural networks, an important downside of such models is that they are unable to quantify their uncertainty. To this end, studying approaches to calibrate deepfake detectors could lead to the development of enhanced models able to indicate when their prediction is uncertain.

## ACKNOWLEDGMENTS

This work was supported by a grant of the Ministry of Research, Innovation and Digitization, CCCDI - UEFIS-CDI, project number PN-IV-P6-6.3-SOL-2024-2-0227, within PNCDI IV.

## REFERENCES

- [1] F.-A. Croitoru, V. Hondru, R. T. Ionescu, *et al.*, "Diffusion Models in Vision: A Survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [2] J. Chen, J. Yu, C. Gu, *et al.*, "Pixart- $\alpha$ : Fast training of diffusion transformer for photorealistic text-to-image synthesis," in *ICLR*, 2024.
- [3] R. Rombach, A. Blattmann, D. Lorenz, *et al.*, "High-Resolution Image Synthesis with Latent Diffusion Models," in *CVPR*, 2022.
- [4] D. Podell, Z. English, K. Lacey, *et al.*, "SDXL: Improving Latent Diffusion Models for High-Resolution Image Synthesis," in *ICLR*, 2024.
- [5] A. Sauer, K. Schwarz, and A. Geiger, "StyleGAN-XL: Scaling StyleGAN to Large Diverse Datasets," in *SIGGRAPH*, 2022.
- [6] P. Esser, R. Rombach, and B. Ommer, "Taming transformers for high-resolution image synthesis," in *CVPR*, 2021.
- [7] Y. Li, C. Ma, Y. Yan, *et al.*, "3D-Aware Face Swapping," in *CVPR*, 2023.
- [8] Z. Liu, M. Li, Y. Zhang, *et al.*, "Fine-Grained Face Swapping Via Regional GAN Inversion," in *CVPR*, 2023.
- [9] J. Kim, J. Lee, and B. Zhang, "Smooth-Swap: A Simple Enhancement for Face-Swapping with Smoothness," in *CVPR*, 2022.
- [10] K. Shiohara, X. Yang, and T. Taketomi, "BlendFace: Re-designing Identity Encoders for Face-Swapping," in *ICCV*, 2023.
- [11] F. Paraperas Papantoniou, A. Lattas, *et al.*, "Arc2Face: A Foundation Model for ID-Consistent Human Faces," in *ECCV*, 2024.
- [12] W. Zhao, Y. Rao, W. Shi, *et al.*, "DiffSwap: High-Fidelity and Controllable Face Swapping via 3D-Aware Masked Diffusion," in *CVPR*, 2023.
- [13] R. Liu, B. Ma, W. Zhang, *et al.*, "Towards a simultaneous and granular identity-expression control in personalized face generation," in *CVPR*, 2024.
- [14] Y. Nirkin, Y. Keller, and T. Hassner, "FSGAN: Subject Agnostic Face Swapping and Reenactment," in *ICCV*, 2019.
- [15] S. Bounareli, C. Tzelepis, V. Argyriou, *et al.*, "HyperReenact: One-Shot Reenactment via Jointly Learning to Refine and Retarget Faces," in *ICCV*, 2023.
- [16] T. Oorloff and Y. Yacoub, "Robust One-Shot Face Video Reenactment using Hybrid Latent Spaces of StyleGAN2," in *ICCV*, 2023.
- [17] Y. Ma, S. Zhang, J. Wang, *et al.*, "DreamTalk: When Expressive Talking Head Generation Meets Diffusion Probabilistic Models," *arXiv:2312.09767*, 2023.
- [18] M. Xu, H. Li, Q. Su, *et al.*, "Hallo: Hierarchical audio-driven visual synthesis for portrait image animation," *arXiv:2406.08801*, 2024.
- [19] K. Shiohara and T. Yamasaki, "Detecting deepfakes with self-blended images," in *CVPR*, 2022.
- [20] L. Chen, Y. Zhang, Y. Song, *et al.*, "Self-supervised learning of adversarial example: Towards good generalizations for deepfake detection," in *CVPR*, 2022.
- [21] S. Dong, J. Wang, R. Ji, *et al.*, "Implicit identity leakage: The stumbling block to improving deepfake detection generalization," in *CVPR*, 2023.
- [22] L. Lin, X. He, Y. Ju, *et al.*, "Preserving fairness generalization in deepfake detection," in *CVPR*, 2024.
- [23] Y. Xu, J. Liang, L. Sheng, *et al.*, "Learning spatiotemporal inconsistency via thumbnail layout for face deepfake detection," *International Journal of Computer Vision*, 2024.
- [24] Y. Zhou and S.-N. Lim, "Joint audio-visual deepfake detection," in *ICCV*, 2021.
- [25] D. Cozzolino, A. Pianese, M. Niefser, *et al.*, "Audio-visual person-of-interest deepfake detection," in *CVPR*, 2023.
- [26] B. Goyal, N. S. Gill, P. Gulia, *et al.*, "Detection of fake accounts on social media using multimodal data with deep learning," *IEEE Transactions on Computational Social Systems*, 2023.
- [27] R. Natsume, T. Yatagawa, and S. Morishima, "RSGAN: Face swapping and editing using face and hair representation in latent spaces," in *SIGGRAPH*, 2018.
- [28] I. Skorokhodov, S. Tulyakov, and M. Elhoseiny, "StyleGAN-V: A Continuous Video Generator with the Price, Image Quality and Perks of StyleGAN2," in *CVPR*, 2022.
- [29] C. Kim, J. Lee, S. Joung, *et al.*, "InstantFamily: Masked Attention for Zero-shot Multi-ID Image Generation," *arXiv:2404.19427*, 2024.
- [30] Y. Wang, W. Zhang, J. Zheng, *et al.*, "High-fidelity person-centric subject-to-image synthesis," in *CVPR*, 2024.
- [31] M. Stypulkowski, K. Vougioukas, S. He, *et al.*, "Diffused Heads: Diffusion Models Beat GANs on Talking-Face Generation," in *WACV*, 2024.
- [32] J. Wang, Z. Wu, W. Ouyang, *et al.*, "M2TR: Multi-modal Multi-scale Transformers for Deepfake Detection," in *ICMR*, 2022.
- [33] X. Dong, J. Bao, D. Chen, *et al.*, "Protecting celebrities from deepfake with identity consistency transformer," in *CVPR*, 2022.
- [34] Z. Cai, S. Ghosh, K. Stefanov, *et al.*, "MARLIN: Masked Autoencoder for facial video Representation Learning," in *CVPR*, 2023.
- [35] Y. Zheng, J. Bao, D. Chen, *et al.*, "Exploring temporal coherence for more general video face forgery detection," in *ICCV*, 2021.
- [36] J. Choi, T. Kim, Y. Jeong, *et al.*, "Exploiting Style Latent Flows for Generalizing Deepfake Video Detection," in *CVPR*, 2024.
- [37] R. Shao, T. Wu, and Z. Liu, "Detecting and recovering sequential deepfake manipulation," in *ECCV*, 2022.
- [38] D. Güera and E. J. Delp, "Deepfake Video Detection Using Recurrent Neural Networks," in *AVSS*, 2018.
- [39] B. Liu, B. Liu, M. Ding, *et al.*, "TI2Net: Temporal Identity Inconsistency Network for Deepfake Detection," in *WACV*, 2023.
- [40] A. Das, K. A. Viji, and L. Sebastian, "A survey on deepfake video detection techniques using deep learning," in *ICNGIS*, 2022.
- [41] A. Heidari, N. J. Navimipour, H. Dag, *et al.*, "Deepfake detection using deep learning methods: A systematic and comprehensive review," *WIREs Data Mining and Knowledge Discovery*, 2024.
- [42] A. Kaur, A. Noori Hoshyar, V. Saikrishna, *et al.*, "Deepfake video detection: challenges and opportunities," *Artificial Intelligence Review*, 2024.
- [43] W. Lei, J. Wang, F. Ma, *et al.*, "A comprehensive survey on human video generation: Challenges, methods, and insights," *arXiv:2407.08428*, 2024.
- [44] M. Li, Y. Ahmadiadli, and X.-P. Zhang, "Audio anti-spoofing detection: A survey," *arXiv:2404.13914*, 2024.
- [45] C. Li, D. Huang, Z. Lu, *et al.*, "A survey on long video generation: Challenges, methods, and prospects," *arXiv:2403.16407*, 2024.
- [46] M. Masood, M. Nawaz, K. M. Malik, *et al.*, "Deepfakes generation and detection: State-of-the-art, open challenges, countermeasures, and way forward," *Applied Intelligence*, 2023.
- [47] Y. Patel, S. Tanwar, R. Gupta, *et al.*, "Deepfake generation and detection case study and challenges," *IEEE Access*, 2023.
- [48] G. Pei, J. Zhang, M. Hu, *et al.*, "Deepfake generation and detection: A benchmark and survey," *arXiv:2403.17881*, 2024.
- [49] J. W. Seow, M. K. Lim, R. Phan, *et al.*, "A comprehensive overview of deepfake: Generation, detection, datasets, and opportunities," *Neurocomputing*, 2022.
- [50] J. Yi, C. Wang, J. Tao, *et al.*, "Audio deepfake detection: A survey," *arXiv:2308.14970*, 2023.
- [51] T. Zhang, "Deepfake generation and detection, a survey," *Multi-media Tools and Applications*, 2022.
- [52] T. Brooks, B. Peebles, C. Holmes, *et al.*, "Video generation models as world simulators," tech. rep., OpenAI, 2024.
- [53] A. Q. Nichol, P. Dhariwal, A. Ramesh, *et al.*, "GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models," in *ICML*, 2022.
- [54] Y. Shen, P. Luo, J. Yan, *et al.*, "FaceID-GAN: Learning a Symmetry Three-Player GAN for Identity-Preserving Face Synthesis," in *CVPR*, 2018.
- [55] T. Karras, T. Aila, S. Laine, *et al.*, "Progressive Growing of GANs for Improved Quality, Stability, and Variation," in *ICLR*, 2018.
- [56] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.

