# Deep Person Generation: A Survey from the Perspective of Face, Pose and Cloth Synthesis

TONG SHA, Beihang University, China

WEI ZHANG, JD AI Research, China

TONG SHEN, JD AI Research, China

ZHOUJUN LI, Beihang University, China

TAO MEI, JD AI Research, China

Deep person generation has attracted extensive research attention due to its wide applications in virtual agents, video conferencing, online shopping and art/movie production. With the advancement of deep learning, visual appearances (face, pose, cloth) of a person image can be easily generated on demand. In this survey, we first summarize the scope of person generation, and then systematically review recent progress and technical trends in identity-preserving deep person generation, covering three major tasks: *talking-head generation* (face), *pose-guided person generation* (pose) and *garment-oriented person generation* (cloth). More than two hundred papers are covered for a thorough overview, and the milestone works are highlighted to witness the major technical breakthrough. Based on these fundamental tasks, many applications are investigated, e.g., virtual fitting, digital human, generative data augmentation. We hope this survey could shed some light on the future prospects of identity-preserving deep person generation, and provide a helpful foundation for full applications towards the digital human.

CCS Concepts: • **Computing methodologies → Computer vision**; **Image manipulation**; **Image-based rendering**.

Additional Key Words and Phrases: Deep Person Generation; Talking-head Generation; Pose-guided Person Generation; Garment-oriented Person Generation; Virtual Try-on; Generative Adversarial Networks; Digital Human

## 1 INTRODUCTION

With the advancement of deep learning, people are no longer satisfied with the visual understanding of camera-taken photos/videos. Visual content generation emerges as another research direction since images and videos are much more efficient for information presentation and exchange. Person generation is to synthesize person images and videos as realistically as possible. The long-term goal is to generate digital humans with life-like appearances, expressions and behaviors as real persons.

Authors' addresses: Tong Sha, tongsha@buaa.edu.cn, Beihang University, Beijing, China; Wei Zhang, wzhang.cu@gmail.com, JD AI Research, Beijing, China; Tong Shen, tshen.st@outlook.com, JD AI Research, Beijing, China; Zhoujun Li, lizj@buaa.edu.cn, Beihang University, Beijing, China; Tao Mei, tmei@live.com, JD AI Research, Beijing, China.
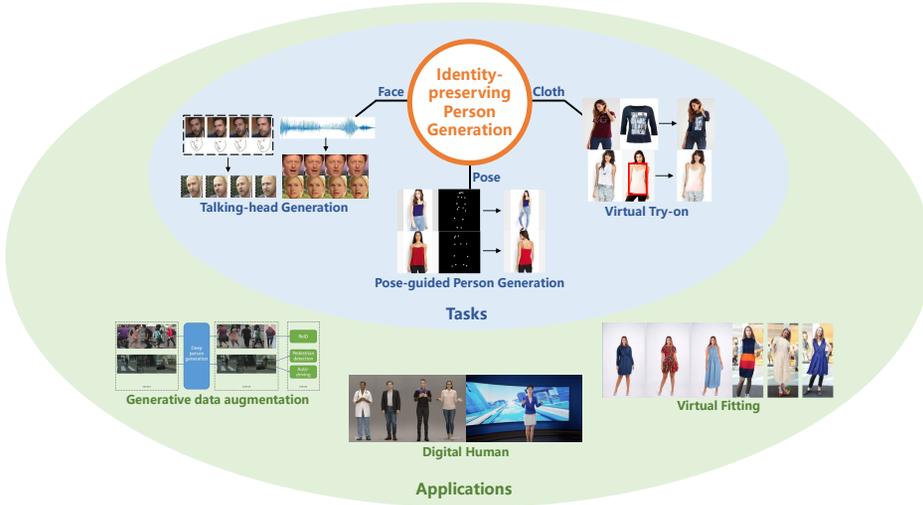
**111**

Fig. 1. The scope of identity-preserving person generation in this paper. From the three key components (face, pose, garment), we choose three mainstream tasks in recent research, namely talking-head generation, pose-guided person generation and virtual try-on.



Fig. 2. An overview of recent works on talking-head generation. Along the time axis, the number of works is increasing rapidly. For audio-driven generation, the trending technique is mixed with 2D and 3D methods. For motion-driven generation, 3D methods are developed earlier, but 2D methods are surging recently, due to the advancement of GAN (Generative Adversarial Network) [52].

As an emerging area, person generation has attracted lots of research attention, due to its wide applications on digital human[1], customer service[2], telepresence, art/fashion design and Metaverse.

For person images, identity, motion and appearance information are three important components. Identity refers to attributes used to recognize the identity of a person, including but not limited to the face, body shape and height (gait). Among them, face serves as the most commonly adopted media. Motion refers to the body gestures and facial expressions conveyed by a person. Appearance refers to the garments and accessories worn by a person. In some cases, non-identity information such as

---

[1]https://www.neon.life/

[2]https://digitalhumans.com/

Fig. 3. An overview of recent works on pose-guided and garment-oriented person generation since 2017. The left and right halves represent pose and garment transfer, respectively. The upper and lower halves denote top-down and bottom-up methods, respectively. The concentric dotted ellipses (inner to outer) are the time axis.

hairstyle and beard style can also be included as part of appearance information. Generally, in many practical applications such as film shooting, fashion performance and social robots, the demand for motion and appearance generation are highly favored. In this survey, we particularly focus on the identity-preserving person generation, where new face/body dynamics are synthesized while preserving the identity information. Importantly, this task is essential to emerging applications such as digital human, social robots.

Generally speaking, there are three key components (face, pose, cloth) to synthesize a person image or video. We choose the mainstream task in recent research from each component, as shown in Fig. 1. *Face generation* is the most popular area and has been extensively studied for years. Facial GANs (Generative Adversarial Networks), including unconditional face generation, facial attributes editing and Deepfakes, are already well reviewed in previous paper [77, 81, 167, 173]. Meanwhile, with the easily accessible tools (e.g., DeepFaceLab [126]), forged facial videos also pose severe social and ethical problems, such as the threats of Deepfakes in the Presidential election. Recently, the identity-preserving facial generation quickly gains popularity, and the most representative branch is the talking-head generation. It has also drawn much attention due to various applications such as video conferencing and customer service. *Pose-guided person generation*[3] is another popular task, where a new person image is generated given a conditioning pose. This task is essential to a number of motion-aware generation tasks, such as dancing and sports synthesis. *Garment-oriented*

---

[3]Also known as "Pose Transfer" or "Gesture-to-Gesture Translation" in other literature.

*generation*[4] is to synthesize new clothes based on conditioning inputs. It is fundamentally important for virtual try-on and other garment manipulation tasks, i.e., text-guided garment manipulation, garment inpainting.

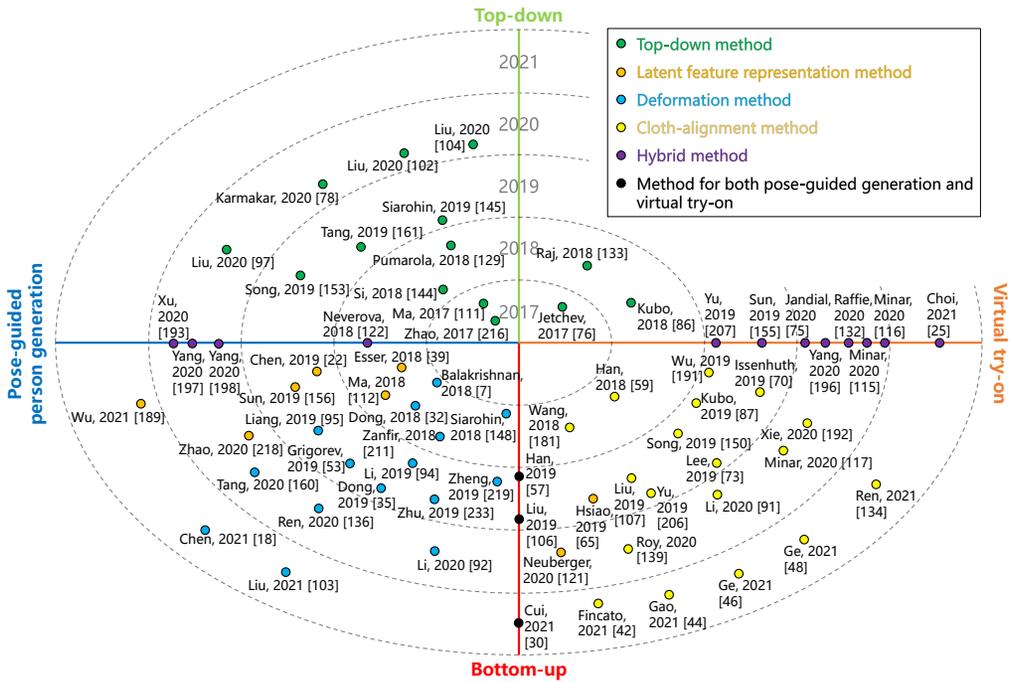Fig. 2 shows the literature map for talking-head generation. Along the time axis, the number of works has increased sharply in recent years. Roughly two branches of techniques are developed, depending on the driving signals, i.e., audio or motion. For both directions, 2D solutions are developed relatively earlier, and thus draw more research attention so far. Meanwhile, 3D-based methods, as an important complement, are also extensively explored recently. More details are covered in Section 2.

For pose and garment generation, their technical routines roughly follow a similar pattern. Fig. 3 plots an overview of existing literature. Overall, the following observations can be clearly identified. (1) The number of works surged rapidly over years (see the time axis). (2) In general, more bottom-up methods are favored over top-down ones. (3) For both pose-oriented and garment-oriented generation, the trend is first from top-down to bottom-up, and then to hybrid methods. (4) The most popular methods for pose and garment transfer are the deformation methods and cloth-alignment methods, respectively. Details are discussed in Section 3 (pose) and Section 4 (garment).

Despite the rapid development and rich literature, there is no systematic survey for identity-preserving person generation. In this paper, we comprehensively review techniques of identity-preserving person generation in terms of face, pose and garment synthesis. Similar works are either focusing on a small sub-area or addressing the topic from a different angle. Liu *et al.* [105] focus on generic image and video synthesis, while Kammoun *et al.* [77] highlight face-oriented generation and reenactment. Ruben *et al.* [167] review facial attribute editing and deepfake techniques. Cheng *et al.* [24] conduct a fashion-related survey covering style transfer and pose transformation. Ghodhbani *et al.* [50] review the image-based virtual try-on from the angle of fashion. However, these works are all with different focuses in partially overlapping domains, and there is no systematic survey with a broad but focused view on identity-preserving person generation and its applications. Generally, our work has the following contributions:

- We provide a systematic survey of identity-preserving person generation from the face, pose, and garment synthesis. To the best of our knowledge, this is the first review of person generation of this kind.
- From the three key components (face, pose, garment), we choose three mainstream tasks in recent research, namely *talking-head generation*, *pose-guided person generation* and *virtual try-on*. We provide a comprehensive and in-depth review of state-of-the-art methods. Meanwhile, we summarize the common points of these three tasks from the viewpoint of decomposition.
- Mainstream benchmarks and metrics are summarized. Meanwhile, we summarize the main applications: generative data augmentation, virtual fitting and digital human.
- We list some possible future directions, to inspire researchers related to person generation.

The remaining parts of this survey are organized as follows. Section 2 reviews talking-head generation driven by audio or motion inputs. Section 3 discusses pose-guided person image/video generation. Section 4 summarizes garment-oriented person generation, namely virtual try-on and other garment manipulation tasks. Section 5 summarizes popular benchmarks and metrics used in person generation. Section 6 concludes the common points of three tasks. Section 7 illustrates major applications of deep person generation. Section 8 discusses possible future directions worth further exploration.

---

[4]Also known as "Garment Transfer" or "Appearance Transfer" in some works.

(a) Motion-driven talking-head generation          (b) Audio-driven talking-head generation
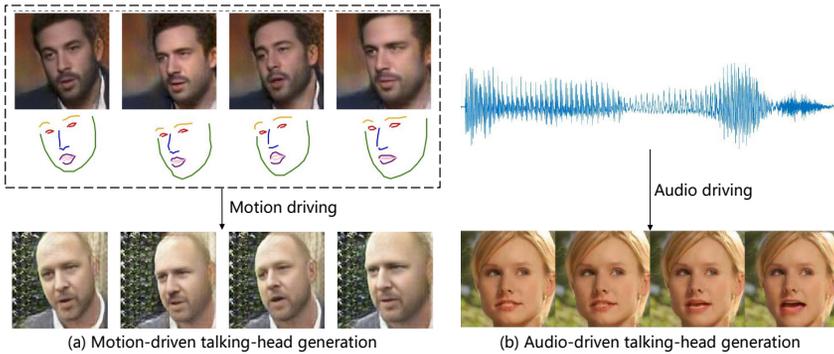
Fig. 4. Illustration of motion-driven (a) and audio-driven (b) talking-head generation.

## 2 TALKING-HEAD GENERATION

Talking-head generation aims to synthesize a person talking image or video, driven by motion, audio or text, which is an important branch of dynamic face generation. As a basis for the subsequent content, we briefly summarise the classic methods on facial image generation.