- [57] T. Karras, S. Laine, M. Aittala, *et al.*, "Analyzing and Improving the Image Quality of StyleGAN," in *CVPR*, 2020.
- [58] T. Karras, M. Aittala, S. Laine, *et al.*, "Alias-free generative adversarial networks," in *NeurIPS*, 2021.
- [59] J. Fu, S. Li, Y. Jiang, *et al.*, "StyleGAN-Human: A Data-Centric Odyssey of Human Generation," in *ECCV*, 2022.
- [60] H. Xu, G. Song, Z. Jiang, *et al.*, "OmniAvatar: Geometry-Guided Controllable 3D Head Synthesis," in *CVPR*, 2023.
- [61] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *ICCV*, 2017.
- [62] A. Sauer, K. Chitta, J. Müller, *et al.*, "Projected GANs converge faster," in *NeurIPS*, 2021.
- [63] P. Dhariwal and A. Nichol, "Diffusion models beat GANs on image synthesis," in *NeurIPS*, 2021.
- [64] I. Goodfellow, J. Pouget-Abadie, M. Mirza, *et al.*, "Generative adversarial nets," in *NeurIPS*, 2014.
- [65] J. Bao, D. Chen, F. Wen, *et al.*, "Towards open-set identity preserving face synthesis," in *CVPR*, 2018.
- [66] R. Chen, X. Chen, B. Ni, *et al.*, "SimSwap: An Efficient Framework For High Fidelity Face Swapping," in *ACMMM*, 2020.
- [67] X. Ren, X. Chen, P. Yao, *et al.*, "Reinforced disentanglement for face swapping without skip connection," in *ICCV*, 2023.
- [68] F. Rosberg, E. Aksoy, F. Alonso-Fernandez, *et al.*, "FaceDancer: Pose- and Occlusion-Aware High Fidelity Face Swapping," in *WACV*, 2023.
- [69] L. Li, J. Bao, H. Yang, *et al.*, "Advancing high fidelity identity swapping for forgery detection," in *CVPR*, 2020.
- [70] H. Zeng, W. Zhang, C. Fan, *et al.*, "FlowFace: semantic flow-guided shape-aware face swapping," in *AAAI*, 2023.
- [71] S.-M. Yoo, T.-M. Choi, J.-W. Choi, *et al.*, "FastSwap: A Lightweight One-Stage Framework for Real-Time Face Swapping," in *WACV*, 2023.
- [72] K. Cui, R. Wu, F. Zhan, *et al.*, "Face Transformer: Towards High Fidelity and Accurate Face Swapping," in *CVPRW*, 2023.
- [73] G. Yuan, M. Li, Y. Zhang, *et al.*, "ReliableSwap: Boosting General Face Swapping Via Reliable Supervision," *arXiv:2306.05356*, 2023.
- [74] W. Cao, T. Wang, A. Dong, *et al.*, "TransFS: Face Swapping Using Transformer," in *FG*, 2023.
- [75] K. He, X. Chen, S. Xie, *et al.*, "Masked autoencoders are scalable vision learners," in *CVPR*, 2022.
- [76] G. Gao, H. Huang, C. Fu, *et al.*, "Information bottleneck disentanglement for identity swapping," in *CVPR*, 2021.
- [77] Y. Zhu, Q. Li, J. Wang, *et al.*, "One shot face swapping on megapixels," in *CVPR*, 2021.
- [78] D. Jiang, D. Song, R. Tong, *et al.*, "StyleIPSB: Identity-Preserving Semantic Basis of StyleGAN for High Fidelity Face Swapping," in *CVPR*, 2023.
- [79] E. R. Chan, M. Monteiro, P. Kellnhofer, *et al.*, "pi-GAN: Periodic Implicit Generative Adversarial Networks for 3D-Aware Image Synthesis," in *CVPR*, 2021.
- [80] J. R. A. Moniz, C. Beckham, S. Rajotte, *et al.*, "Unsupervised Depth Estimation, 3D Face Rotation and Replacement," in *NeurIPS*, 2018.
- [81] Q. Sun, A. Tewari, W. Xu, *et al.*, "A hybrid model for identity obfuscation by face replacement," in *ECCV*, 2018.
- [82] Z. Chen, L. Xie, S. Pang, *et al.*, "MagDR: Mask-guided Detection and Reconstruction for Defending Deepfakes," in *CVPR*, 2021.
- [83] J.-Y. Zhu, T. Park, P. Isola, *et al.*, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *ICCV*, 2017.
- [84] Y. Choi, M. Choi, M. Kim, *et al.*, "StarGAN: Unified Generative Adversarial Networks for Multi-domain Image-to-Image Translation," in *CVPR*, 2018.
- [85] Y. Choi, Y. Uh, J. Yoo, *et al.*, "StarGAN v2: Diverse Image Synthesis for Multiple Domains," in *CVPR*, 2020.
- [86] G.-S. Hsu, C.-H. Tsai, and H.-Y. Wu, "Dual-generator face reenactment," in *CVPR*, 2022.
- [87] O. Tov, Y. Alaluf, Y. Nitzan, *et al.*, "Designing an encoder for StyleGAN image manipulation," *ACM Transactions on Graphics*, 2021.
- [88] A. Suwała, B. Wójcik, M. Proszewska, *et al.*, "Face Identity-Aware Disentanglement in StyleGAN," in *WACV*, 2024.
- [89] C. Saharia, W. Chan, S. Saxena, *et al.*, "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding," in *NeurIPS*, 2022.
- [90] N. Ruiz, Y. Li, V. Jampani, *et al.*, "DreamBooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation," in *CVPR*, 2023.
- [91] Z. Chen, S. Fang, W. Liu, *et al.*, "DreamIdentity: Enhanced Editability for Efficient Face-Identity Preserved Image Generation," in *AAAI*, 2024.
- [92] X. Peng, J. Zhu, B. Jiang, *et al.*, "PortraitBooth: A Versatile Portrait Model for Fast Identity-Preserved Personalization," in *CVPR*, 2024.
- [93] J. Ma, J. Liang, C. Chen, *et al.*, "Subject-Diffusion: Open Domain Personalized Text-to-Image Generation without Test-time Fine-tuning," in *SIGGRAPH*, 2024.
- [94] Y. Liu, C. Yu, L. Shang, *et al.*, "FaceChain: A Playground for Human-centric Artificial Intelligence Generated Content," *arXiv:2308.14256*, 2023.
- [95] F. Boutros, J. H. Grebe, A. Kuijper, *et al.*, "IDiff-Face: Synthetic-based Face Recognition through Fizzy Identity-Conditioned Diffusion Model," in *ICCV*, 2023.
- [96] Y. Han, J. Zhang, J. Zhu, *et al.*, "A Generalist FaceX via Learning Unified Facial Representation," *arXiv:2401.00551*, 2023.
- [97] H. Lin, "DreamSalon: A Staged Diffusion Framework for Preserving Identity-Context in Editable Face Generation," in *CVPR*, 2024.
- [98] Q. Wang, X. Bai, H. Wang, *et al.*, "InstantID: Zero-shot Identity-Preserving Generation in Seconds," *arXiv:2401.07519*, 2024.
- [99] Y. Wu, Z. Li, H. Zheng, *et al.*, "Infinite-ID: Identity-preserved Personalization via ID-semantics Decoupling Paradigm," *arXiv:2403.11781*, 2024.
- [100] Q. Wang, X. Jia, X. Li, *et al.*, "StableIdentity: Inserting Anybody into Anywhere at First Sight," *arXiv:2401.15975*, 2024.
- [101] Z. Guo, Y. Wu, Z. Chen, *et al.*, "PuLID: Pure and Lightning ID Customization via Contrastive Alignment," in *NeurIPS*, 2024.
- [102] J. Gu, Y. Wang, N. Zhao, *et al.*, "SwapAnything: Enabling Arbitrary Object Swapping in Personalized Visual Editing," in *ECCV*, 2024.
- [103] K.-C. Wang, D. Ostashev, Y. Fang, *et al.*, "MoA: Mixture-of-Attention for Subject-Context Disentanglement in Personalized Image Generation," *arXiv:2404.11565*, 2024.
- [104] Q. Wang, B. Li, X. Li, *et al.*, "CharacterFactory: Sampling Consistent Characters with GANs for Diffusion Models," *arXiv:2404.15677*, 2024.
- [105] W. Chen, J. Zhang, J. Wu, *et al.*, "ID-Aligner: Enhancing Identity-Preserving Text-to-Image Generation with Reward Feedback Learning," *arXiv:2404.15449*, 2024.
- [106] Z. Huang, H. Fan, L. Wang, *et al.*, "From parts to whole: A unified reference framework for controllable human image generation," *arXiv:2404.15267*, 2024.
- [107] J. Huang, X. Dong, W. Song, *et al.*, "ConsistentID: Portrait Generation with Multimodal Fine-Grained Identity Preserving," *arXiv:2404.16771*, 2024.
- [108] J. He, Y. Geng, and L. Bo, "UniPortrait: A Unified Framework for Identity-Preserving Single- and Multi-Human Image Personalization," *arXiv:2408.05939*, 2024.
- [109] Z. Li, M. Cao, X. Wang, *et al.*, "PhotoMaker: Customizing Realistic Human Photos via Stacked ID Embedding," in *CVPR*, 2024.
- [110] Y. Wei, Z. Ji, J. Bai, *et al.*, "MasterWeaver: Taming Editability and Face Identity for Personalized Text-to-Image Generation," in *ECCV*, 2024.
- [111] R. Gal, Y. Alaluf, Y. Atzmon, *et al.*, "An image is worth one word: Personalizing text-to-image generation using textual inversion," in *ICLR*, 2023.
- [112] N. Ruiz, Y. Li, V. Jampani, *et al.*, "HyperDreamBooth: Hyper-Networks for Fast Personalization of Text-to-Image Models," in *CVPR*, 2024.
- [113] W. Chen, H. Hu, Y. Li, *et al.*, "Subject-driven text-to-image generation via apprenticeship learning," in *NeurIPS*, 2023.
- [114] E. J. Hu, Y. Shen, P. Wallis, *et al.*, "LoRA: Low-Rank Adaptation of Large Language Models," in *ICLR*, 2022.
- [115] D. Bitouk, N. Kumar, S. Dhillon, *et al.*, "Face swapping: automatically replacing faces in photographs," in *SIGGRAPH*, 2008.
- [116] I. Korshunova, W. Shi, J. Dambre, *et al.*, "Fast face-swap using convolutional neural networks," in *ICCV*, 2017.
- [117] Y. Wang, X. Chen, J. Zhu, *et al.*, "HifiFace: 3D Shape and Semantic Prior Guided High Fidelity Face Swapping," in *IJCAI*, 2021.
- [118] J. Bao, D. Chen, F. Wen, *et al.*, "CVAE-GAN: Fine-Grained Image Generation through Asymmetric Training," in *ICCV*, 2017.

- [119] X. Li, X. Hou, and C. C. Loy, "When StyleGAN Meets Stable Diffusion: a  $W_+$  Adapter for Personalized Image Generation," in *CVPR*, 2024.
- [120] Z. Xu, H. Zhou, Z. Hong, *et al.*, "StyleSwap: Style-Based Generator Empowers Robust Face Swapping," in *ECCV*, 2022.
- [121] Y. Gao, Y. Zhou, J. Wang, *et al.*, "High-Fidelity and Freely Controllable Talking Head Video Generation," in *CVPR*, 2023.
- [122] Y. Tian, J. Ren, M. Chai, *et al.*, "A Good Image Generator Is What You Need for High-Resolution Video Synthesis," in *ICLR*, 2021.
- [123] S. Tulyakov, M.-Y. Liu, X. Yang, *et al.*, "MoCoGAN: Decomposing Motion and Content for Video Generation," in *CVPR*, 2018.
- [124] R. Abdal, Y. Qin, and P. Wonka, "Image2StyleGAN: How to Embed Images Into the StyleGAN Latent Space?," in *ICCV*, 2019.
- [125] Z. Wu, D. Lischinski, and E. Shechtman, "StyleSpace Analysis: Disentangled Controls for StyleGAN Image Generation," in *CVPR*, 2021.
- [126] S. Yu, J. Tack, S. Mo, *et al.*, "Generating Videos with Dynamics-aware Implicit Generative Adversarial Networks," in *ICLR*, 2022.
- [127] V. Sitzmann, J. Martel, A. Bergman, *et al.*, "Implicit Neural Representations with Periodic Activation Functions," in *NeurIPS*, 2020.
- [128] T. Brooks, J. Hellsten, M. Aittala, *et al.*, "Generating Long Videos of Dynamic Scenes," in *NeurIPS*, 2022.
- [129] J. Ho, T. Salimans, A. Gritsenko, *et al.*, "Video Diffusion Models," *arXiv:2204.03458*, 2022.
- [130] A. Blattmann, R. Rombach, H. Ling, *et al.*, "Align your Latents: High-Resolution Video Synthesis with Latent Diffusion Models," in *CVPR*, 2023.
- [131] A. Blattmann, T. Dockhorn, S. Kulal, *et al.*, "Stable Video Diffusion: Scaling Latent Video Diffusion Models to Large Datasets," *arXiv:2311.15127*, 2023.
- [132] J. Z. Wu, Y. Ge, X. Wang, *et al.*, "Tune-A-Video: One-Shot Tuning of Image Diffusion Models for Text-to-Video Generation," in *ICCV*, 2023.
- [133] F. Bao, C. Xiang, G. Yue, *et al.*, "Vidu: a Highly Consistent, Dynamic and Skilled Text-to-Video Generator with Diffusion Models," *arXiv:2405.04233*, 2024.
- [134] Y. Guo, C. Yang, A. Rao, *et al.*, "AnimateDiff: Animate Your Personalized Text-to-Image Diffusion Models without Specific Tuning," in *ICLR*, 2024.
- [135] X. Ma, Y. Wang, G. Jia, *et al.*, "Latte: Latent Diffusion Transformer for Video Generation," *arXiv:2401.03048*, 2024.
- [136] W. Wang, J. Liu, Z. Lin, *et al.*, "MagicVideo-V2: Multi-Stage High-Aesthetic Video Generation," *arXiv:2401.04468*, 2024.
- [137] U. Singer, A. Polyak, T. Hayes, *et al.*, "Make-A-Video: Text-to-Video Generation without Text-Video Data," in *ICLR*, 2022.
- [138] A. Ramesh, P. Dhariwal, A. Nichol, *et al.*, "Hierarchical Text-Conditional Image Generation with CLIP Latents," *arXiv:2204.06125*, 2022.
- [139] J. Ho, W. Chan, C. Saharia, *et al.*, "Imagen Video: High Definition Video Generation with Diffusion Models," *arXiv:2210.02303*, 2022.
- [140] S. Bounareli, C. Tzelepis, V. Argyriou, *et al.*, "DiffusionAct: Controllable Diffusion Autoencoder for One-shot Face Reenactment," *arXiv:2403.17217*, 2024.
- [141] Y. Ma, H. Liu, H. Wang, *et al.*, "Follow-Your-Emoji: Fine-Controllable and Expressive Freestyle Portrait Animation," *arXiv:2406.01900*, 2024.
- [142] S. Yang, H. Li, J. Wu, *et al.*, "MegActor: Harness the Power of Raw Video for Vivid Portrait Animation," *arXiv:2405.20851*, 2024.
- [143] Y. Zhang, J. Gu, L.-W. Wang, *et al.*, "MimicMotion: High-Quality Human Motion Video Generation with Confidence-aware Pose Guidance," *arXiv:2406.19680*, 2024.
- [144] A. Rochow, M. Schwarz, and S. Behnke, "FSRT: Facial Scene Representation Transformer for Face Reenactment from Factorized Appearance Head-pose and Facial Expression Features," in *CVPR*, 2024.
- [145] R. Villegas, M. Babaeizadeh, P.-J. Kindermans, *et al.*, "Phenaki: Variable Length Video Generation From Open Domain Textual Description," in *ICLR*, 2023.
- [146] A. Van Den Oord, O. Vinyals, and K. Kavukcuoglu, "Neural discrete representation learning," in *NeurIPS*, 2017.
- [147] L. Yu, Y. Cheng, K. Sohn, *et al.*, "MAGVIT: Masked Generative Video Transformer," in *CVPR*, 2023.
- [148] W. Hong, M. Ding, W. Zheng, *et al.*, "CogVideo: Large-scale Pretraining for Text-to-Video Generation via Transformers," in *ICLR*, 2023.
- [149] Y. Jiang, S. Yang, T. L. Koh, *et al.*, "Text2Performer: Text-Driven Human Video Generation," in *ICCV*, 2023.
- [150] W. Yan, Y. Zhang, P. Abbeel, *et al.*, "VideoGPT: Video Generation using VQ-VAE and Transformers," *arXiv:2104.10157*, 2021.
- [151] B. Mildenhall, P. P. Srinivasan, M. Tancik, *et al.*, "NeRF: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, 2021.
- [152] J. Thies, M. Zollhöfer, and M. Nießner, "Deferred Neural Rendering: Image Synthesis using Neural Textures," *ACM Transactions on Graphics*, 2019.
- [153] J. Thies, M. Zollhofer, M. Stamminger, *et al.*, "Face2Face: Real-time Face Capture and Reenactment of RGB Videos," in *CVPR*, 2016.
- [154] K. Yang, K. Chen, D. Guo, *et al.*, "Face2Face  $\rho$ : Real-Time High-Resolution One-Shot Face Reenactment," in *ECCV*, 2022.
- [155] B. Zhang, C. Qi, P. Zhang, *et al.*, "MetaPortrait: Identity-Preserving Talking Head Generation with Fast Personalized Adaptation," in *CVPR*, 2023.
- [156] W. Li, L. Zhang, D. Wang, *et al.*, "One-Shot High-Fidelity Talking-Head Synthesis with Deformable Neural Radiance Field," in *CVPR*, 2023.
- [157] J. Kim, J. Kong, and J. Son, "Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech," in *ICML*, 2021.
- [158] E. Casanova, J. Weber, C. D. Shulby, *et al.*, "YourTTS: Towards zero-shot multi-speaker TTS and zero-shot voice conversion for everyone," in *ICML*, 2022.
- [159] X. Tan, J. Chen, H. Liu, *et al.*, "NaturalSpeech: End-to-End Text-to-Speech Synthesis With Human-Level Quality," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [160] S.-g. Lee, W. Ping, B. Ginsburg, *et al.*, "BigVGAN: A Universal Neural Vocoder with Large-Scale Training," in *ICLR*, 2023.
- [161] Y. Ren, Y. Ruan, X. Tan, *et al.*, "FastSpeech: Fast, Robust and Controllable Text to Speech," in *NeurIPS*, 2019.
- [162] Z. Jiang, Y. Ren, Z. Ye, *et al.*, "Mega-TTS: Zero-Shot Text-to-Speech at Scale with Intrinsic Inductive Bias," *arXiv:2306.03509*, 2023.
- [163] C. Wang, S. Chen, Y. Wu, *et al.*, "Neural codec language models are zero-shot text to speech synthesizers," *arXiv:2301.02111*, 2023.
- [164] E. Kharitonov, D. Vincent, Z. Borsos, *et al.*, "Speak, read and prompt: High-fidelity text-to-speech with minimal supervision," *Transactions of the Association for Computational Linguistics*, 2023.
- [165] D. Yang, J. Tian, X. Tan, *et al.*, "UniAudio: An Audio Foundation Model Toward Universal Audio Generation," *arXiv:2310.00704*, 2023.
- [166] R. Huang, M. W. Lam, J. Wang, *et al.*, "FastDiff: A fast conditional diffusion model for high-quality speech synthesis," in *IJCAI*, 2022.
- [167] Y. Ren, C. Hu, X. Tan, *et al.*, "FastSpeech 2: Fast and High-Quality End-to-End Text to Speech," in *ICLR*, 2020.
- [168] R. Huang, Z. Zhao, H. Liu, *et al.*, "ProDiff: Progressive fast diffusion model for high-quality text-to-speech," in *ACMMM*, 2022.
- [169] K. Shen, Z. Ju, X. Tan, *et al.*, "NaturalSpeech 2: Latent Diffusion Models are Natural and Zero-Shot Speech and Singing Synthesizers," in *ICLR*, 2024.
- [170] Z. Ju, Y. Wang, K. Shen, *et al.*, "NaturalSpeech 3: Zero-Shot Speech Synthesis with Factorized Codec and Diffusion Models," *arXiv:2403.03100*, 2024.
- [171] C. Du, Y. Guo, F. Shen, *et al.*, "UniCATS: A unified context-aware text-to-speech framework with contextual vq-diffusion and vocoding," in *AAAI*, 2024.
- [172] D. Yang, S. Liu, R. Huang, *et al.*, "InstructTTS: Modelling Expressive TTS in Discrete Latent Space with Natural Language Style Prompt," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024.
- [173] Y. Liu, M. Ott, N. Goyal, *et al.*, "RoBERTa: A robustly optimized BERT pretraining approach," *arXiv:1907.11692*, 2019.
- [174] Y. Jang, J.-H. Kim, J. Ahn, *et al.*, "Faces that speak: Jointly synthesising talking face and speech from text," in *CVPR*, 2024.
- [175] K. Cheng, X. Cun, Y. Zhang, *et al.*, "VideoReTalking: Audio-based Lip Synchronization for Talking Head Video Editing In the Wild," in *SIGGRAPH*, 2022.
- [176] S. Wang, L. Li, Y. Ding, *et al.*, "One-shot talking face generation from single-speaker audio-visual correlation learning," in *AAAI*, 2022.