Goodfellow *et al.* [52] propose the Generative Adversarial Network (GAN), where the idea of game adversarial is adopted in network training. Deep Convolutional GAN (DCGAN) [131] introduces deep convolutions in GAN, which stabilizes the training process. These methods can be directly adopted for unconditional face generation. For conditional face generation, Conditional GAN (CGAN) [118] introduces additional inputs to guide the generator for image-to-image translation [60, 123] and face reenactment [49, 190]. Furthermore, controllable GANs are introduced to control specific attributes of faces. For example, StyleGAN [80] uses style codes to control the overall style of the facial image. These methods serve as the foundation for the subsequent talking-head generation.

Talking-head generation is crucial for several applications, including video conferencing, virtual anchors and customer services. As shown in Fig. 4, these methods can generally be divided into two categories, motion-driven and audio-driven, depending on the driving signal. Note that the "text driven" (text-to-video) generation is a natural extension of "audio-driven" (speech-to-video) since the only difference is the well-developed text-to-speech technique [89, 93]. Therefore in this survey, we consider the text-driven branch as a variant of the audio-driven problem. Tab. 1 summarizes the representative works on talking-head generation.

### 2.1 Motion-driven Talking-Head Generation

Motion-driven branch adopts the driving factor of motions, in which the motion includes two parts: head pose and facial expression. This task is also known as "face reenactment", which usually manipulates head pose and facial expression simultaneously. There are also methods to manipulate only facial expressions, known as "Expression Swap". Motion-driven talking-head generation has many applications in telepresence (e.g., video conferencing, multiplayer online games). Technically, two lines of research can be identified: 2D-based and 3D-based methods, depending on their internal representation of the head model.

*2.1.1 2D-based Methods.* Existing 2D-based works can be roughly grouped into three categories, depending on their intermediate facial representation, i.e., facial landmarks, latent features and action units.

Table 1. Summary of talking-head generation methods. ID D / I: the method is "Identity-dependent" or "Identity-independent". D: Dependent. I: Independent. H: Hybrid. 3D model: the method use the 3D head model or not. Change pose: the method changes the head poses of source images/videos or not.

| References | Key idea | Driving factor | ID D/I | 3D model | Change pose |
|---|---|---|---|---|---|
| Bansal [8] | Cycle-consistency video-to-video generation | Video | D | × | √ |
| [6, 49] | Landmarks-guided warping and detail refinement | Image or video | I | × | √ |
| Wu [190] | Video to landmarks to video | Video | H | × | √ |
| Zakharov [209, 210] | Landmark-driven few-shot adversarial learning | Image or video | I | × | √ |
| Gu [54] | Landmark-driven GAN with warping and appearance streams | Image or video | I | × | √ |
| [60, 123] | Landmark-driven GAN | Image or video | I | × | √ |
| Wang [184] | Landmark-driven few-shot video-to-video | Video | I | × | √ |
| Chen [23] | Landmarks to semantic map to result | Image or video | I | × | √ |
| Tripathy [168] | Action units based face reenactment | Image or video | I | × | √ |
| Tripathy [169] | Action units based landmarks transformer | Image or video | I | × | √ |
| [12, 56, 188, 212] | Identity and pose features extraction and fusion | Image or video | I | × | √ |
| Siarohin [147] | Flow warping based video-to-video generation | Video | I | × | √ |
| [119, 124] | 3D reconstruction based GAN | Video | I | √ | × |
| Thies [164] | 3D reconstruction and parameter replacement | Video | I | √ | × |
| Kim [83] | 3D reconstruction and parameter replacement | Video | I | √ | √ |
| Kim [82] | 3D reconstruction and parameter replacement | Video | I | √ | × |
| Koujan [36, 85] | 3D reconstruction and parameter replacement | Video | I | √ | √ |
| [151, 186] | 3D Keypoints extraction and flow warping | Video | I | √ | √ |
| [20, 74] | Audio and identity features extraction and fusion | Audio | I | × | × |
| Zhou [225] | Person and word features extraction and fusion | Video or audio | I | × | × |
| Zhou [226] | Contrastive-learning based feature extraction | Video and audio | I | × | × |
| Chen [21] | Audio to landmarks to video | Audio | I | × | × |
| Yu [205] | Audio & text to landmarks to video | Audio and text | I | × | × |
| Vougioukas [177] | Aduio-driven GAN | Audio | I | × | × |
| [154, 178] | RNN-based GAN with three discriminators | Audio | I | × | × |
| [38, 140] | Audio & emotion driven GAN | Audio | I | × | × |
| Eskimez [37] | RNN-based GAN | Audio | I | × | × |
| Zhou [229] | Audio-driven landmark prediction | Audio | I | × | √ |
| Prajwal [128] | Wav2Lip: GAN + pre-trained lip-sync expert | Audio | I | × | √ |
| Kumar [89] | Text to audio to keypoints to video | Text | D | × | × |
| Suwajanakorn [158] | Audio to shapes to mouth images to video | Audio | D | √ | × |
| Karras [79] | Audio & emotion to 3D model | Audio | D | √ | × |
| Cudeiro [29] | VOCA: speech to 3D model network | Audio | I | √ | × |
| Thies [163] | Audio to expressions to 3D model to video | Audio | H | √ | × |
| Song [152] | 3D reconstruction and parameter replacement | Audio | I | √ | × |
| Yi [200] | 3D reconstruction and parameter replacement | Audio | I | √ | × |
| Guo [55] | Audio-driven neural radiance fields | Audio | D | √ | √ |
| Fried [43] | 3D reconstruction and parameter recombination | Text | D | √ | × |
| Li [93] | Text-driven 3D parameter generation | Text | D | √ | √ |

Facial landmark driving methods [6, 23, 49, 54, 60, 123, 147, 169, 184, 190, 209, 210] adopt explicit facial landmarks to encode motions. Averbuch-Elor *et al.* [6] and Geng *et al.* [49] warp faces based on landmarks for a coarse result, and adopt a generative adversarial network for refinement. Furthermore, Gu *et al.* [54], Hao *et al.* [60] and Nirkin *et al.* [123] warp faces on the feature-map level, also based on landmarks. Chen *et al.* [23] introduce facial semantic maps to further improve the visual quality. Siarohin *et al.* [147] estimate flow between the source and driving frames, and then warp source image to target frames based on these flows. However, facial landmarks contain additional identity information. Some works improve the methods based on this point. Zakharov *et al.* [209, 210] directly map facial landmarks to images, modulated by the identity features. Besides, their few-shot adversarial learning enables training with even fewer images. Wu *et al.* [190] adopt a person-specific transformer to warp landmarks into a specific identity.

Latent feature driving methods [12, 56, 188, 212, 225, 226] are to extract the identity feature from the source video and the motion feature from the driving video, and then fuse them together to generate the target video. The identity feature could be explicit. X2Face [188] adopts a neutral face

image as an explicit identity representation and then warps this neutral face with driving motions. Meanwhile, some works extract implicit identity features. Ha *et al.* [56] design an attention block for effective extraction of identity features. Zhou *et al.* [226] use a contrastive learning strategy to decompose the identity feature and non-identity feature (pose and facial movements) from talking-head videos.

Action units are the coding systems for describing facial expressions. Some methods [168, 169] use action units to represent identity-independent facial motions. Tripathy et al. [168] use the source image and action units extracted from the driving image to generate the result. This method has a good decoupling of identity and motion, but the precision is limited. They also propose FACEGAN [169] to transform landmarks using the action units. The transformed landmarks replace the driving landmarks as the intermediate facial representation. This method ensures that the generated results do not carry the driving face identity information. Meanwhile, it controls the motion more accurately.

*2.1.2 3D-based Methods.* Different from the 2D branch, 3D methods are mainly based on 3D face modeling. Vlasic *et al.* [175] and Garrido *et al.* [45] use dubber video to guide lip motions based on 3D face models. Thies *et al.* [165, 166] adopt video face tracking to extract 3D face models, and then apply expression transfer for lip motion generation. Note that Face2Face [166] later becomes the prototype for many subsequent 3D methods.

The above methods need to build special 3D face models manually in advance. Recent works [36, 82, 83, 85, 119, 124, 164] are mostly based on monocular 3D reconstruction, following a similar schema as [166]. First, monocular 3D reconstruction is adopted to obtain face model parameters for both the source and driving videos. Then, these 3D model parameters are combined to generate the target 3D model. Finally, a video rendering module is applied for video generation. Kim *et al.* [82, 83] modify the pose, expression, and eye parameters for fully controllable faces, or modify only the expression parameters to preserve styles. Koujan *et al.* [36, 85] preserve the scale parameter of the source video, to guarantee proper sizes of generated faces.

Different from the above methods, some works [151, 186] utilize 3D facial landmarks to control facial motions. Wang *et al.* [186] propose to extract 3D feature points from videos, instead of using 3D reconstruction. Furthermore, Song *et al.* [151] animate illusory faces with control points extracted from talking video frames.

*2.1.3 Facial Representation Comparison.* The methods summarized in previous sections use many intermediate facial representations. There are five main intermediate facial representations: 2D-based facial landmarks, latent features, action units, 3D-based face model parameters and 3D landmarks. These forms have their own advantages and disadvantages. For example, Tripathy *et al.* [169] analyze the pros and cons for landmarks and Wang *et al.* [186] for the 3D face model. We systematically summarize these five intermediate facial representations in Tab. 2.

## 2.2 Audio-driven Talking-Head Generation

Audio-driven talking-head generation is to synthesize a realistic animated video driven by a piece of audio (e.g., speech, singing audio). Recently, this area has been blooming with extensive research attention. Similarly, there are two streams of research, 2D-based and 3D-based methods.

*2.2.1 2D-based Methods.* 2D methods mostly adopt landmarks, semantic maps, or other image-like representations during synthesis, which dated back to Bregler *et al.* 1997 [11]. Early works [11, 16, 17, 40, 99, 180] mainly use traditional learning methods, such as Hidden Markov Model (HMM) [11, 180], LSTM [40] and frame retrieval [17]. Due to the restrictions on method, hardware

Table 2. Comparison of intermediate facial representations

| Facial representation | Demonstration | Advantages | Disadvantages | Reference |
|---|---|---|---|---|
| Facial landmarks |  Face Image  Landmarks | • Strong interpretability and controllability.<br>• Easy to obtain.<br>• The training model has high generalization. | • There is some identity information as a distraction, including the contour shape of the face.<br>• The pose and expression cannot be decoupled.<br>• Poor expression of some details. For example, it is difficult to convey the subtle expression. | [6, 23, 210] |
| Latent feature |  | • Have stronger decomposition than facial landmarks.<br>• Can handle some more delicate decomposition, such as the decomposition between pose and expression. | • The interpretability is poor, especially the learning based latent feature.<br>• It is difficult to ensure that the features have sufficient decomposition. The motion features inevitably contain some identity information.<br>• The generalization of training model is generally poor.<br>• Some attribute information will be lost inevitably, resulting in slightly inaccurate or fuzzy results. | [188, 225, 226] |
| Action units |  | • Have strong decomposition and almost no identity information.<br>• Have a strong ability to express expression and can express the subtle differences between expressions. | • Poor interpretability and controllability making the low quality results. | [168, 169] |
| 3D model parameters |  Face Image  3DMM | • Strong interpretability and controllability.<br>• Have strong decomposition. Identity and motion can be well decoupled, and the pose and expression can also be decoupled.<br>• Results can be generated from almost any view. | • It is difficult to obtain 3D parameters. Even with mature monocular 3D reconstruction technology, it takes a certain amount of time to obtain parameters.<br>• The training model does not have strong generalization and is usually applied to a single identity or a small range of identities.<br>• Poor ability to express some details, including appearance details and micro-expressions. | [83, 85, 164] |
| 3D landmarks |  Face Image  3D Keypoints | • Relatively easy to obtain.<br>• Compared with the learning based latent feature, it is more explanatory.<br>• Strong decomposition, basically not affected by identity information.<br>• The training model is highly generalized. | • Only suitable for local range of views. | [151, 186] |

and data collection, these works apply only to specific identities, and their results are rather preliminary.

Since 2017, GAN-like approaches gradually become popular due to their superior visual quality and strong generalization for identities. Person identities (identity-independent) become popular. Identity-independent methods are mainly divided into two categories, latent feature based and facial landmarks based.

Latent feature based approaches extract audio and identity features with two encoders and fuse them to generate talking heads with new lip motions. Jamaludin *et al.* [74] propose the first identity-independent architecture by disentangling audio and identity features with two encoders. However, this method does not take into account lip synchronization and video temporal coherence. To solve these problems, Song *et al.* [154] and Vougioukas *et al.* [178] both propose multiple discriminators: frame (image-level fidelity), sequence (video-level fidelity) and syncing (lip reading) discriminators, to improve the visual quality and temporal coherence in video generation.

Facial landmarks based approaches use explicit facial landmarks to connect the audio with lip motions. Chen *et al.* [21] and Yu *et al.* [205] utilize facial landmarks as the internal representation to bridge audio and facial images. Furthermore, Sadoughi *et al.* [140] and Eskimez *et al.* [38] introduce

emotions as another conditioning input for diverse facial expressions. Zhou *et al.* [229] further apply a speaker-aware animation model to predict spontaneous head poses alongside the audio. To further improve lip synchronization, Prajwal *et al.* [128] introduce a pre-trained SyncNet as lip-sync discriminator.

Some methods use text as the driving factor. Kumar *et al.* [89] use Char2Wav to transfer textual input into audio. Yu *et al.* [205] fuse text and audio together to generate corresponding mouth landmarks.

*2.2.2  3D-based Methods.* Early 3D-based methods pre-build a 3D face model of a specific person, and then render faces based on the 3D model. Compared to 2D solutions, 3D methods are better at motion controlling, especially when synthesizing novel motions and views. But the drawback is the cost of delicate 3D model construction. Suwajanakorn *et al.* [158] and Karras *et al.* [79] both adopt a pre-built 3D face model, and then drive the model by learning a sequence (a piece of audio) to sequence (motions of the 3D model) mapping.

Recent methods tend to directly reconstruct the 3D face model out of the training images or videos. Thies *et al.* [163] design an identity-independent Audio2ExpressionNet and an identity-dependent 3D face construction model. Song *et al.* [152] replace the 3D expression parameters from the source video with those generated by the input audio. Yi *et al.* [200] improve over [152] by introducing pose parameters for generating natural head motions. To get more refined rendering results, Guo *et al.* [55] train two conditional Neural Radiance Fields (NeRFs) [114] to render the head and torso parts, respectively.

Some methods use text instead of audio to drive faces. Fried *et al.* [43] align phonemes with the 3DMM [10] parameters to generate audio out of input text. Li *et al.* [93] directly generate 3D head pose, upper face and mouth shape animation parameters according to the input text, skipping the process of generating audio.

## 2.3  Discussion

Facial expressions are becoming vital in recent 2D-based talking-head works. Fig. 5 shows the vividness difference among the results generated by different methods. As is shown in Fig. 5, early methods only generate and modify the mouth area [11, 74, 154, 225], resulting in limited liveliness. Recently, some methods [38, 229] start to explore rich emotions and spontaneous head motions. Especially, Rotger *et al.* [138] bring the 2D facial expression transfer problem into the 3D domain. Based on 3D triangle meshes, they achieved sufficiently detailed expression generation. In the near future, modifying micro-expressions could be possible to appear.

Compared with motion-driven ones, audio-driven tasks are more challenging due to the large domain gap and non-deterministic mapping. It is inherently difficult to establish cross-modal connections among different forms of data: audio, text and video. Furthermore, the same input audio might be interpreted as diverse head motions and facial expressions. In addition, human eyes are highly sensitive to visual artifacts in frames and temporal jitters in videos, making this problem even more challenging.

Either audio- or motion-driven talking-head generation is now leaning towards a solution combining 2D and 3D techniques. In general, GAN-based 2D solutions give more visually plausible results, but 3D solutions are better at motion control. Furthermore, recent advancements in monocular 3D reconstruction also provide convenient 3D models during generation. Regarding performance, this task is still far from mature. For example, generating high-resolution video with micro-expressions is still a major challenge to current methods.
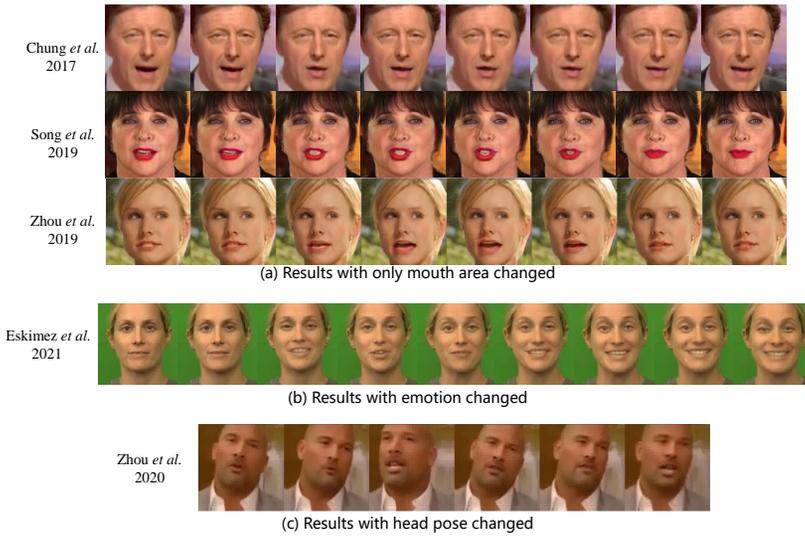
Fig. 5. The audio-driven talking-head generation results with (a) only mouth area changed, (b) emotion changed and (c) head pose changed. These results are from Chung *et al.* [74], Song *et al.* [154], Zhou *et al.* [225], Eskimez *et al.* [38] and Zhou *et al.* [229].



Fig. 6. Illustration of pose-guided person generation.

## 3 POSE-GUIDED PERSON GENERATION

Pose-guided person generation aims to generate full-body person images or videos guided by target poses as realistic as possible, as shown in Fig. 6. The fundamental difference between talking-head is the emphasis on body poses and motions. This task has many potential applications in movie production, online fashion shopping, etc. We focus on two major branches: pose-guided person image and video generation.

### 3.1 Pose-guided Person Image Generation

Pose-guided person image generation aims to transform a source person image according to a given target pose. The main challenge lies in the misalignment between the source and target poses.

The summary on pose-guided person image generation is shown in Tab. 3. Existing methods can be roughly grouped into three categories: top-down, bottom-up and hybrid methods. Concretely, top-down methods directly learn a mapping from the input to desired output image, mostly with a GAN-based network. Bottom-up approaches decompose the whole pipeline into several

Table 3. Summary of pose-guided person image generation works. Use parsing: use human parsing or not. Separate FG: Separate the foreground of images or not.

| References | Main idea | Use 3D pose | Use parsing | Separate FG |
|---|---|---|---|---|
| Top-down methods | | | | |
| [111, 216] | Coarse-to-fine generation | × | × | × |
| Si [144] | Direct generation | × | × | √ |
| Karmakar [78] | Direct generation | × | × | × |
| Tang [161] | Cycle consistency | × | × | × |
| Liu [102] | Semantic-guided , attention mechanism | × | √ | × |
| Liu [104] | Mask-guided generation | × | × | √ |
| Siarohin [145] | Multi-source generation, attention mechanism | × | × | × |
| Liu [97] | High-resolution progressive training | × | × | √ |
| Pumarola [129] | Cycle consistency | × | × | × |
| Song [153] | Cycle consistency | × | √ | × |
| Bottom-up latent feature representation methods | | | | |
| Esser [39] | Feature representation | × | × | × |
| Chen [22] | Cycle consistency and feature representation | × | × | × |
| Ma [112] | Feature representation | × | × | √ |
| Sun [156] | Multi-source feature representation | × | × | × |
| Zhao [218] | Pose serialization | × | × | × |
| Wu [189] | Feature representation | × | × | × |
| Bottom-up deformation methods | | | | |
| Siarohin [148] | Local deformation | × | × | × |
| Liang [95] | Local deformation | × | × | √ |
| Balakrishnan [7] | Local deformation | × | √ | √ |
| Dong [32] | Global deformation | × | √ | × |
| Zheng [219] | Flow warping | √ | √ | √ |
| Han [57] | Flow warping | × | √ | × |
| [53, 211] | Global deformation | √ | × | √ |
| [30, 94] | Flow warping | × | × | √ |
| Liu [106] | Flow warping | √ | × | √ |
| Dong [35] | Decompose human and pose into parts | × | √ | × |
| [92, 136, 160, 233] | Attention mechanism | × | × | × |
| Chen [18] | Multi-attention mechanism | × | × | × |
| Liu [103] | Semantic parsing attention | × | √ | × |
| Hybrid methods | | | | |
| Neverova [122] | Direct generation and global deformation | √ | × | × |
| Yang [197, 198] | Attention mechanism and refinement | × | × | × |
| Xu [193] | Deformation and refinement in a non-iconic view | × | × | √ |

intermediate components and build up the final result step by step. Hybrid approaches take advantage of both sides. As an overview, Fig. 7 shows the representative works from each category, and Fig. 8 compares the main network structures.

*3.1.1 Top-down Methods.* Inspired by the GAN-based view synthesis [216], Ma *et al.* [111] apply conditional GAN [118] for pose-guided person generation, and propose the pioneer top-down solution PG$^2$. PG$^2$ directly concatenates the source image and target pose as input to synthesize the target image, plus a refinement model to improve the generation details.

The results generated by PG$^2$ still lose many appearance details. For better performance, several improvements are proposed subsequently, including data augmentation [78] and more informative inputs [102, 104, 145]. Specifically, Siarohin *et al.* [145] propose a multi-source method, where an attention mechanism is adopted to assign weights to different source images. Liu *et al.* [102, 104] replace poses with segmentation masks and parsing maps to guide image generation.

Besides the aforementioned supervised branch, some works explore the unsupervised setting. Pumarola *et al.* [129] introduce the first unsupervised method, which uses the generated target person image with the source pose to regenerate the source image similar to CycleGAN [231]. Song *et al.* [153] operates on the semantic maps to bypass the requirement of paired data and generate more
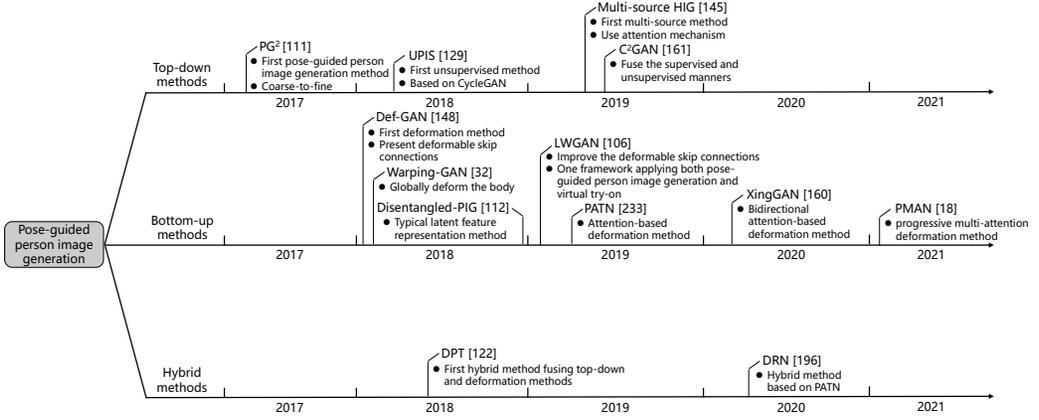
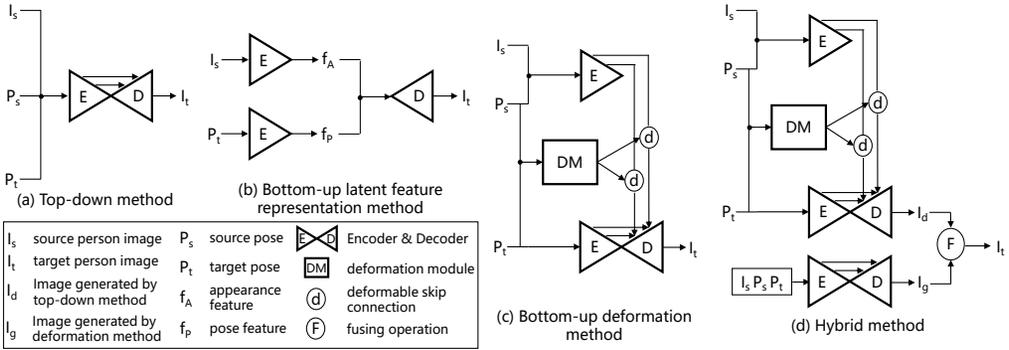Fig. 7. Representative works for pose-guided image person generation.



Fig. 8. Four typical network structures for pose-guided person image generation. (a) Top-down methods directly concatenate $I_s$, $P_s$ and $P_t$ as input to synthesize $I_t$. (b) Bottom-up latent feature representation methods extract the source appearance feature $f_A$ and target pose feature $f_P$ as internal representations. (c) Bottom-up deformation methods try to deform the feature maps of $I_s$ and $P_s$ into the target pose $P_t$. (d) Hybrid methods fuse the top-down and bottom-up deformation methods into a framework.

appearance details. Tang *et al.* [161] introduce multiple cyclic losses, i.e., $1 \times$ image→image→image cycle and $2 \times$ pose→image→pose cycles, to introduce more self-supervisions.

Top-down methods can not well preserve the appearance details of the source person, and it is difficult to ensure the identity consistency of the generated results. To solve these problems, bottom-up methods are proposed.

*3.1.2 Bottom-up Methods.* Bottom-up methods tend to decompose the whole process into components or steps, where the intermediate results are usually essential for the generation. In general, these methods can be grouped into two categories: latent feature representation and deformation methods.

**Latent Feature Representation Methods** [22, 39, 112, 156, 189, 218] extracts latent features from the source image and target pose as the intermediate results, to control the generation results. The works based on latent feature representation aim at extracting accurate and pure feature

information. A typical way [22, 39] is to extract appearance features from the source image via VAE [84]. Chen *et al.* [22] use cycle-consistency [231] to support unpaired training data. Zhao *et al.* [218] extract pose features from an interpolated pose sequence, from source pose to target pose. Sun *et al.* [156] extract appearance features from a set of source images by bidirectional convolutional LSTM. Ma *et al.* [112] use adversarial training to disentangle the input into three factors: foreground, background and pose.

Compared to top-down methods, latent feature representation methods have better performance and preserve the identity information to a certain extent.

**Deformation Methods** are another line of research. Top-down and latent feature representation methods usually show inferior results for large differences between source and target poses. Deformation methods are proposed to address the above deficiency, by transferring the feature maps of the source person image into the target pose. These methods are better at preserving texture details. Deformation methods can be further divided into three developmental branches according to different deformation strategies: local deformation, global deformation and attention-based deformation.