- [177] J. Ling, X. Tan, L. Chen, *et al.*, "StableFace: Analyzing and Improving Motion Stability for Talking Face Generation," *IEEE Journal of Selected Topics in Signal Processing*, 2023.
- [178] Y. Gan, Z. Yang, X. Yue, *et al.*, "Efficient emotional adaptation for audio-driven talking-head generation," in *ICCV*, 2023.
- [179] H. Wei, Z. Yang, and Z. Wang, "AniPortrait: Audio-Driven Synthesis of Photorealistic Portrait Animation," *arXiv:2403.17694*, 2024.
- [180] C. Wang, K. Tian, J. Zhang, *et al.*, "V-Express: Conditional Dropout for Progressive Training of Portrait Video Generation," *arXiv:2406.02511*, 2024.
- [181] Z. Chen, J. Cao, Z. Chen, *et al.*, "EchoMimic: Lifelike Audio-Driven Portrait Animations through Editable Landmark Conditions," *arXiv:2407.08136*, 2024.
- [182] S. Xu, G. Chen, Y.-X. Guo, *et al.*, "VASA-1: Lifelike Audio-Driven Talking Faces Generated in Real Time," *arXiv:2404.10667*, 2024.
- [183] L. Tian, Q. Wang, B. Zhang, *et al.*, "EMO: Emote Portrait Alive – Generating Expressive Portrait Videos with Audio2Video Diffusion Model under Weak Conditions," *arXiv:2402.17485*, 2024.
- [184] Z. Peng, W. Hu, Y. Shi, *et al.*, "SyncTalk: The Devil is in the Synchronization for Talking Head Synthesis," in *CVPR*, 2024.
- [185] Z. Ye, J. He, Z. Jiang, *et al.*, "GeneFace++: Generalized and Stable Real-Time Audio-Driven 3D Talking Face Generation," *arXiv:2305.00787*, 2023.
- [186] S. Zhang, J. Yuan, M. Liao, *et al.*, "Text2Video: Text-driven Talking-head Video Synthesis with Personalized Phoneme-Pose Dictionary," in *ICASSP*, 2022.
- [187] M. C. Doukas, S. Zafeiriou, and V. Sharmanska, "HeadGAN: One-shot Neural Head Synthesis and Editing," in *ICCV*, 2021.
- [188] J. Wang, X. Qian, M. Zhang, *et al.*, "Seeing what you said: Talking face generation guided by a lip reading expert," in *CVPR*, 2023.
- [189] X. Liu, Q. Wu, H. Zhou, *et al.*, "Audio-driven co-speech gesture video generation," in *NeurIPS*, 2022.
- [190] Y. Lu, J. Chai, and X. Cao, "Live Speech Portraits: Real-Time Photorealistic Talking-Head Animation," *ACM Transactions on Graphics*, 2021.
- [191] S. Gururani, A. Mallya, T.-C. Wang, *et al.*, "Space: Speech-driven portrait animation with controllable expression," in *ICCV*, 2023.
- [192] W. Zhang, X. Cun, X. Wang, *et al.*, "SadTalker: Learning Realistic 3D Motion Coefficients for Stylized Audio-Driven Single Image Talking Face Animation," in *CVPR*, 2023.
- [193] G. Hwang, S. Hong, S. Lee, *et al.*, "DisCoHead: Audio-and-Video-Driven Talking Head Generation by Disentangled Control of Head Pose and Facial Expressions," in *ICASSP*, 2023.
- [194] Z. Peng, H. Wu, Z. Song, *et al.*, "EmoTalk: Speech-Driven Emotional Disentanglement for 3D Face Animation," in *ICCV*, 2023.
- [195] W. Zhong, C. Fang, Y. Cai, *et al.*, "Identity-preserving talking face generation with landmark and appearance priors," in *CVPR*, 2023.
- [196] Z. Yan, Y. Luo, S. Lyu, *et al.*, "Transcending Forgery Specificity with Latent Space Augmentation for Generalizable Deepfake Detection," in *CVPR*, 2024.
- [197] C. Tan, H. Liu, Y. Zhao, *et al.*, "Rethinking the Up-Sampling Operations in CNN-Based Generative Network for Generalizable Deepfake Detection," in *CVPR*, 2024.
- [198] D. Nguyen, N. Mejri, I. Singh, *et al.*, "LAA-Net: Localized Artifact Attention Network for Quality-Agnostic and Generalizable Deepfake Detection," in *CVPR*, 2024.
- [199] K. Yao, J. Wang, B. Diao, *et al.*, "Towards understanding the generalization of deepfake detectors from a game-theoretical view," in *ICCV*, 2023.
- [200] Z. Yan, Y. Zhang, Y. Fan, *et al.*, "UCF: Uncovering Common Features for Generalizable Deepfake Detection," in *ICCV*, 2023.
- [201] L. Chen, Y. Zhang, Y. Song, *et al.*, "OST: Improving Generalization of DeepFake Detection via One-Shot Test-Time Training," in *NeurIPS*, 2022.
- [202] C. Tan, Y. Zhao, S. Wei, *et al.*, "Frequency-aware deepfake detection: Improving generalizability through frequency space domain learning," in *AAAI*, 2024.
- [203] B. M. Le and S. S. Woo, "ADD: Frequency attention and multi-view based knowledge distillation to detect low-quality compressed deepfake images," in *AAAI*, 2022.
- [204] M. Kim, S. Tariq, and S. S. Woo, "FReTAL: Generalizing Deepfake Detection using Knowledge Distillation and Representation Learning," in *CVPRW*, 2021.
- [205] Y. Xu, K. Raja, L. Verdoliva, *et al.*, "Learning pairwise interaction for generalizable deepfake detection," in *WACVW*, 2023.
- [206] M. Du, S. Pentyala, Y. Li, *et al.*, "Towards generalizable deepfake detection with locality-aware autoencoder," in *CIKM*, 2019.
- [207] B. Huang, Z. Wang, J. Yang, *et al.*, "Implicit identity driven deepfake face swapping detection," in *CVPR*, 2023.
- [208] H. Zhao, T. Wei, W. Zhou, *et al.*, "Multi-attentional deepfake detection," in *CVPR*, 2021.
- [209] S. Dong, J. Wang, J. Liang, *et al.*, "Explaining deepfake detection by analysing image matching," in *ECCV*, 2022.
- [210] B. M. Le and S. S. Woo, "Quality-agnostic deepfake detection with intra-model collaborative learning," in *ICCV*, 2023.
- [211] N. Larue, N.-S. Vu, V. Struc, *et al.*, "SeeABLE: Soft Discrepancies and Bounded Contrastive Learning for Exposing Deepfakes," in *ICCV*, 2023.
- [212] Z. Sun, S. Chen, T. Yao, *et al.*, "Contrastive pseudo learning for open-world deepfake attribution," in *ICCV*, 2023.
- [213] T. Zhao, X. Xu, M. Xu, *et al.*, "Learning self-consistency for deepfake detection," in *ICCV*, 2021.
- [214] A. Hooda, N. Mangaokar, R. Feng, *et al.*, "D4: Detection of adversarial diffusion deepfakes using disjoint ensembles," in *WACV*, 2024.
- [215] Y. Ju, S. Hu, S. Jia, *et al.*, "Improving fairness in deepfake detection," in *WACV*, 2024.
- [216] D. Tantarú, E. Oneata, and D. Oneata, "Weakly-supervised deepfake localization in diffusion-generated images," in *WACV*, 2024.
- [217] L. Trinh, M. Tsang, S. Rambhatla, *et al.*, "Interpretable and trustworthy deepfake detection via dynamic prototypes," in *WACV*, 2021.
- [218] Z. Ba, Q. Liu, Z. Liu, *et al.*, "Exposing the deception: Uncovering more forgery clues for deepfake detection," in *AAAI*, 2024.
- [219] T. Yang, Z. Huang, J. Cao, *et al.*, "Deepfake network architecture attribution," in *AAAI*, 2022.
- [220] Y. Nirkin, L. Wolf, Y. Keller, *et al.*, "Deepfake detection based on discrepancies between faces and their context," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [221] R. Lanzino, F. Fontana, A. Diko, *et al.*, "Faster than lies: Real-time deepfake detection using binary neural networks," in *CVPRW*, 2024.
- [222] A. Ciamarra, R. Caldelli, F. Becattini, *et al.*, "Deepfake detection by exploiting surface anomalies: The surfake approach," in *WACVW*, 2024.
- [223] Y. Jeong, D. Kim, S. Min, *et al.*, "BiHPF: Bilateral High-Pass Filters for Robust Deepfake Detection," in *WACV*, 2022.
- [224] M. Tan, "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," in *ICML*, 2019.
- [225] F. Chollet, "Xception: Deep Learning With Depthwise Separable Convolutions," in *CVPR*, 2017.
- [226] K. He, X. Zhang, S. Ren, *et al.*, "Deep residual learning for image recognition," in *CVPR*, 2015.
- [227] Z. Liu, X. Qi, and P. H. Torr, "Global texture enhancement for fake face detection in the wild," in *CVPR*, 2020.
- [228] Y. Luo, Y. Zhang, J. Yan, and W. Liu, "Generalizing face forgery detection with high-frequency features," in *CVPR*, 2021.
- [229] Y. Qian, G. Yin, L. Sheng, *et al.*, "Thinking in frequency: Face forgery detection by mining frequency-aware clues," in *ECCV*, 2020.
- [230] S. Hussain, P. Neekhar, M. Jere, *et al.*, "Adversarial deepfakes: Evaluating vulnerability of deepfake detectors to adversarial examples," in *WACV*, 2021.
- [231] Z. Yang, J. Liang, Y. Xu, *et al.*, "Masked relation learning for deepfake detection," *IEEE Transactions on Information Forensics and Security*, 2023.
- [232] Y. Wang, K. Yu, C. Chen, *et al.*, "Dynamic graph learning with content-guided spatial-frequency relation reasoning for deepfake detection," in *CVPR*, 2023.
- [233] J. Cao, C. Ma, T. Yao, *et al.*, "End-to-end reconstruction-classification learning for face forgery detection," in *CVPR*, 2022.
- [234] A. Dosovitskiy, L. Beyer, A. Kolesnikov, *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *ICLR*, 2021.
- [235] A. Aghasanli, D. Kangin, and P. Angelov, "Interpretable-through-prototypes deepfake detection for diffusion models," in *ICCVW*, 2023.
- [236] C. Hong, Y. Hsu, and T. Liu, "Contrastive learning for deepfake classification and localization via multi-label ranking," in *CVPR*, 2024.