*Local deformation methods* [7, 35, 95, 148] adopt the strategy of disintegrating the body and partial deformation. The first method dates back to Deformable GAN (Def-GAN). Inspired by the spatial transformer networks [72], Siarohin *et al.* [148] propose to locally deform near-rigid body parts at the feature map level. A new Deformable Skip Connection (DSC) is proposed to replace standard skip operation, where the deformed feature and the target pose are concatenated as a comprehensive input. DSC uses a rough deformation strategy, so some works [7, 95] improve over DSC. Especially, Liang *et al.* [95] further consider the modification of the body shape in the local deformation process. Besides, a few methods choose not to explicitly deform body parts at the input side, but to deform implicitly with an encoder-decoder structure during generation. Dong *et al.* [35] introduce Part-Preserving GAN (PP-GAN), which directly takes as input the decomposed body parts and implicitly deforms with a parsing-consistent loss. Although the local deformation strategy retains enough appearance details, it ignores the connections between the local parts, resulting in a poor generation effect in these connection areas.

*Global deformation methods* [30, 32, 53, 57, 94, 106, 211, 219] adopt global deformation to warp the human body as a whole. Dong *et al.* [32] propose a soft-gated warping-block to predict the transformation grid between parsing maps. Some works [30, 57, 94, 106, 219] estimate the appearance flow from source pose to target pose, to transform the full body. Furthermore, 3D models are introduced to improve warping accuracy. Zanfir *et al.* [211] fill in the dense 3D mesh of the source person with textures and then transfer the mesh into the target result. Grigorev *et al.* [53] first extract the 3D DensePose [5] for target pose, and then inpaint the surfaces with a coordinate-based method. Global deformation methods can preserve enough appearance details, but they are hard to deal with the appearance details of unseen areas.

*Attention-based deformation methods* [18, 92, 103, 136, 160, 233] introduce the attention mechanism for implicit deformation. Zhu *et al.* [233] propose Pose-Attentional Transfer Network (PATN) to assign weights to different image patches according to the pose. Ren *et al.* [136] generate flow fields to calculate local attention for feature maps extracted from the source image, and thus spatial transformation can be locally operated. The above two methods only contain the pose-based attention branch. Furthermore, Tang *et al.* [160] design a novel attention-based network XingGAN including two branches: Shape-guided Appearance-based generation (SA) and Appearance-guided Shape-based generation (AS). A novel crossing connection is introduced to capture the joint attention between image and pose modalities. Attention-based deformation methods can well predict the appearance details of unseen areas and show great potential.

*3.1.3  Hybrid Methods.* Top-down methods are better at preserving the overall appearance but not at retaining the local textures, especially there is a large gap between source and target poses. Bottom-up methods deal well with large gaps, but have difficulties in handling occlusions. Therefore, the hybrid frameworks are explored to take advantage of both sides.

Neverova *et al.* [122] propose the first hybrid method, Dense Pose Transfer (DPT), which fuses the top-down result and deformation output for target image generation. Yang *et al.* [197, 198] combine PATN [233] with image refinement [111] and apply an alternate updating strategy to facilitate mutual guidance between two modules for better appearance and shape consistency. Xu *et al.* [193] propose MR-Net to synthesize person images with distorted poses and cluttered scenes.

## 3.2  Pose-guided Person Video Generation

It is natural to extend pose-guided generation from image to video domain, by further considering temporal coherence in generated videos. Similar to image generation, we categorize related works into two branches: top-down and bottom-up methods.

*3.2.1  Top-down Methods.* Top-down methods directly map a target pose sequence to a person video via GAN-based networks. The target pose sequence is either given [194] or predicted by another network [174, 179, 195, 217]. Yang *et al.* [195] predict the pose sequence via Pose Sequence GAN (PSGAN). Walker *et al.* [179] combine VAE [84] with LSTM to learn the distribution of future poses. Meanwhile, pose flow is introduced to improve temporal coherence. Furthermore, Villegas *et al.* [174] propose an analogy generating method to synthesize future frames. Zhao *et al.* [217] improve over [174] with motion refinement network for better temporal coherence.

*3.2.2  Bottom-up Methods.* To retain more temporal coherence, bottom-up methods are explored to inject more structured information during video generation. Wang *et al.* [185] propose the first video-to-video synthesis approach, where optical flow is adopted to improve temporal coherence. Chan *et al.* [15] introduce a typical pose-guided video generation method to transfer body motions from another video. A sequence of frames, instead of a single image, is adopted to preserve temporal coherence. Liu *et al.* [101] better transfer poses on 3D model.

The above methods [15, 101, 185] are limited in identity dependency. That is, one has to re-train the model when applying it to a new person. Some researchers attempt to design identity-independent methods [135, 146, 147, 157, 184, 199, 204, 208, 230]. Siarohin *et al.* [146, 147] add a source image as input to generate video for different identities. Furthermore, optical flow between source and driving frames is estimated to preserve temporal coherence. Wang *et al.* [184] also add identify images as additional inputs to control the appearance of the generated video. Zhou *et al.* [230] present a local deformation method that deforms body parts into target pose. Ren *et al.* [135] remould image generation network [136] to synthesize continuous frames, and introduce a Motion Extraction Network to improve temporal consistency. Yoon *et al.* [204] use the 3D human model to transfer the pose in a temporally consistent way, which can be applied to in-the-wild images.

## 3.3  Discussion

Pose-guided person generation is challenging mainly due to the diversity and complexity of human poses. For large gaps between source and target poses, unseen areas are extremely difficult to synthesize. Top-down methods adopt GAN-based networks to hallucinate unseen areas, but sacrifice the texture details. Bottom-up methods deform the source image into the target pose, which better preserves texture details but is sensitive to occluded areas.

So far, most pose-guided generation methods belong to the bottom-up branch, while top-down approaches are also evolving fast in visual quality. Among bottom-up methods, latent feature representation methods are pure implicit learning methods with poor interpretation, but they

Fig. 9. Illustration of virtual try-on. The target cloth is given as an in-shop image (a), or worn by a model person (b).

have a strong ability to predict invisible regions according to visible regions. On the contrary, explicit deformation methods use many strong explanatory strategies including disentangled affine transformation, appearance flow warping and 3D-based deformation, But it is difficult to predict the details of the appearance of invisible areas. While, deformation methods based on the attention mechanism, which implicitly deform the pose in latent feature space, take the advantage of both latent feature representation and explicit deformation methods. They can predict the details of the visible area well and have shown great potential in further improving the texture details. However, large differences between poses, e.g., front to back, are still difficult to handle for now. The hybrid methods have more potentials in the future in addressing the difficulties from occlusions and texture details.

Recently, monocular 3D human reconstruction has been developed gradually [141, 142, 223]. However, only reconstruction based on the full-body image is supported. For pose-guided person generation, many source images do not show complete bodies, so it is difficult to apply the 3D reconstruction idea. 3D models have strong controllability. 3D models have great potential in future pose-guided person generation.

## 4 GARMENT-ORIENTED PERSON GENERATION

Garment-oriented generation concentrates on generating new clothing of a person image, which has wide applications in the fashion domain. In this section, we focus on two mainstream tasks: virtual try-on and garment manipulation.

### 4.1 Virtual try-on

Virtual try-on, as shown in Fig. 9, is to transfer a specific garment onto a person image, which allows customers to virtually try on garments before online shopping. The early method [28] focuses on 3D body and clothing reconstruction, mostly via manual modeling. This method is complex in procedure, and its results are less visually realistic. Since 2017, GAN-based approaches start to emerge as a complement to the 3D solution. As a comparison, 2D methods have much simple architecture and produce more realistic results.

A brief summary of related works is shown in Tab. 4. Deep-learning based virtual try-on can also be grouped into three categories: top-down, bottom-up and hybrid methods. Fig. 10 shows the representative methods, and Fig. 11 highlights the differences among the three branches.
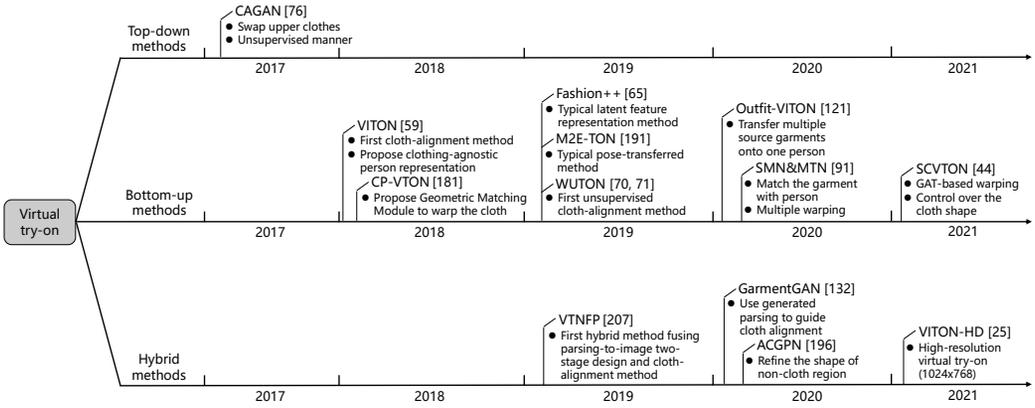
Fig. 10. Representative works on virtual try-on.



Fig. 11. Comparison among four virtual try-on methods. (a) Top-down methods directly input $I_s$, $I_g$, $S_t$ to synthesize output $I_t$. (b) Bottom-up latent feature representation methods devote to represent the shape, garment and other features for synthesis. (c) Bottom-up cloth-alignment methods warp $I_g$ into the shape of $I_s$ via a warping model. (d) Hybrid methods combine top-down with the cloth-alignment method.

*4.1.1 Top-down Methods.* Top-down methods deal with virtual try-on from a global perspective. The straightforward way is to directly map the source person image to the final result, conditioning on the target garment. However, this approach only obtains a rough result. To better preserve shape and texture, human parsing maps are added as another conditioning input [76, 86, 133].

In general, paired data is necessary for decent results, i.e., two images with the same identity and pose but with different garments. This assumption poses a big barrier in real scenarios. Therefore, unsupervised methods that operate on unpaired data, are more practical in real cases. Raj *et al.* [133] propose SwapNet that swaps garments between a pair of person images while preserving the pose and body shape. Self-supervised training is also applied to address the lack of paired data. Jetchev and Bergmann [76] introduce a Conditional Analogy GAN (CAGAN) to transfer the upper clothes of a person into the target one, where cycle-consistency [231] is adopted in training. Kubo *et al.* [86] improve over CAGAN [76] by segmenting the clothing region for better precision.

Table 4. Summary of virtual try-on methods. Applicable garment: applicable types of input garment.

| References | Main idea | Applicable garment |
|---|---|---|
| Top-down methods | | |
| Raj [133] | Parsing-to-image two-step generation | Garment of model person image |
| [76, 86] | Cycle consistency | In-shop upper cloth image |
| Bottom-up latent feature representation methods | | |
| Hsiao [65] | Garment and pose feature representation and parsing-to-image two-step generation | Multiple garment images |
| Neuberger [121] | Shape and garment feature representation parsing-to-image two-step generation | Multiple garment images |
| Bottom-up cloth-alignment methods | | |
| [59, 150] | Cloth-alignment coarse-to-fine generation | In-shop upper cloth image |
| [117, 181] | Cloth alignment | In-shop upper cloth image |
| Lee [73] | Cloth alignment and coarse-to-fine generation | In-shop upper cloth image |
| Kubo [87] | Coarse-to-fine 3D pose surface painting | In-shop upper cloth image |
| Han [57] | Segmentation-based clothing flow estimation | In-shop upper cloth image |
| Issenhuth [70, 71] | Cloth alignment and knowledge distillation | In-shop upper cloth image |
| Li [91] | Garment-person pair matching and multi-warp | Upper cloth image |
| Xie [192] | Decomposed cloth alignment | In-shop upper cloth image |
| Fincato [42] | Two-stage cloth alignment | In-shop upper cloth image |
| Gao [44] | Graph Attention Network based cloth alignment | In-shop upper cloth image |
| Ge [48] | Appearance flow based cloth alignment | In-shop upper cloth image |
| Ge [46] | Cycle consistency and cloth alignment | In-shop upper cloth image |
| Ren [134] | Transformer based cloth alignment | In-shop upper cloth image |
| Liu [107] | Human pose transfer and cycle consistency | Garment of model person image |
| [106, 191] | Deformation method of human pose transfer | Garment of model person image |
| Yu [206] | Pose-transferred parsing-to-image generation | Garment of model person image |
| Roy [139] | Pose transformation and coarse-to-fine generation | Garment of model person image |
| Cui [30] | Global flow field based cloth alignment | Garment of model person image |
| Hybrid methods | | |
| [132, 207] | Cloth-alignment parsing-to-image generation | In-shop upper cloth image |
| [75, 115, 116, 155, 196] | Cloth-alignment parsing-to-image generation and non-cloth region refinement | In-shop upper cloth image |
| Choi [25] | High-resolution cloth-alignment parsing-to-image generation | In-shop upper cloth image |

Similar to the pose-guided person generation, top-down virtual try-on methods are also difficult to preserve sufficient appearance details and garment texture details. Therefore, more people focus on bottom-up methods.

*4.1.2 Bottom-Up Methods.* Top-down methods suffer from poor texture details. Thus two typical solutions are explored to improve the clothing details, namely latent feature representation and cloth-alignment.

**Latent Feature Representation Methods** [58, 65, 90, 121] represent clothing attributes (e.g. shape and appearance) with latent features to better control the process of generation. These methods are widely adopted in tasks involving multiple garments. Neuberger *et al.* [121] propose Outfit-VITON to try on multiple garments simultaneously, where the shape and appearance features are extracted to represent selected clothes. Hsiao *et al.* [65] propose Fashion++ to slightly adjust full-body clothing outfits, where multiple texture and shape features encode different parts of the garment and body. Although the strong decomposition of latent feature representation methods is suitable for virtual try-on of multiple garments, some details are inevitably lost in the process of feature extraction.

**Cloth-alignment Methods** warp the target cloth into the shape of the source person. Due to its simplicity and high performance, cloth alignment makes the most popular bottom-up branch. According to the alignment strategies, cloth-alignment methods can be further divided into cloth-based warping and pose-based warping methods.

*Cloth-based warping methods* directly deform cloth images with geometric transformation, which are generally applied to in-shop cloth images as shown in Fig. 9(a). The first method is introduced in VITON [59], which aligns clothes considering the pose, body shape and head region. Wang *et al.* [181] improve VITON by introducing Characteristic-Preserving Virtual Try-On Network (CP-VTON), which applies a Geometric Matching Module (GMM) to better align clothes. Subsequent methods are mostly following this paradigm.

Several works improve VITON and CP-VTON in different ways. Some of these methods focus on improving the accuracy of warping. Lee *et al.* [73] and Fincato *et al.* [42] extend GMM to align clothes by a two-step warping. CP-VTON+ [117] explicitly regresses the warped cloth mask to improve the precision of alignment. Li *et al.* [91] propose a Shape-Matching-Net to choose shape-wise compatible "garment-person" pairs. Xie *et al.* [192] estimate the landmarks of the target cloth and warp the left sleeve, right sleeve and middle body, respectively. Gao *et al.* [44] transform the cloth image into the mesh and warp it with Graph Attention network (GAT) [172]. This method gives control over the cloth shape, including length and tightness. Some works introduce new techniques to optimize the results. Issenhuth *et al.* [70, 71] propose a Warping U-net for virtual Try-On Net (WUTON), which does not need the ground-truth target image during training. Meanwhile, they train a parsing-free student WUTON based on teacher-student knowledge distillation [64] strategy. Furthermore, Ge *et al.* [48] add a tutor model based on WUTON to train a more precise parsing-free virtual try-on model named PFAFN. Ren *et al.* [134] introduce Transformer [171] to cloth alignment to capture long-range relation between garment and person.

Cloth-based warping methods can only deform clothing images based on body pose and shape, so it is difficult to handle the situations of complex body poses and occlusions.

*Pose-based warping methods* take cloth images as a part of a person body, to deform based on body pose. UVTON [87] deforms on the 3D dense pose [5] and then fills the surface with clothing textures. Han *et al.* [57] propose to warp based on clothing flow.

This idea is applicable to the situation where reference cloth is worn by a model person as shown in Fig. 9(b). It greatly expands the application scope of cloth-alignment methods. Pose-guided generation (Sec. 3.1) is often adopted to register the model and source cloth. Liu *et al.* [107] introduce SwapGAN to directly transfer the model person to the target pose. Wu *et al.* [191] introduce an M2E Try-On Net (M2E-TON) to deform based on pose supervision. Roy *et al.* [139] introduce LGVTON to warp upper cloth using both pose and fashion landmarks [109]. I-VTON [206] combines warped textures of multiple body parts for virtual try-on and then inpaints the missing parts with realistic appearances. Cui *et al.* [30] estimate the global flow field (between source person and model person poses) to warp the garment, which can put on multiple garments in a certain order.

### 4.1.3 Hybrid Methods.
Top-down methods can well capture the overall structure of a person and garments, but lose a lot of texture details. Bottom-up cloth-alignment methods easily preserve details on the garments, but the overall structure looks less natural. Therefore, some works seek a hybrid solution combining both perspectives.

Some methods [25, 132, 207] take as input the human parsing map and aligned clothes, such that the top-down model is also aware of the internal structure used in bottom-up approaches. Moreover, Choi *et al.* [25] design an ALIAS generator to synthesize high-resolution images (1024×768). Other methods [75, 115, 116, 155, 196] preserve details in the non-cloth region, also guided by the parsing map. Yang *et al.* [196] propose an Adaptive Content Generating and Preserving Network (ACGPN) that adds a non-target composition to refine the shape of the non-cloth region. Minar *et al.* [115, 116] align clothes based on a 3D clothing model. Meanwhile, they further consider non-cloth regions (e.g., skin, lower body) to comprehensively preserve details on the whole body.

## 4.2 Garment Manipulation

Besides virtual try-on, there are several tasks for garment manipulation, such as garment synthesis, text-guided garment manipulation and garment inpainting.

Garment synthesis generates person images with new garments. ClothNet [90] generates photo-realistic clothing images from sketch with image-to-image translation [69]. Text-guided garment manipulation aims to modify clothes with the guidance of text. Zhu *et al.* [232] propose FashionGAN to generate clothing according to text descriptions. Garment inpainting uses in-painting methods to add garment textures. Han *et al.* [58] present the FiNet to inpaint the missing garments such as upper/lower garments and shoes.

## 4.3 Discussion

Virtual try-on is challenging for two main reasons. The first issue comes from the diversity of human poses, as well as the mismatch between garments and poses. Top-down methods can generate photo-realistic results for a wide range of poses, but can not well preserve the texture details. Latent feature representation methods support multi-garment virtual try-on, but are difficult to keep enough cloth textures. Cloth-alignment methods preserve most cloth textures, but can hardly handle the large mismatch between clothes and poses. In the future, we can combine latent feature representation with cloth alignment for predicting the clothing area as accurately as possible while retaining more details. The second issue is the lack of proper datasets for fully supervised training. Paired data is difficult to collect in the scenario of virtual try-on. But, the obstacle in supervised data gives rise to unsupervised methods such as [70, 76, 121].

So far, bottom-up methods are more popular. Due to the advantage of retaining texture details, cloth-alignment methods are developing rapidly. To improve the generation quality, many deformation methods have been explored, including decomposed TPS [192] and ClothFlow [57]. Besides, the variety of garments for try-on is getting much wider recently. Early works only support upper clothes. Recent works [65, 121] begin to try on clothes including upper clothes, trousers, skirts, hats, etc. However, current methods can not work well on wild datasets and multi-garments.

## 5 BENCHMARKS

We first review the datasets and evaluation metrics in this section, and then benchmark recent performances where applicable.

## 5.1 Datasets

Details of popular datasets for talking-head generation are summarized in Tab. 5, and those for pose and garment-oriented generation are in Tab. 6. Here we only briefly highlight the datasets adopted by the major researchers.

For talking-head generation, GRID [27] and LRW [26] are adopted to evaluate methods without head motion; CREMA-D [14] validates on videos with spontaneous motions; VoxCeleb2 [149] and LRS3-TED [2] focus on in-the-wild videos. Alternatively, Chen *et al.* [19] also suggest several evaluation protocols for evaluating talking heads generation.

For pose-guided person image generation, the most commonly adopted datasets are Market-1501 [220] and DeepFasion [108]. For garment-based generation, Zalando [59] and DeepFasion [108] datasets are popular for evaluation. Recently, MPV [33] dataset is also used to evaluate garment and pose-based generation.

Table 5. Summary of talking-head video datasets.

| Name | Year | Data scale (in hours) | # Speaker | # Sentence | Head movement | Emotion | Open source |
|---|---|---|---|---|---|---|---|
| GRID | 2006 [27] | 27.5 | 33 | 33k | × | × | Link |
| TCD-TIMIT | 2015 [61] | 11.1 | 62 | 6.9k | × | × | Link |
| LRW | 2016 [26] | 173 | 1k+ | 539k | × | × | Link |
| MODALITY | 2017 [31] | 31 | 35 | 5.8k | × | × | Link |
| CREMA-D | 2014 [14] | 11.1 | 91 | 12 | √ | √ | Link |
| MSP-IMPROV | 2016 [13] | 18 | 12 | 652 | √ | √ | Link |
| ObamaSet | 2017 [158] | 14 | 1 | — | √ | × | Link |
| VoxCeleb | 2017 [120] | 352 | 1.2k | 153.5k | √ | × | Link |
| VoxCeleb2 | 2018 [149] | 2.4k | 6.1k | 1.1m | √ | × | Link |
| RAVDESS | 2018 [110] | 7 | 24 | 2 | √ | √ | Link |
| LRS2-BBC | 2018 [1] | 224.5 | 500+ | 140k+ | √ | × | Link |
| LRS3-TED | 2018 [2] | 438 | 5k+ | 152k+ | √ | × | Link |
| MELD | 2018 [127] | 13.7 | 407 | 13.7k | √ | √ | Link |
| Lombard | 2018 [4] | 3.6 | 54 | 5.4k | √ | √ | Link |
| Faceforensics++ | 2019 [137] | 5.7 | 1k | 1k+ | √ | × | Link |
| MEAD | 2020 [183] | 39 | 60 | 20 | √ | √ | Link |

Table 6. Datasets for pose and garment-oriented person generation. PGPIG: Pose-Guided Person Image Generation. PGPVG: Pose-Guided Person Video Generation. VTON: Virtual Try-on. VF: Virtual Fitting. VVTON: Video Virtual Try-on.

| Name | Year | Data scale | Applicable fields | Open source |
|---|---|---|---|---|
| Market-1501 | 2015 [220] | 32,668 images of 1,501 persons | PGPIG | Link |
| DeepFasion | 2016 [108] | 52,712 in-shop cloth images and over 200,000 cross-pose / scale pairs | PGPIG and VTON | Link |
| MVC | 2016 [100] | 161,638 clothing images of 37,499 items | PGPIG and VTON | Link |
| Human3.6M | 2013 [68] | 3,578,080 images of 11 persons | PGPIG | Link |
| Chictopia10k | 2015 [96] | 17,706 images | VTON | Link |
| Zalando | 2018 [59] | 16,253 person-cloth pairs | VTON | NA |
| LookBook | 2016 [203] | 75,016 person images and 9,732 in-shop cloth images | PGPIG and VTON | NA |
| MPV | 2019 [33] | 35,687 person images and 13,524 in-shop cloth images | PGPIG, VTON and VF | Link |
| FashionOn | 2019 [67] | 21,790 person images and 10,895 in-shop cloth images | PGPIG, VTON and VF | NA |
| FashionTryOn | 2019 [221] | 57,428 person images and 28,714 in-shop cloth images | PGPIG, VTON and VF | Link |
| Penn Action | 2013 [214] | 2,326 videos | PGPIG and PGPVG | Link |
| Tai-Chi | 2018 [170] | 4,500 videos | PGPIG and PGPVG | Link |
| Fashion | 2019 [208] | 600 videos | PGPIG and PGPVG | Link |
| iPER | 2019 [106] | 206 videos of 30 persons | PGPIG, VTON and PG-PVG | Link |
| VVT | 2019 [34] | 791 videos, 791 person images and 791 cloth images | PGPIG, VTON, PGPVG and VVTON | Link |

## 5.2 Evaluation Metrics

Evaluating generation tasks are known difficult. As a common practice, multiple objectives (e.g, Inception Score, SSIM) and subjectives (e.g., Amazon Mechanical Turk, User study) metrics are adopted for a comprehensive evaluation. Subjective evaluation involves humans in the loop, which is often applied to compare the perceptual visual quality of generated content. However, due to the subjective factors and higher costs during evaluation, most works also seek quantitative evaluations with objective metrics. For person generation, the common objective metrics are summarized below.

**SSIM** (Structural Similarity) [187] measures the quality of generated image compared with the original image, which is widely used in image synthesis. Specifically, SSIM calculates the similarity between the synthesized image and the ground-truth real image in three dimensions luminance, contrast and structure.

**IS** (Inception Score) [143] evaluates the quality of generated image in terms of clear objects and high diversity. IS calculates the distribution of generated images via a pre-trained Inception

v3 [159] network. This metric is not that informative, and Barratt and Sharma [9] present some shortcomings of IS.

**FID** (Fréchet Inception Distance) [63] calculates the Fréchet distance between the distribution of generated images and real images to capture their similarity.

**FReID** uses a pre-trained person re-ID model to estimate the Gaussian distribution of generated images and real images, which calculates the Fréchet distance of these distributions.

**LPIPS** (Learned Perceptual Image Patch Similarity) [213] uses a pre-trained deep network to learn deep perceptual features of the image patch and then computes the average $\ell_2$ distance between features of two images.

**DS** (Detection Score) measures the confidence of a pre-trained person detector in a generated person image, by computing the average person-class detection scores on generated images.

**AttrRec-k** (Top-k Clothing Attribute Retaining Rate) uses a pre-trained clothing attribute recognition model to predict clothing attributes from a generated image. It utilizes the top-k recall rate as the final score.

### 5.3 Performance Comparison

Talking-head video generation is mostly evaluated by user studies. Recently, Chen *et al.* [19] propose several objective metrics to evaluate talking-head videos in terms of identity preserving, lip synchronization, video quality, and spontaneous head movements. Meanwhile, they present a new benchmark with standardized dataset pre-processing strategies for evaluating talking-head generation methods.

For pose-guided person generation and virtual try-on, SSIM [187], IS [143] and their variants are widely applied. But these metrics are not perfect [9, 213]. People gradually turn to FID [63] and LPIPS [213] for evaluation. In particular, user studies are more widely used in virtual try-on.