- [237] S. Kamat, S. Agarwal, T. Darrell, *et al.*, "Revisiting generalizability in deepfake detection: Improving metrics and stabilizing transfer," in *ICCVW*, 2023.
- [238] Y. Jeong, D. Kim, Y. Ro, *et al.*, "FrePGAN: Robust Deepfake Detection Using Frequency-Level Perturbations," in *AAAI*, 2022.
- [239] F. Lugstein, S. Baier, G. Bachinger, *et al.*, "PRNU-based Deepfake Detection," in *IH&MMSec*, 2021.
- [240] L. Guarnera, O. Giudice, and S. Battiato, "Deepfake detection by analyzing convolutional traces," in *CVPRW*, 2020.
- [241] S. Agarwal, H. Farid, T. El-Gaaly, *et al.*, "Detecting Deep-Fake Videos from Appearance and Behavior," in *WIFS*, 2020.
- [242] D. Cozzolino, A. Rössler, J. Thies, *et al.*, "ID-Reveal: Identity-aware DeepFake Video Detection," in *ICCV*, 2021.
- [243] Z. Gu, Y. Chen, T. Yao, *et al.*, "Delving into the Local: Dynamic Inconsistency Learning for DeepFake Video Detection," in *AAAI*, 2022.
- [244] Z. Gu, T. Yao, Y. Chen, *et al.*, "Hierarchical Contrastive Inconsistency Learning for Deepfake Video Detection," in *ECCV*, 2022.
- [245] Y. Zhao, W. Ge, W. Li, *et al.*, "Capturing the Persistence of Facial Expression Features for Deepfake Video Detection," in *ICICS*, 2020.
- [246] A. Haliassos, K. Vougioukas, S. Petridis, *et al.*, "Lips don't lie: A generalisable and robust approach to face forgery detection," in *CVPR*, 2021.
- [247] A. Haliassos, R. Mira, S. Petridis, *et al.*, "Leveraging real talking faces via self-supervision for robust forgery detection," in *CVPR*, 2022.
- [248] I. Demir and U. A. Çiftçi, "How Do Deepfakes Move? Motion Magnification for Deepfake Source Detection," in *WACV*, 2024.
- [249] E. Sabir, J. Cheng, A. Jaiswal, *et al.*, "Recurrent Convolutional Strategies for Face Manipulation Detection in Videos," in *CVPR*, 2019.
- [250] J. Hu, X. Liao, J. Liang, *et al.*, "FlNfer: Frame Inference-Based Deepfake Detection for High-Visual-Quality Videos," in *AAAI*, 2022.
- [251] D. M. Montserrat, H. Hao, S. K. Yarlagadda, *et al.*, "Deepfakes Detection with Automatic Face Weighting," in *CVPR*, 2020.
- [252] I. Amerini and R. Caldelli, "Exploiting Prediction Error Inconsistencies through LSTM-based Classifiers to Detect Deepfake Videos," in *IH&MMSec*, 2020.
- [253] I. Masi, A. Killekar, R. M. Mascarenhas, *et al.*, "Two-branch Recurrent Network for Isolating Deepfakes in Videos," in *ECCV*, 2020.
- [254] J. Deng, J. Guo, N. Xue, *et al.*, "ArcFace: Additive Angular Margin Loss for Deep Face Recognition," in *CVPR*, 2019.
- [255] N. Bonettini, E. D. Cannas, S. Mandelli, *et al.*, "Video face manipulation detection through ensemble of CNNs," in *ICPR*, 2021.
- [256] T. Wang and K. P. Chow, "Noise Based Deepfake Detection via Multi-Head Relative-Interaction," in *AAAI*, 2023.
- [257] Z. Yan, Y. Zhao, S. Chen, *et al.*, "Generalizing Deepfake Video Detection with Plug-and-Play: Video-Level Blending and Spatiotemporal Adapter Tuning," *arXiv:2408.17065*, 2024.
- [258] A. Coccomini, N. Messina, C. Gennaro, *et al.*, "Combining EfficientNet and Vision Transformers for Video Deepfake Detection," in *ICIAP*, 2022.
- [259] A. Vaswani, N. Shazeer, N. Parmar, *et al.*, "Attention is all you need," in *NeurIPS*, 2017.
- [260] J. Guan, H. Zhou, Z. Hong, *et al.*, "Delving into sequential patches for deepfake detection," in *NeurIPS*, 2022.
- [261] E. Richardson, Y. Alaluf, O. Patashnik, *et al.*, "Encoding in Style: a StyleGAN Encoder for Image-to-Image Translation," in *CVPR*, 2021.
- [262] L. Tan, Y. Wang, J. Wang, *et al.*, "Deepfake Video Detection via Facial Action Dependencies Estimation," in *AAAI*, 2023.
- [263] Y. Xu, J. Liang, G. Jia, *et al.*, "TALL: Thumbnail Layout for Deepfake Video Detection," in *ICCV*, 2023.
- [264] Z. Liu, Y. Lin, Y. Cao, *et al.*, "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows," in *ICCV*, 2021.
- [265] H. Tak, J. Patino, M. Todisco, *et al.*, "End-to-end anti-spoofing with RawNet2," in *ICASSP*, 2021.
- [266] C. Wang, J. Yi, J. Tao, *et al.*, "TO-RawNet: improving RawNet with TCN and orthogonal regularization for fake audio detection," in *INTERSPEECH*, 2023.
- [267] E. Conti, D. Salvi, C. Borrelli, *et al.*, "Deepfake speech detection through emotion recognition: a semantic approach," in *ICASSP*, 2022.
- [268] G. Hua, A. B. J. Teoh, and H. Zhang, "Towards end-to-end synthetic speech detection," *IEEE Signal Processing Letters*, 2021.
- [269] H. Tak, J.-w. Jung, J. Patino, *et al.*, "Graph attention networks for anti-spoofing," in *INTERSPEECH*, 2021.
- [270] H. Tak, J.-w. Jung, J. Patino, *et al.*, "End-to-end spectro-temporal graph attention networks for speaker verification anti-spoofing and speech deepfake detection," in *ASVSPOOF*, 2021.
- [271] J.-w. Jung, H.-S. Heo, H. Tak, *et al.*, "AASIST: Audio Anti-Spoofing using Integrated Spectro-Temporal Graph Attention Networks," in *ICASSP*, 2022.
- [272] F. Chen, S. Deng, T. Zheng, *et al.*, "Graph-based spectro-temporal dependency modeling for anti-spoofing," in *ICASSP*, 2023.
- [273] E. R. Bartusiak and E. J. Delp, "Synthesized speech detection using convolutional transformer-based spectrogram analysis," in *ACSSC*, 2021.
- [274] E. R. Bartusiak and E. J. Delp, "Transformer-based speech synthesizer attribution in an open set scenario," in *JCMMLA*, 2022.
- [275] J. M. Martín-Doñas and A. Álvarez, "The Vicomtech Audio Deepfake Detection System Based on Wav2vec2 for the 2022 ADD Challenge," in *ICASSP*, 2022.
- [276] Z. Cai, W. Wang, and M. Li, "Waveform boundary detection for partially spoofed audio," in *ICASSP*, 2023.
- [277] Z. Zhang, X. Yi, and X. Zhao, "Fake speech detection using residual network with transformer encoder," in *IH&MMSec*, 2021.
- [278] X. Liu, M. Liu, L. Wang, *et al.*, "Leveraging positional-related local-global dependency for synthetic speech detection," in *ICASSP*, 2023.
- [279] R. Wang, F. Juefei-Xu, Y. Huang, *et al.*, "DeepSonar: Towards Effective and Robust Detection of AI-Synthesized Fake Voices," in *ACMMM*, 2020.
- [280] X. Zhang, J. Yi, C. Wang, *et al.*, "What to remember: Self-adaptive continual learning for audio deepfake detection," in *AAAI*, 2024.
- [281] M. A. Raza and K. M. Malik, "Multimodaltrace: Deepfake detection using audiovisual representation learning," in *CVPR*, 2023.
- [282] M. Kihal and L. Hamza, "Robust multimedia spam filtering based on visual, textual, and audio deep features and random forest," *Multimedia Tools and Applications*, 2023.
- [283] T. Oorloff, S. Koppiseti, N. Bonettini, *et al.*, "AVFF: Audio-Visual Feature Fusion for Video Deepfake Detection," in *CVPR*, 2024.
- [284] D. Salvi, H. Liu, S. Mandelli, *et al.*, "A robust approach to multimodal deepfake detection," *Journal of Imaging*, 2023.
- [285] H. Ilyas, A. Javed, and K. M. Malik, "AVFakeNet: A unified end-to-end Dense Swin Transformer deep learning model for audio-visual deepfakes detection," *Applied Soft Computing*, 2023.
- [286] S. Asha, P. Vinod, and V. G. Menon, "A defensive attention mechanism to detect deepfake content across multiple modalities," *Multimedia Systems*, 2024.
- [287] X. Liu, Y. Yu, X. Li, *et al.*, "Magnifying multimodal forgery clues for deepfake detection," *Signal Processing: Image Communication*, 2023.
- [288] C. Feng, Z. Chen, and A. Owens, "Self-supervised video forensics by audio-visual anomaly detection," in *CVPR*, 2023.
- [289] H. Zou, M. Shen, Y. Hu, *et al.*, "Cross-modality and within-modality regularization for audio-visual deepfake detection," in *ICASSP*, 2024.
- [290] F. Nie, J. Ni, J. Zhang, *et al.*, "FRADE: Forgery-aware Audio-distilled Multimodal Learning for Deepfake Detection," in *ACMMM*, 2024.
- [291] Y. Zhang, W. Lin, and J. Xu, "Joint audio-visual attention with contrastive learning for more general deepfake detection," *ACM Transactions on Multimedia Computing, Communications and Applications*, 2024.
- [292] H. Dang, F. Liu, J. Stehouwer, *et al.*, "On the detection of digital face manipulation," in *CVPR*, 2020.
- [293] R. Wang, F. Juefei-Xu, L. Ma, *et al.*, "FakeSpotter: A Simple yet Robust Baseline for Spotting AI-Synthesized Fake Faces," in *IJCAI*, 2020.
- [294] Y. He, B. Gan, S. Chen, *et al.*, "ForgeryNet: A Versatile Benchmark for Comprehensive Forgery Analysis," in *CVPR*, 2021.
- [295] Z. Chen, K. Sun, Z. Zhou, *et al.*, "DiffusionFace: Towards a Comprehensive Dataset for Diffusion-Based Face Forgery Analysis," *arXiv:2403.18471*, 2024.
- [296] L. Verdoliva, C. Riess, J. Thies, *et al.*, "FaceForensics++: Learning to Detect Manipulated Facial Images," in *ICCV*, 2019.
- [297] L. Jiang, R. Li, W. Wu, *et al.*, "DeeperForensics-1.0: A Large-Scale Dataset for Real-World Face Forgery Detection," in *CVPR*, 2020.