Tab. 7 shows the quantitative results of major state-of-the-art pose-guided person image generation methods. On average, deformation and hybrid methods are better than top-down and latent feature representation methods. For a more intuitive comparison, Fig. 12 shows the qualitative

Table 7. Comparison of state-of-the-art methods for pose-guided person image generation. ↑: larger is better.

| Methods | Year | Market-1501 | | | | | DeepFashion | | |
|---|---|---|---|---|---|---|---|---|---|
| | | SSIM↑ | IS↑ | mask-SSIM↑ | mask-IS↑ | FID↓ | SSIM↑ | IS↑ | FID↓ |
| Zhao [216] | 2017 | — | — | — | — | — | 0.62 | 3.03 [122] | — |
| Ma [111] | 2017 | 0.253 | 3.460 | 0.792 | 3.435 | — | 0.762 | 3.090 | — |
| Pumarola [129] | 2018 | — | — | — | — | — | 0.747 | 2.97 | — |
| Esser [39] | 2018 | 0.353 | 3.214 | — | — | 20.144 [136] | 0.786 | 3.087 | 23.667 [136] |
| Ma [112] | 2018 | 0.099 | 3.483 | 0.614 | 3.491 | — | 0.614 | 3.228 | — |
| Siarohin [148] | 2018 | 0.290 | 3.185 | 0.805 | 3.502 | 25.364 [136] | 0.756 | 3.439 | 18.457 [136] |
| Dong [32] | 2018 | 0.356 | 3.409 | — | — | — | 0.793 | 3.314 | — |
| Neverova [122] | 2018 | — | — | — | — | — | 0.785 | 3.61 | — |
| Sun [156] | 2019 | 0.344 | 3.291 | — | — | — | 0.789 | 3.006 | — |
| Tang [161] | 2019 | 0.282 | 3.349 | 0.811 | 3.510 | — | — | — | — |
| Song [153] | 2019 | 0.203 | 3.499 | 0.758 | 3.680 | — | 0.736 | 3.441 | — |
| Siarohin [145] | 2019 | 0.326 | 3.613 | 0.806 | 3.814 | — | 0.774 | 3.421 | — |
| Liang [95] | 2019 | — | 3.657 | — | 3.614 | 20.355 | — | 3.536 | 29.684 |
| Han [57] | 2019 | — | — | — | — | — | 0.771 | 3.88 | — |
| Li [94] | 2019 | — | — | — | — | 27.163 [136] | 0.778 | 3.338 | 16.314 [136] |
| Dong [35] | 2019 | 0.396 | 3.581 | — | — | — | 0.782 | 3.595 | — |
| Zhu [233] | 2019 | 0.311 | 3.323 | 0.811 | 3.773 | 22.657 [136] | 0.773 | 3.209 | 20.739 [136] |
| Karmakar [78] | 2020 | 0.302 | 3.488 | — | — | — | 0.781 | 3.238 | — |
| Ren [136] | 2020 | — | — | — | — | 19.751 | — | — | 10.573 |
| Yang [197] | 2020 | — | — | — | — | — | 0.774 | 3.125 | 14.611 |
| Li [92] | 2020 | 0.315 | 3.487 | 0.814 | 3.867 | — | 0.775 | 3.338 | — |
| Chen [18] | 2021 | 0.306 | 3.827 | 0.804 | 3.809 | — | 0.760 | 3.348 | — |
| Liu [103] | 2021 | 0.732 | 3.703 | 0.821 | 3.750 | 16.142 | 0.786 | 3.736 | 8.732 |

Table 8. Comparison of state-of-the-art virtual try-on methods. ↑: larger is better.

| Methods | Year | Zalando | | | DeepFashion | | |
|---|---|---|---|---|---|---|---|
| | | SSIM↑ | IS↑ | FID↓ | SSIM↑ | IS↑ | FID↓ |
| Han [59] | 2018 | 0.786 [57] | 2.514 | 41.80 [206] | 0.71 [139] | 2.40 [107] | 78.45 [139] |
| Raj [133] | 2018 | 0.83 | 2.631 [206] | 114.50 [206] | — | — | — |
| Wang [181] | 2018 | 0.792 [57] | 2.748 [150] | 23.60 [206] | 0.72 [139] | 2.41 [139] | 72.95 [139] |
| Song [150] | 2019 | — | 2.656 | — | — | — | — |
| Han [57] | 2019 | 0.803 | — | — | — | — | — |
| Yu [207] | 2019 | 0.803 [196] | 2.784 [196] | — | — | — | — |
| Liu [107] | 2019 | — | — | — | 0.717 | 2.65 | — |
| Wu [191] | 2019 | — | 2.510 [206] | 33.28 [206] | — | — | — |
| Yu [206] | 2019 | — | 2.708 | 29.48 | — | — | — |
| Roy [139] | 2020 | — | — | — | 0.86 | 2.71 | 56.85 |
| Jandial [75] | 2020 | 0.766 | 2.82 | 14.65 | — | — | — |
| Raffie [132] | 2020 | — | 2.774 | 16.578 | — | — | — |
| Yang [196] | 2020 | 0.845 | 2.829 | 15.67 [48] | — | — | — |
| Minar [117] | 2020 | 0.817 | 3.074 | — | — | — | — |
| Fincato [42] | 2021 | 0.886 | 2.76 | 12.45 | — | — | — |
| Ge [48] | 2021 | — | — | 10.09 | — | — | — |
| Ge [46] | 2021 | 0.83 | 2.85 | 14.82 | — | — | — |
| Ren [134] | 2021 | 0.827 | 3.060 | — | — | — | — |



Fig. 12. The qualitative comparisons with state-of-the-art pose-guided person image generation methods on DeepFashion [108] dataset, including PG$^2$ [111], UPIS [129], VU-Net [39], Def-GAN [148], DPT [122], Song *et al.* [153], ClothFlow [57], PATN [233], FHPT [197] and SPAN [103].

comparisons with some state-of-the-art methods. Meanwhile, Fig. 13 shows some failure cases. It demonstrates that the present pose-guided person image generation methods are still difficult to deal with complex poses and in-the-wild images.

Similarly, Tab. 8 gives the comparison for virtual try-on methods. Hybrid methods are more competitive than other methods. Fig. 14 shows the visual comparison with some state-of-the-art methods. Meanwhile, Fig. 15 shows some failure cases. It indicates that the present virtual try-on methods are still difficult to handle different garment shapes and occlusion.

## 6 DISCUSSION

In Sections 2, 3 and 4, we have reviewed talking-head generation, pose-guided person generation and virtual try-on, respectively. In this section, we will discuss the commonalities between these three tasks.
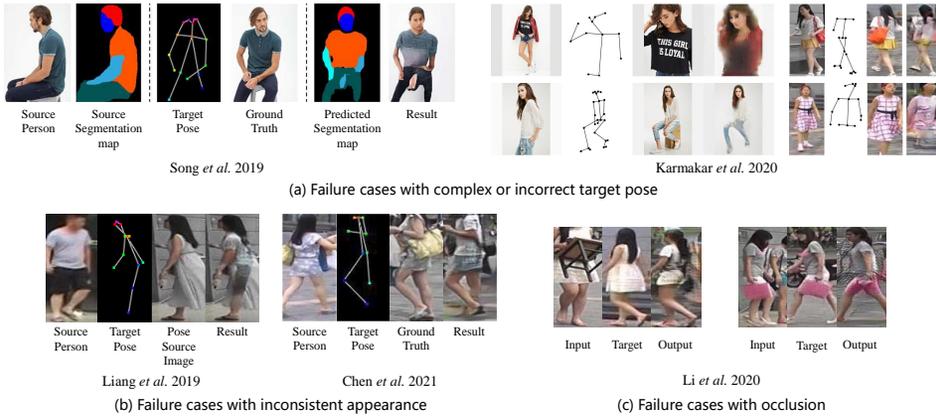
Fig. 13. Some failure cases with (a) complex or incorrect target pose, (b) inconsistent appearance and (c) occlusion. These cases are from Song *et al.* [153], Karmakar *et al.* [78], Liang *et al.* [95], Chen *et al.* [18] and Li *et al.* [92].



Fig. 14. The qualitative comparisons with state-of-the-art virtual try-on methods on Zalando [59] dataset, including VITON [59], CP-VTON [181], VTNFP [207], I-VTON [206], ACGPN [196], CP-VTON+ [117], PF-AFN [48] and CIT [134].

- **Same high-level framework.** The three tasks all firstly construct a representation of the source person identity, extract and encode the modified attributes, and finally combine the identity representation and attribute encoding to obtain the target result. Fig. 16 illustrates this common framework.
- **The idea of deformation.** There are many talking-head generation methods [6, 49, 54, 60, 123, 190, 229] that deform the source face image based on the target facial landmarks. Flow warping is the major deformation method.

  Deformation methods in pose-guided person generation deform the source person according to the target pose. The deformation tools include decoupled deformation, global flow warping and deformation based on the attention mechanism.

(a) Failure cases with wrong arm shape



(b) Failure cases with occlusion

Fig. 15. Some failure cases with (a) wrong arm shape and (b) occlusion. These cases are from Wang *et al.* [181], Fincato *et al.* [42], Minar *et al.* [117] and Ren *et al.* [134].
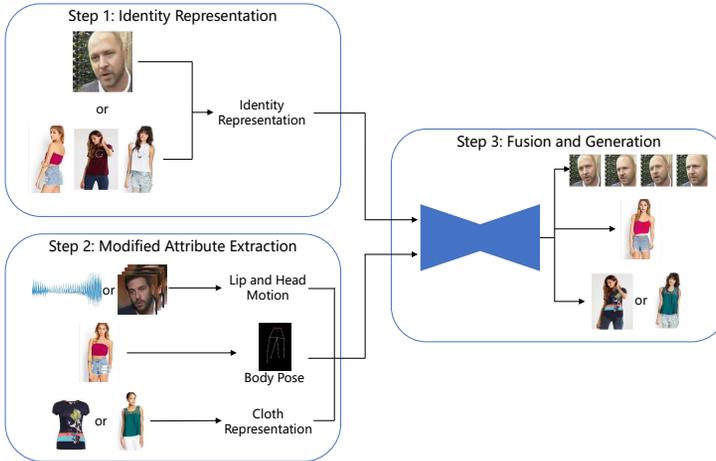


Fig. 16. The framework of talking-head generation, pose-guided person generation and virtual try-on.

Cloth-alignment methods in virtual try-on also deform the target cloth according to the source person pose. TPS and flow warping are two major deformation tools.

These three types of methods are essentially the same: they all deform the image appearance according to the pose. As evidence, some methods [136, 161] use the same model for talking-head generation and pose-guided person generation. Meanwhile, some works [26, 57] use the same methods for pose-guided person generation and virtual try-on simultaneously.

• **The idea of feature representation and fusion.** Since the three tasks all need to express the person identity and modified attribute, latent feature representation methods appear in all three tasks. To preserve the identity information, latent feature representation methods focus on feature decoupling to extract pure identity information and pure attribute information as far as possible.

Talking-head generation pays more attention to the decomposition degree, while pose-guided person generation and virtual try-on pay less attention. Talking-head generation needs to

generate rich facial details. In order to preserve more identity information, interference from other attributes should be eliminated as much as possible. Meanwhile, pose-guided person generation and virtual try-on do not need to generate rich details related to identity information. However, with the development of research and the improvement of demand, it is inevitable to study how to improve the decomposition degree in pose-guided person generation and virtual try-on.

## 7 APPLICATIONS

We have reviewed deep person generation from the perspective of face, pose and cloth synthesis. This section presents typical applications based on the aforementioned fundamental tasks. We mainly focus on three typical applications: Generative Data Augmentation (how generated data can help machines), Virtual Fitting and Digital Human (how generated data can help human beings).

### 7.1 Generative Data Augmentation

Most deep learning models are data-hungry for superior performance. Deep person generation is often adopted as data augmentation in person-related tasks such as person ReID (Re-identification), pedestrian detection and autonomous driving.

Person ReID requires massive images of the same person from different views, but existing datasets can only provide limited images for each ID. Zheng *et al.* [224] generate person images with different views and propose a semi-supervised learning method to utilize the generated unlabeled data. Liu *et al.* [98] use pose-guided generation to augment the ReID. Zhang *et al.* [215] use an improved PG$^2$ [111] network to synthesize person images with diverse poses, where the manual labeling is omitted. There are also methods integrating the generation process into the ReID pipeline [47, 130, 222]. Ge *et al.* [47] propose a Feature Distilling GAN (FD-GAN), where pose-guided generation is used to learn the robust identity-related and pose-unrelated features. Zheng *et al.* [222] present a DG-Net to tactfully combine the re-ID discriminative module with a pose-guided generation module, which generates new person images in new poses and discriminates the ID of the generated images simultaneously.

For tasks in pedestrian detection and autonomous driving, person generation is also utilized for data augmentation. Pedestrian Synthesis GAN (PS-GAN) [125] generates new pedestrian images to enrich the dataset and stably improve the pedestrian detection model. Similarly, Vobecký *et al.* [176] propose a GAN-based framework to augment the pedestrian dataset in autonomous driving. Moreover, they also take human poses as another input to synthesize persons with required poses.

### 7.2 Virtual Fitting

Virtual try-on (Sec. 4.1) emphasizes cloth synthesis, while virtual fitting requires both cloth and pose changes in real scenarios, as shown in Fig.17. For example in online shopping, a big requirement is to generate try-on images based on novel garments and poses. People attempt to combine virtual try-on with pose-guided generation for this purpose. There have been some virtual fitting apps recently, such as Zeekit[5], FXMirror[6] and Magic Mirror[7].

Some approaches are based on top-down methods. FashionOn [67] transfers the source human parsing into the target pose, and then globally fills the transformed parsing map with clothing textures. Another stream of works is based on latent feature representation. Yildirim *et al.* [201] use latent vectors to represent human pose, garment color and garment parts (e.g., shirt, coat, trousers,

---

[5]https://zeekit.me/

[6]http://www.fxmirror.net/

[7]https://www.magicmirror.me/

(a) Illustration of virtual fitting



Zeekit                              FXMirror                          Magic Mirror
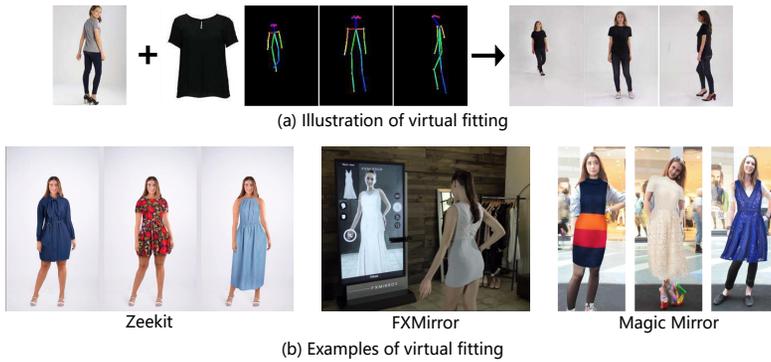(b) Examples of virtual fitting

Fig. 17.  Illustration (a) and example applications (b) of virtual fitting. The target fitting images are controlled by the target cloth and pose.



(a) Video conferencing          (b) Virtual anchor               (c) Customer service
      (Zoombot)                        (Sogou)                          (UneeQ)

(d) Autonomous human           (e) Neon humans                 (f) Sales associate
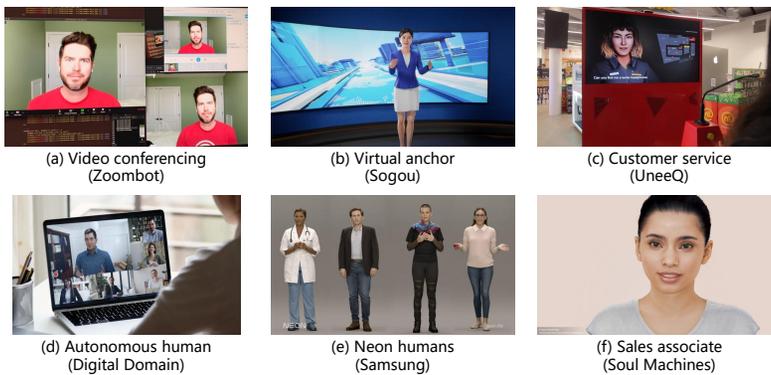   (Digital Domain)                  (Samsung)                       (Soul Machines)

Fig. 18.  Applications of digital human in different scenarios. The company names are noted in parentheses.

skirt, shoes). Men *et al.* [113] meticulously separate a person image into several parts (e.g., pose, head, upper clothes, pants, arms, legs) and then encode them respectively.

Other applications are based on cloth-alignment [33, 66, 182, 221]. By extending CP-VTON, MG-VTON [33] deforms the parsing map by the target pose, and takes the target pose as input to generate the result in the target pose. Wang *et al.* [182] replace standard ResNet Block with tree dilated fusion blocks (tree-blocks) to enrich texture details. Zheng *et al.* [221] design an attentive bidirectional GAN to highlight clothing textures. Besides, Dong *et al.* [34] propose a video-based fitting model (FW-GAN: Flow-navigated Warping GAN) to generate videos with new garments.

### 7.3  Digital Human

Recently, digital human has attracted lots of public attention, which aims to construct a virtual human with realistic appearance and behaviors. Different from robots that have physical body, digital human only exists in the digital realm. At present, digital human has been applied to

many scenarios, such as Zoombot video conferencing AI[8], Sogou AI virtual anchor[9] and customer service[10], which are shown in Fig. 18.

A digital human is expected to interact with people with natural language, facial expressions and body gestures, just like real human beings. While a holistic approach is preferred, a feasible way is to combine talking-head and pose-guided generation. Given a speech signal, talking-head generation can synthesize natural lip motions, facial expressions and head movements. Meanwhile, the body can also be synthesized according to the signal to deliver speech-relevant gestures or spontaneous movements [3, 41, 51, 62, 88]. The models are trained on large-scale datasets to capture generic speech-to-status (poses, expressions, etc.) mapping and personalized styles.

## 8 FUTURE DIRECTIONS

In this paper, we have reviewed deep person generation from three components: face, pose and garment. Thanks to the surprising evolution of deep learning, we have witnessed a rapid development of person generation, from generating low-resolution and rough images to producing high-resolution, detailed and realistic images. However, person generation is still far from mature in generating visually plausible person images/videos on demand. Here we list several future directions worthy of further investigation.

- **Convergence of computer graphics and computer vision.** Computer Graphics (CG) has a mature procedure for creating virtual characters in movies and games. Meanwhile, Computer Vision (CV) has supported the synthesis of photo-realistic human face or body images. Both of them have their unique advantages. Specifically, CG is good at motion control and appearance editing (with explicit mesh, 3DMM, SMPL), and CV produces more realistic appearances (with GAN, Diffusion Model). Recently, there have been some early attempts to combine the advantages of CG and CV, e.g., Neural rendering [162], NeRF [114] based talking-head generation [55], pose-guided generation based on 3D human mesh [101]. For deep person generation, 3D face reconstruction has shown great potential in improving the accuracy of expressions and head motions.

- **Trustworthy contents.** The growing maturity of person generation poses increasing threats to society. Abusing fake person images might cause serious ethical and legal problems, especially in videos of celebrities or politicians. For example, Deepfakes (e.g, head swap, face reenactment) can now produce realistic forged images or videos of celebrities or politicians. Person-related forensics, on the contrary, aims at detecting forged images or videos, which has attracted increasing attention. However, existing works are mostly performance-driven, and thus ignore the model explainability and efficiency. Moreover, most methods are only tuned on a fixed dataset, and the generalization ability is also limited for practical usage. Robust and trustworthy-oriented forgery detection, especially on person-related images or videos, plays an essential role in both accelerating the technical evolution and preventing fake material from being abused.

- **Emerging tasks.** Several derived tasks start to emerge with promising prospects.
  *Conversational head generation.* So far, most people focus on how to make digital human speak according to audio and other conditions. Recently, to make more vivid and natural interaction between digital human and real people, Zhou *et al.* [227] propose responsive listening head generation, which allows digital human to respond as a listener, such as nodding and responsive expressions. Listener-centric or even conversational-centric (both talking

---

and listening) generation emerge and has considerable potential for two-way engagement between a virtual agent and human.

*Person-in-context synthesis.* It aims to generate multiple person instances in complex contexts defined by bounding boxes (layout) [202]. Simply borrowing single-person models for multi-person cases might be sub-optimal. As a more challenging problem, person-in-context synthesis starts to draw attention recently.

*Text-guided person generation.* Text descriptions provide a natural way to interact with humans during generation. Person-related appearance manipulation [228] is particularly helpful for interactive generation based on human request.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Triantafyllos Afouras, Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman. 2018. Deep audio-visual speech recognition. *IEEE transactions on pattern analysis and machine intelligence* (2018).

[2] Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman. 2018. LRS3-TED: a large-scale dataset for visual speech recognition. *arXiv e-prints* (2018), arXiv–1809.

[3] Chaitanya Ahuja, Shugao Ma, Louis-Philippe Morency, and Yaser Sheikh. 2019. To react or not to react: End-to-end visual pose forecasting for personalized avatar during dyadic conversations. In *2019 International Conference on Multimodal Interaction*. 74–84.

[4] Najwa Alghamdi, Steve Maddock, Ricard Marxer, Jon Barker, and Guy J Brown. 2018. A corpus of audio-visual Lombard speech with frontal and profile views. *The Journal of the Acoustical Society of America* 143, 6 (2018), EL523–EL529.

[5] Rıza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. 2018. Densepose: Dense human pose estimation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 7297–7306.

[6] Hadar Averbuch-Elor, Daniel Cohen-Or, Johannes Kopf, and Michael F Cohen. 2017. Bringing portraits to life. *ACM Transactions on Graphics (TOG)* 36, 6 (2017), 1–13.

[7] Guha Balakrishnan, Amy Zhao, Adrian V Dalca, Fredo Durand, and John Guttag. 2018. Synthesizing images of humans in unseen poses. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8340–8348.

[8] Aayush Bansal, Shugao Ma, Deva Ramanan, and Yaser Sheikh. 2018. Recycle-gan: Unsupervised video retargeting. In *Proceedings of the European conference on computer vision (ECCV)*. 119–135.

[9] Shane Barratt and Rishi Sharma. 2018. A Note on the Inception Score. *arXiv e-prints* (2018), arXiv–1801.

[10] Volker Blanz and Thomas Vetter. 1999. A morphable model for the synthesis of 3D faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*. 187–194.

[11] Christoph Bregler, Michele Covell, and Malcolm Slaney. 1997. Video rewrite: Driving visual speech with audio. In *Proceedings of the 24th annual conference on Computer graphics and interactive techniques*. 353–360.

[12] Egor Burkov, Igor Pasechnik, Artur Grigorev, and Victor Lempitsky. 2020. Neural Head Reenactment with Latent Pose Descriptors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13786–13795.

[13] Carlos Busso, Srinivas Parthasarathy, Alec Burmania, Mohammed AbdelWahab, Najmeh Sadoughi, and Emily Mower Provost. 2016. MSP-IMPROV: An acted corpus of dyadic interactions to study emotion perception. *IEEE Transactions on Affective Computing* 8, 1 (2016), 67–80.

[14] Houwei Cao, David G Cooper, Michael K Keutmann, Ruben C Gur, Ani Nenkova, and Ragini Verma. 2014. CREMA-D: Crowd-sourced emotional multimodal actors dataset. *IEEE transactions on affective computing* 5, 4 (2014), 377–390.

[15] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. 2019. Everybody dance now. In *Proceedings of the IEEE International Conference on Computer Vision*. 5933–5942.

[16] Yao-Jen Chang and Tony Ezzat. 2005. Transferable videorealistic speech animation. In *Proceedings of the 2005 ACM SIGGRAPH/Eurographics symposium on Computer animation*. 143–151.

[17] James Charles, Derek Magee, and David Hogg. 2016. Virtual immortality: Reanimating characters from tv shows. In *European Conference on Computer Vision*. Springer, 879–886.

[18] Baoyu Chen, Yi Zhang, Hongchen Tan, Baocai Yin, and Xiuping Liu. 2021. PMAN: Progressive Multi-Attention Network for Human Pose Transfer. *IEEE Transactions on Circuits and Systems for Video Technology* (2021).

[19] Lele Chen, Guofeng Cui, Ziyi Kou, Haitian Zheng, and Chenliang Xu. 2020. What comprises a good talking-head video generation?. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*.

[20] Lele Chen, Zhiheng Li, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. 2018. Lip movements generation at a glance. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 520–535.

[21] Lele Chen, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. 2019. Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 7832–7841.

[22] Xu Chen, Jie Song, and Otmar Hilliges. 2019. Unpaired pose guided human image generation. In *Conference on Computer Vision and Pattern Recognition (CVPR 2019)*. Computer Vision Foundation (CVF).

[23] Zhuo Chen, Chaoyue Wang, Bo Yuan, and Dacheng Tao. 2020. PuppeteerGAN: Arbitrary Portrait Animation with Semantic-aware Appearance Transformation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13518–13527.

[24] Wen-Huang Cheng, Sijie Song, Chieh-Yun Chen, Shintami Chusnul Hidayati, and Jiaying Liu. 2021. Fashion meets computer vision: A survey. *ACM Computing Surveys (CSUR)* 54, 4 (2021), 1–41.

[25] Seunghwan Choi, Sunghyun Park, Minsoo Lee, and Jaegul Choo. 2021. VITON-HD: High-Resolution Virtual Try-On via Misalignment-Aware Normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 14131–14140.

[26] Joon Son Chung and Andrew Zisserman. 2016. Lip reading in the wild. In *Asian Conference on Computer Vision*. Springer, 87–103.

[27] Martin Cooke, Jon Barker, Stuart Cunningham, and Xu Shao. 2006. An audio-visual corpus for speech perception and automatic speech recognition. *The Journal of the Acoustical Society of America* 120, 5 (2006), 2421–2424.

[28] Frédéric Cordier, WonSook Lee, Hyewon Seo, and Nadia Magnenat-Thalmann. 2001. Virtual-try-on on the web. *Laval Virtual* (2001).

[29] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael J Black. 2019. Capture, learning, and synthesis of 3d speaking styles. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 10101–10111.

[30] Aiyu Cui, Daniel McKee, and Svetlana Lazebnik. 2021. Dressing in Order: Recurrent Person Image Generation for Pose Transfer, Virtual Try-On and Outfit Editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3940–3945.

[31] Andrzej Czyzewski, Bozena Kostek, Piotr Bratoszewski, Jozef Kotus, and Marcin Szykulski. 2017. An audio-visual corpus for multimodal automatic speech recognition. *Journal of Intelligent Information Systems* 49, 2 (2017), 167–192.