- [298] Y. Li, X. Yang, P. Sun, *et al.*, "Celeb-DF: A Large-scale Challenging Dataset for DeepFake Forensics," in *CVPR*, 2020.
- [299] B. Zi, M. Chang, J. Chen, *et al.*, "WildDeepfake: A Challenging Real-World Dataset for Deepfake Detection," in *ACMMM*, 2020.
- [300] P. Korshunov and S. Marcel, "DeepFakes: a New Threat to Face Recognition? Assessment and Detection," *arXiv:1812.08685*, 2018.
- [301] X. Yang, Y. Li, and S. Lyu, "Exposing deep fakes using inconsistent head poses," in *ICASSP*, 2019.
- [302] H. Chen, Y. Hong, Z. Huang, *et al.*, "DeMamba: AI-Generated Video Detection on Million-Scale GenVideo Benchmark," *arXiv:2405.19707*, 2024.
- [303] J. Frank and L. Schönherr, "WaveFake: A Data Set to Facilitate Audio Deepfake Detection," in *NeurIPS*, 2021.
- [304] X. Wang, J. Yamagishi, M. Todisco, *et al.*, "ASVspoof 2019: A large-scale public database of synthesized, converted and re-played speech," *Computer Speech & Language*, 2020.
- [305] J. Yamagishi, X. Wang, M. Todisco, *et al.*, "ASVspoof 2021: accelerating progress in spoofed and deepfake speech detection," in *ASVSPPOOF*, 2021.
- [306] N. Müller, P. Czempin, F. Dieckmann, *et al.*, "Does audio deepfake detection generalize?," *INTERSPEECH*, 2022.
- [307] J. Yi, R. Fu, J. Tao, *et al.*, "ADD 2022: The First Audio Deep Synthesis Detection Challenge," in *ICASSP*, 2022.
- [308] J. Yi, J. Tao, R. Fu, *et al.*, "ADD 2023: the Second Audio Deepfake Detection Challenge," *arXiv:2305.13774*, 2023.
- [309] R. Reimao and V. Tzerpos, "FoR: A Dataset for Synthetic Speech Detection," in *SpED*, 2019.
- [310] N. Müller, P. Kawa, W. Choong, *et al.*, "MLAAD: The Multi-Language Audio Anti-Spoofing Dataset," *IJCNN*, 2024.
- [311] H. Khalid, S. Tariq, M. Kim, *et al.*, "FakeAVCeleb: A novel audio-video multimodal deepfake dataset," in *NeurIPS*, 2021.
- [312] Z. Cai, K. Stefanov, A. Dhall, *et al.*, "Do You Really Mean That? Content Driven Audio-Visual Deepfake Dataset and Multimodal Method for Temporal Forgery Localization," in *DICTA*, 2022.
- [313] B. Dolhansky, J. Bitton, B. Pfau, *et al.*, "The DeepFake Detection Challenge (DFDC) Dataset," *arXiv:2006.07397*, 2020.
- [314] W. Yang, X. Zhou, Z. Chen, *et al.*, "AVoid-DF: Audio-Visual Joint Learning for Detecting Deepfake," *IEEE Transactions on Information Forensics and Security*, 2023.
- [315] B. Hosler, D. Salvi, A. Murray, *et al.*, "Do deepfakes feel emotions? a semantic approach to detecting deepfakes via emotional inconsistencies," in *CVPR*, 2021.
- [316] H. Tak, M. Todisco, X. Wang, *et al.*, "Automatic speaker verification spoofing and deepfake detection using wav2vec 2.0 and data augmentation," *arXiv:2202.12233*, 2022.
- [317] V. S. Katamneni and A. Rattani, "MIS-AVoiDD: Modality invariant and specific representation for audio-visual deepfake detection," in *ICMLA*, 2023.
- [318] Y. Yu, X. Liu, R. Ni, *et al.*, "PVASS-MDD: Predictive visual-audio alignment self-supervision for multimodal deepfake detection," *IEEE Transactions on Circuits and Systems for Video Technology*, 2023.
- [319] T. Liu, F. Chen, F. S., *et al.*, "AniTalker: Animate Vivid and Diverse Talking Faces through Identity-Decoupled Facial Motion Encoding," *arXiv:2405.03121*, 2024.
- [320] K. Cho, J. Lee, H. Yoon, *et al.*, "GaussianTalker: Real-Time Talking Head Synthesis with 3D Gaussian Splatting," in *ACMMM*, 2024.
- [321] Y. Zheng, H. Yang, T. Zhang, *et al.*, "General facial representation learning in a visual-linguistic manner," in *CVPR*, 2022.
- [322] Z. Zhang, L. Li, Y. Ding, *et al.*, "Flow-guided one-shot talking face generation with a high-resolution audio-visual dataset," in *CVPR*, 2021.
- [323] I. Demirsahin, O. Kjartansson, A. Gutkin, *et al.*, "Open-source multi-speaker corpora of the English accents in the British isles," in *LREC*, 2020.
- [324] Y. A. Li, C. Han, V. Raghavan, *et al.*, "StyleTTS 2: Towards Human-Level Text-to-Speech through Style Diffusion and Adversarial Training with Large Speech Language Models," in *NeurIPS*, 2023.
- [325] H. Wang, M. Yu, J. Hai, *et al.*, "SSR-Speech: Towards Stable, Safe and Robust Zero-shot Text-based Speech Editing and Synthesis," *arXiv:2409.07556*, 2024.
- [326] H. Zen, V. Dang, R. Clark, *et al.*, "LibriTTS: A Corpus Derived from LibriSpeech for Text-to-Speech," in *INTERSPEECH*, 2019.
- [327] T.-C. Wang, A. Mallya, and M.-Y. Liu, "One-shot free-view neural talking-head synthesis for video conferencing," in *CVPR*, 2021.
- [328] Z. Yan, Y. Zhang, X. Yuan, *et al.*, "DeepfakeBench: A Comprehensive Benchmark of Deepfake Detection," in *NeurIPS*, 2023.

[329] D. P. Kingma and M. Welling, "Auto-Encoding Variational Bayes," *arXiv*, 2013.

[330] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," in *NeurIPS*, 2020.



**Florinel-Alin Croitoru** is a PhD student at the University of Bucharest, Romania. In 2021, he obtained his masters degree in Artificial Intelligence with a thesis on action spotting in football videos. His domains of interest include machine learning, computer vision and deep learning. He published several studies in top-tier conferences and journals, such as CVPR, ECAI, TPAMI, IJCV, CVIU.



**Andrei Hiji** is a PhD student at the University of Bucharest, Romania. He obtained his master's degree in Information Security from the Military Technical Academy "Ferdinand I" in 2019. His research interests include machine learning, anomaly detection and cybersecurity.



**Vlad Hondru** is a PhD student at the University of Bucharest, Romania. He obtained his bachelor's degree from the University of Manchester in Mechatronic Engineering, then he graduated from Imperial College London, studying towards an MSc in Computing Science, with a Visual Computing and Robotics specialization, focusing on Artificial Intelligence.



**Nicolae-Cătălin Ristea** graduated as valedictorian from the Faculty of Electronics, Telecommunications and Information Technology, NUST Politehnica Bucharest, in 2019. He completed his PhD in 2024 at the same university. Nicolae is co-author of multiple papers accepted at top-tier conferences and journals, such as CVPR, WACV and TPAMI. His research interests include deep learning, computer vision, machine learning and signal processing.



**Paul Irofti** is an Associate Professor within the Computer Science Department of the Faculty of Mathematics and Computer Science at the University of Bucharest. He is the co-author of the book "Dictionary Learning Algorithms and Applications" (Springer 2018). He is PhD in Systems Engineering at the Politehnica University of Bucharest since 2016. His interests are anomaly detection, signal processing, numerical algorithms and optimization.



**Marius Popescu** is associate professor at the University of Bucharest, Romania. He defended his PhD in 2004. His domains of interest are: AI, ML, computational linguistics, computer vision. His achievements in these fields include methods that ranked 3rd in the NLI Shared Task of BEA-8, 4th in the FER Challenge of WREPL 2013, 2nd in the ADI Shared Task of VarDial 2016, 1st in the NLI Shared Task of BEA-12.



**Cristian Rusu** is associate professor within the Computer Science Department of the Faculty of Mathematics and Computer Science at the University of Bucharest. He received his PhD in Systems Engineering at the Politehnica University of Bucharest in 2012. His research interests include signal processing, numerical linear algebra, machine learning, and deep learning.



**Radu Ionescu** is full professor at the University of Bucharest, Romania. He completed his PhD at the University of Bucharest in 2013. His research interests include machine learning, computer vision, image processing, computational linguistics and medical imaging. He published over 140 articles at international venues, including CVPR, NeurIPS, ICCV, ACL, EMNLP, TPAMI, IJCV, CVIU.



**Fahad Khan** is a faculty member at MBZ University of AI (MBZUAI), UAE and Linköping University, Sweden. He received the M.Sc. degree in Intelligent Systems Design from Chalmers University of Technology, Sweden and a Ph.D. degree in Computer Vision from Autonomous University of Barcelona, Spain. His research interests include a wide range of topics within computer vision, such as object recognition, object detection, action recognition and visual tracking.



**Mubarak Shah** is the UCF Trustee chair professor and the founding director of the Center for Research in Computer Vision at the University of Central Florida (UCF). He is a fellow of the NAI, IEEE, AAAS, IAPR and SPIE. He is an editor of an international book series on video computing, was editor-in-chief of *Machine Vision and Applications* and an associate editor of *ACM Computing Surveys* and *IEEE TPAMI*.

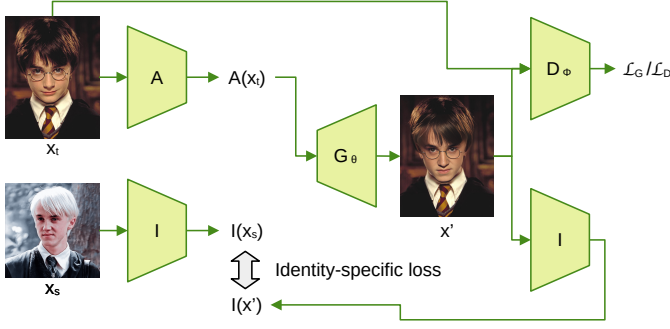


Fig. 4. An overview of face swapping based on GANs. The generative process is conditioned on an identity encoder  $I$  and an attribute encoder  $A$ , aiming to preserve the target attributes while replacing the target identity with the source identity.

## 8 SUPPLEMENTARY: DEEFAKE GENERATION TUTORIAL

Various deep generative models are actively being used to successfully generate deepfake content. Among the deepfake generative methods, we next explain in detail how the most popular and interesting developments work, such as generative adversarial networks, variational autoencoders, diffusion models, as well as Neural Radiance Fields.

### 8.1 Generative Adversarial Networks

Generative adversarial networks (GANs) [64] consist of two neural networks, called the generator and the discriminator. The generator, denoted as  $G_\theta(z)$ , transforms input Gaussian noise  $z \sim p(z)$  into a sample from the data distribution. The discriminator, represented as  $D_\phi(x)$ , outputs a single scalar value that predicts the probability of a given sample  $x$  to be real, rather than being generated by  $G_\theta$ . Therefore, the discriminator  $D_\phi$  is trained as a binary classifier, where the real samples are labeled with 1 and the generated ones with 0, as follows:

$$\min_{D_\phi} \mathcal{L}_D = -\mathbb{E}_{x \sim p(x)} [\log(D_\phi(x))] - \mathbb{E}_{z \sim p(z)} [\log(1 - D_\phi(G_\theta(z)))], \quad (1)$$

where  $p(x)$  and  $p(z)$  represent the real data distribution and the Gaussian distribution, respectively, with  $p(z)$  serving as a source for sampling varied inputs for the generator. The generator is trained to deceive the discriminator. Thus, its training objective maximizes the probability assigned by the discriminator to the generated samples  $D_\phi(G_\theta(z))$ :

$$\min_{G_\theta} \mathcal{L}_G = -\mathbb{E}_{z \sim p(z)} [\log(D_\phi(G_\theta(z)))]. \quad (2)$$

We can also express the whole training framework through a single mini-max optimization objective, as follows:

$$\min_{G_\theta} \max_{D_\phi} \mathbb{E}_{x \sim p(x)} [\log(D_\phi(x))] + \mathbb{E}_{z \sim p(z)} [\log(1 - D_\phi(G_\theta(z)))], \quad (3)$$

Note that optimizing for  $G_\theta$  does not influence the first term of the objective, as it depends only on  $D_\phi$ .

**Usage example.** In Figure 4, we illustrate a typical face swapping pipeline powered by GANs. The generator  $G_\theta$  is conditioned on features derived from two sources. The identity encoder  $I$  extracts features from the identity source image  $x_s$ , while the attribute encoder  $A$  extracts features from the target image  $x_t$ . Commonly, the encoder  $I$  is

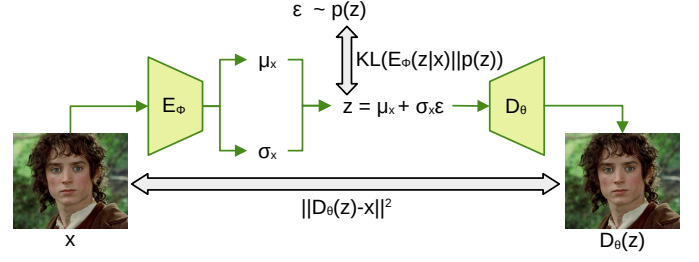


Fig. 5. An overview of face synthesis based on VAEs. The KL divergence is used to minimize the distribution gap between the distribution of  $z$  and the standard Gaussian distribution  $p(z)$ .

a (pre-trained) face recognition model, while  $A$  is a randomly initialized encoder trained along with the rest of the pipeline. The generator harnesses the input features to produce an image  $x'$  that retains the attributes of  $x_t$ , while swapping the target identity with the source identity from  $x_s$ . Aside from the previously discussed adversarial losses, the pipeline incorporates a reconstruction loss between  $x_t$  and  $x'$  to ensure attribute preservation. Additionally, an identity-specific loss is employed, typically calculated as the cosine similarity between the feature vectors extracted by the identity encoder  $I$  from  $x'$  and  $x_s$ .

### 8.2 Variational Autoencoders

A Variational Autoencoder (VAE) [329] is a modified version of the classic autoencoder, designed to support generative modeling by learning a probabilistic latent space. In the classic autoencoder case, during training, an encoder maps a sample  $x$  to a latent representation  $z$  from which a decoder is tasked to reconstruct the original input  $x$ . This approach is unsuitable for generative modeling because  $z$  follows an arbitrary complex distribution, making it impossible to directly sample from it. Kingma *et al.* [329] address this limitation by enforcing  $z$  to follow a standard Gaussian distribution. They achieved this by adding a Kullback-Leibler (KL) divergence regularization term to the loss function, which minimizes the divergence between the distribution of  $z$  and the standard Gaussian distribution. To model the distribution of  $z$ , the encoder outputs a mean  $\mu_x$  and a variance  $\sigma_x$  that describe a Gaussian distribution. During training  $z$  is sampled from this distribution. The loss function is as follows:

$$\mathcal{L} = \mathbb{E}_{z \sim E_\phi(z|x)} [\|D_\theta(z) - x\|_2^2] + \text{KL}(E_\phi(z|x) \| p(z)), \quad (4)$$

where  $E_\phi(z|x)$  denotes the encoder,  $D_\theta(x|z)$  is the decoder and  $p(z)$  represents the standard Gaussian distribution.

**Usage example.** In Figure 5, we showcase the training process of a VAE for face synthesis. To enable gradient backpropagation through the encoder, the reparameterization trick is applied to sample  $z \sim \mathcal{N}(\mu_x, \sigma_x \mathbf{I})$ . During inference,  $z$  is directly sampled from  $p(z)$  and fed into the decoder.

### 8.3 Diffusion Models

Diffusion models [1] consist of two diffusion processes, the forward diffusion process and the reverse diffusion process. In the forward process, a data sample is progressively transformed over  $T$  steps by adding Gaussian noise at each step, eventually converting it into an approximate standard Gaussian distribution. The reverse process operates in the



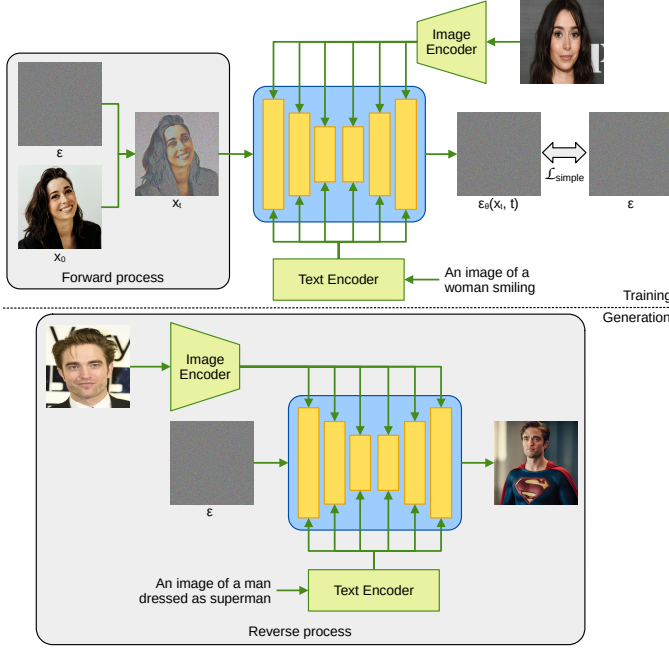


Fig. 6. An overview of text-conditional personalized generation based on diffusion models. The model aims to generate images of a source identity conditioned on a text prompt.

opposite direction, serving as the generative mechanism. It learns to map a standard Gaussian sample back to a data sample.

**Forward Process.** The forward process is a Markov chain  $x_t \sim q(x_t|x_{t-1}) = \mathcal{N}(\sqrt{1-\beta_t} \cdot x_{t-1}, \beta_t \cdot \mathbf{I})$ , which gradually adds Gaussian noise that depends on a variance schedule  $\{\beta_t\}_{t=1}^T$ , where  $x_0 \sim q(x_0)$  represents a data sample and  $x_T$  is approximately a standard Gaussian sample. This formulation supports an efficient sampling for an arbitrary  $x_t$  during training:

$$x_t \sim \mathcal{N}(\sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) \cdot \mathbf{I}), \quad (5)$$

where  $\alpha_t = 1 - \beta_t$ ,  $\bar{\alpha}_t = \prod_{i=0}^t \alpha_i$ .

**Reverse Process.** The reverse process starts with  $x_T \sim \mathcal{N}(0, \mathbf{I})$  and follows the learned Gaussian transitions denoted as  $p_\theta(x_{t-1}|x_t) = \mathcal{N}(\mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$  to recover a data sample  $x_0$ . Commonly, in practice, the variance  $\Sigma_\theta(x_t, t)$  is approximated with  $\beta_t$ . Thus, the only learnable component that remains is the mean  $\mu_\theta(x_t, t)$ , which can be rewritten as a function of noise [330]:

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\bar{\alpha}_t}} \left( x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \cdot \epsilon_\theta(x_t, t) \right). \quad (6)$$

As a consequence, during training, the neural network, denoted by  $\epsilon_\theta(x_t, t)$ , learns to approximate the noise  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$  added at arbitrary steps  $t \sim \mathcal{U}(1, \dots, T)$  to the data samples  $x_0 \sim q(x_0)$ . Formally, the optimization objective for the reverse process is defined as:

$$\min_{\epsilon_\theta} \mathcal{L}_{simple} = \mathbb{E}_{x_0 \sim q(x_0), t \sim \mathcal{U}(1, \dots, T), \epsilon \sim \mathcal{N}(0, \mathbf{I})} \|\epsilon - \epsilon_\theta(x_t, t)\|_2^2. \quad (7)$$

**Usage example.** In Figure 6, we illustrate the training and inference processes of a text-conditional personalized generative pipeline for a given identity. During training, pairs of images representing the same identity are used. One image undergoes the forward process, and the model is tasked with predicting the added noise. The second image in the

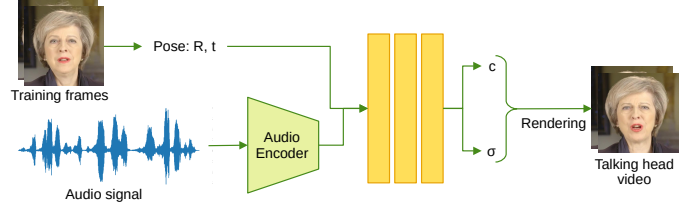


Fig. 7. An overview of talking-head synthesis based on NeRFs. The model learns to predict color and density values, which further enable the rendering of the deepfake talking head video.

pair, along with a text description of the first image, serve as conditional inputs to guide the model. The generative (reverse) process begins with standard Gaussian noise and progressively generates an image that aligns with the provided text description, while preserving the identity from the given source image.

## 8.4 Neural Radiance Fields

Neural Radiance Fields (NeRFs) [151] represent a method for synthesizing novel views from a sparse set of input images of a scene. The core idea is to train a single neural network to overfit to a specific scene, with the weights encoding the detailed information and structure of the scene. The input to the neural network is a 5D vector consisting of three spatial coordinates  $(x, y, z)$ , which represent a point in 3D space of the scene, and two angles  $(\theta, \phi)$ , which define the viewing direction. The network outputs the density  $\sigma$  at the given spatial point and its color  $c$ , represented as an RGB vector. However, datasets typically lack direct annotations for the density and color of each spatial point, so the network optimization is performed in the pixel space instead. This process involves selecting a viewing direction and casting a ray through the scene, along which multiple spatial points are sampled. These points are processed by the neural network to compute their corresponding density and color values. The outputs are then combined using volume rendering to generate an approximate pixel value. This predicted pixel value is compared with the ground-truth pixel using a mean squared error loss function to optimize the network. Formally, the loss function is defined as:

$$\mathcal{L} = \sum_{r \sim \mathcal{R}} \|\hat{C}(r) - C(r)\|_2^2, \quad (8)$$

where  $\mathcal{R}$  is the set of rays in the current batch,  $\hat{C}(r)$  is the predicted color for a particular ray  $r$  and  $C(r)$  is the ground-truth color.

In practice, NeRFs have two implementation details for a better optimization. First, spatial positions are projected into a higher-dimensional space using a positional encoding, similar to the approach used for vision transformers. Second, NeRFs use Hierarchical Volume Sampling, which involves training two neural networks. The first network samples coarse spatial points along a ray to estimate the overall structure of the scene. It then identifies regions of greater importance. The second network focuses on these regions, performing finer sampling to achieve more accurate results.

**Usage example.** NeRFs are used for talking head synthesis, as illustrated in Figure 7. The main idea is to use the audio as the driving signal. More precisely, along with the information about the pose, the neural network processes audio

features from a pre-trained encoder to output the color and density values from which the rendering is possible.