[32] Haoye Dong, Xiaodan Liang, Ke Gong, Hanjiang Lai, Jia Zhu, and Jian Yin. 2018. Soft-gated warping-gan for pose-guided person image synthesis. In *Advances in neural information processing systems*. 474–484.

[33] Haoye Dong, Xiaodan Liang, Xiaohui Shen, Bochao Wang, Hanjiang Lai, Jia Zhu, Zhiting Hu, and Jian Yin. 2019. Towards multi-pose guided virtual try-on network. In *Proceedings of the IEEE International Conference on Computer Vision*. 9026–9035.

[34] Haoye Dong, Xiaodan Liang, Xiaohui Shen, Bowen Wu, Bing-Cheng Chen, and Jian Yin. 2019. FW-GAN: Flow-navigated Warping GAN for Video Virtual Try-on. In *Proceedings of the IEEE International Conference on Computer Vision*. 1161–1170.

[35] Haoye Dong, Xiaodan Liang, Chenxing Zhou, Hanjiang Lai, Jia Zhu, and Jian Yin. 2019. Part-preserving pose manipulation for person image synthesis. In *2019 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 1234–1239.

[36] Michail Christos Doukas, Mohammad Rami Koujan, Viktoriia Sharmanska, Anastasios Roussos, and Stefanos Zafeiriou. 2021. Head2head++: Deep facial attributes re-targeting. *IEEE Transactions on Biometrics, Behavior, and Identity Science* 3, 1 (2021), 31–43.

[37] Sefik Emre Eskimez, Ross K Maddox, Chenliang Xu, and Zhiyao Duan. 2020. End-To-End Generation of Talking Faces from Noisy Speech. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 1948–1952.

[38] Sefik Emre Eskimez, You Zhang, and Zhiyao Duan. 2021. Speech driven talking face generation from a single image and an emotion condition. *IEEE Transactions on Multimedia* (2021).

[39] Patrick Esser, Ekaterina Sutter, and Björn Ommer. 2018. A variational u-net for conditional appearance and shape generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8857–8866.

[40] Bo Fan, Lijuan Wang, Frank K Soong, and Lei Xie. 2015. Photo-real talking head with deep bidirectional LSTM. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 4884–4888.

[41] Ylva Ferstl, Michael Neff, and Rachel McDonnell. 2019. Multi-objective adversarial gesture generation. In *Motion, Interaction and Games*. 1–10.

[42] Matteo Fincato, Federico Landi, Marcella Cornia, Fabio Cesari, and Rita Cucchiara. 2021. VITON-GT: An Image-based Virtual Try-On Model with Geometric Transformations. In *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 7669–7676.

[43] Ohad Fried, Ayush Tewari, Michael Zollhöfer, Adam Finkelstein, Eli Shechtman, Dan B Goldman, Kyle Genova, Zeyu Jin, Christian Theobalt, and Maneesh Agrawala. 2019. Text-based editing of talking-head video. *ACM Transactions on Graphics (TOG)* 38, 4 (2019), 1–14.

[44] Xin Gao, Zhenjiang Liu, Zunlei Feng, Chengji Shen, Kairi Ou, Haihong Tang, and Mingli Song. 2021. Shape Controllable Virtual Try-on for Underwear Models. In *Proceedings of the 29th ACM International Conference on Multimedia*. 563–572.

[45] Pablo Garrido, Levi Valgaerts, Hamid Sarmadi, Ingmar Steiner, Kiran Varanasi, Patrick Perez, and Christian Theobalt. 2015. Vdub: Modifying face video of actors for plausible visual alignment to a dubbed audio track. In *Computer graphics forum*, Vol. 34. Wiley Online Library, 193–204.

[46] Chongjian Ge, Yibing Song, Yuying Ge, Han Yang, Wei Liu, and Ping Luo. 2021. Disentangled Cycle Consistency for Highly-realistic Virtual Try-On. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16928–16937.

[47] Yixiao Ge, Zhuowan Li, Haiyu Zhao, Guojun Yin, Shuai Yi, Xiaogang Wang, et al. 2018. Fd-gan: Pose-guided feature distilling gan for robust person re-identification. In *Advances in neural information processing systems*. 1222–1233.

[48] Yuying Ge, Yibing Song, Ruimao Zhang, Chongjian Ge, Wei Liu, and Ping Luo. 2021. Parser-Free Virtual Try-on via Distilling Appearance Flows. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8485–8493.

[49] Jiahao Geng, Tianjia Shao, Youyi Zheng, Yanlin Weng, and Kun Zhou. 2018. Warp-guided gans for single-photo facial animation. *ACM Transactions on Graphics (TOG)* 37, 6 (2018), 1–12.

[50] Hajer Ghodhbani, Mohamed Neji, Imran Razzak, and Adel M Alimi. 2022. You can try without visiting: a comprehensive survey on virtually try-on outfits. *Multimedia Tools and Applications* (2022), 1–32.

[51] Shiry Ginosar, Amir Bar, Gefen Kohavi, Caroline Chan, Andrew Owens, and Jitendra Malik. 2019. Learning individual styles of conversational gesture. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3497–3506.

[52] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *Advances in neural information processing systems*. 2672–2680.

[53] Artur Grigorev, Artem Sevastopolsky, Alexander Vakhitov, and Victor Lempitsky. 2019. Coordinate-Based Texture Inpainting for Pose-Guided Human Image Generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 12135–12144.

[54] Kuangxiao Gu, Yuqian Zhou, and Thomas S Huang. 2020. FLNet: Landmark Driven Fetching and Learning Network for Faithful Talking Facial Animation Synthesis. In *AAAI*. 10861–10868.

[55] Yudong Guo, Keyu Chen, Sen Liang, Yong-Jin Liu, Hujun Bao, and Juyong Zhang. 2021. Ad-nerf: Audio driven neural radiance fields for talking head synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5784–5794.

[56] Sungjoo Ha, Martin Kersner, Beomsu Kim, Seokjun Seo, and Dongyoung Kim. 2020. Marionette: Few-shot face reenactment preserving identity of unseen targets. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 10893–10900.

[57] Xintong Han, Xiaojun Hu, Weilin Huang, and Matthew R Scott. 2019. Clothflow: A flow-based model for clothed person generation. In *Proceedings of the IEEE International Conference on Computer Vision*. 10471–10480.

[58] Xintong Han, Zuxuan Wu, Weilin Huang, Matthew R Scott, and Larry S Davis. 2019. Finet: Compatible and diverse fashion image inpainting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4481–4491.

[59] Xintong Han, Zuxuan Wu, Zhe Wu, Ruichi Yu, and Larry S Davis. 2018. Viton: An image-based virtual try-on network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7543–7552.

[60] Hanxiang Hao, Sriram Baireddy, Amy R Reibman, and Edward J Delp. 2020. FaR-GAN for One-Shot Face Reenactment. *arXiv e-prints* (2020), arXiv–2005.

[61] Naomi Harte and Eoin Gillen. 2015. TCD-TIMIT: An audio-visual corpus of continuous speech. *IEEE Transactions on Multimedia* 17, 5 (2015), 603–615.

[62] Dai Hasegawa, Naoshi Kaneko, Shinichi Shirakawa, Hiroshi Sakuta, and Kazuhiko Sumi. 2018. Evaluation of speech-to-gesture generation using bi-directional LSTM network. In *Proceedings of the 18th International Conference on Intelligent Virtual Agents*. 79–86.

[63] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in neural information processing systems*. 6626–6637.

[64] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the Knowledge in a Neural Network. *arXiv e-prints* (2015), arXiv–1503.

[65] Wei-Lin Hsiao, Isay Katsman, Chao-Yuan Wu, Devi Parikh, and Kristen Grauman. 2019. Fashion++: Minimal edits for outfit improvement. In *Proceedings of the IEEE International Conference on Computer Vision*. 5047–5056.

[66] Chia-Wei Hsieh, Chieh-Yun Chen, Chien-Lung Chou, Hong-Han Shuai, and Wen-Huang Cheng. 2019. Fit-me: Image-based virtual try-on with arbitrary poses. In *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 4694–4698.

[67] Chia-Wei Hsieh, Chieh-Yun Chen, Chien-Lung Chou, Hong-Han Shuai, Jiaying Liu, and Wen-Huang Cheng. 2019. FashionOn: Semantic-guided Image-based Virtual Try-on with Detailed Human and Clothing Information. In *Proceedings of the 27th ACM International Conference on Multimedia*. 275–283.

[68] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. 2013. Human3. 6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE transactions on pattern analysis and machine intelligence* 36, 7 (2013), 1325–1339.

[69] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. 2017. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1125–1134.

[70] Thibaut Issenhuth, Jérémie Mary, and Clément Calauzènes. 2019. End-to-End Learning of Geometric Deformations of Feature Maps for Virtual Try-On. *arXiv e-prints* (2019), arXiv–1906.

[71] Thibaut Issenhuth, Jérémie Mary, and Clément Calauzènes. 2020. Do not mask what you do not need to mask: a parser-free virtual try-on. In *European Conference on Computer Vision*. Springer, 619–635.

[72] Max Jaderberg, Karen Simonyan, and Andrew Zisserman. 2015. Spatial transformer networks. In *Advances in neural information processing systems*. 2017–2025.

[73] Hyug Jae Lee, Rokkyu Lee, Minseok Kang, Myounghoon Cho, and Gunhan Park. 2019. LA-VITON: A Network for Looking-Attractive Virtual Try-On. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 0.

[74] Amir Jamaludin, Joon Son Chung, and Andrew Zisserman. 2019. You said that?: Synthesising talking faces from audio. *International Journal of Computer Vision* 127, 11 (2019), 1767–1779.

[75] Surgan Jandial, Ayush Chopra, Kumar Ayush, Mayur Hemani, Balaji Krishnamurthy, and Abhijeet Halwai. 2020. SieveNet: A Unified Framework for Robust Image-Based Virtual Try-On. In *The IEEE Winter Conference on Applications of Computer Vision*. 2182–2190.

[76] Nikolay Jetchev and Urs Bergmann. 2017. The conditional analogy gan: Swapping fashion articles on people images. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 2287–2292.

[77] Amina Kammoun, Rim Slama, Hedi Tabia, Tarek Ouni, and Mohmed Abid. 2022. Generative Adversarial Networks for face generation: A survey. *ACM Computing Surveys (CSUR)* (2022).

[78] Arnab Karmakar and Deepak Mishra. 2019. A Robust Pose Transformational GAN for Pose Guided Person Image Synthesis. In *National Conference on Computer Vision, Pattern Recognition, Image Processing, and Graphics*. Springer, 89–99.

[79] Tero Karras, Timo Aila, Samuli Laine, Antti Herva, and Jaakko Lehtinen. 2017. Audio-driven facial animation by joint end-to-end learning of pose and emotion. *ACM Transactions on Graphics (TOG)* 36, 4 (2017), 1–12.

[80] Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4401–4410.

[81] Jan Kietzmann, Linda W Lee, Ian P McCarthy, and Tim C Kietzmann. 2020. Deepfakes: Trick or treat? *Business Horizons* 63, 2 (2020), 135–146.

[82] Hyeongwoo Kim, Mohamed Elgharib, Michael Zollhöfer, Hans-Peter Seidel, Thabo Beeler, Christian Richardt, and Christian Theobalt. 2019. Neural style-preserving visual dubbing. *ACM Transactions on Graphics (TOG)* 38, 6 (2019), 1–13.

[83] Hyeongwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng Xu, Justus Thies, Matthias Niessner, Patrick Pérez, Christian Richardt, Michael Zollhöfer, and Christian Theobalt. 2018. Deep video portraits. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–14.

[84] Diederik P Kingma and Max Welling. 2013. Auto-Encoding Variational Bayes. *arXiv e-prints* (2013), arXiv–1312.

[85] Mohammad Rami Koujan, Michail Christos Doukas, Anastasios Roussos, and Stefanos Zafeiriou. 2020. Head2head: Video-based neural head synthesis. In *2020 15th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2020)*. IEEE, 16–23.

[86] Shizuma Kubo, Yusuke Iwasawa, and Yutaka Matsuo. 2018. Generative adversarial network-based virtual try-on with clothing region. (2018).

[87] Shizuma Kubo, Yusuke Iwasawa, Masahiro Suzuki, and Yutaka Matsuo. 2019. UVTON: UV Mapping to Consider the 3D Structure of a Human in Image-Based Virtual Try-On Network. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 0.

[88]  Taras Kucherenko, Dai Hasegawa, Naoshi Kaneko, Gustav Eje Henter, and Hedvig Kjellström. 2021. Moving fast and slow: Analysis of representations and post-processing in speech-driven automatic gesture generation. *International Journal of Human–Computer Interaction* 37, 14 (2021), 1300–1316.

[89]  Rithesh Kumar, Jose Sotelo, Kundan Kumar, Alexandre de Brebisson, and Yoshua Bengio. 2017. ObamaNet: Photo-realistic lip-sync from text. *arXiv e-prints* (2017), arXiv–1801.

[90]  Christoph Lassner, Gerard Pons-Moll, and Peter V Gehler. 2017. A generative model of people in clothing. In *Proceedings of the IEEE International Conference on Computer Vision*. 853–862.

[91]  Kedan Li, Min Jin Chong, Jingen Liu, and David Forsyth. 2020. Toward Accurate and Realistic Virtual Try-on Through Shape Matching and Multiple Warps. *arXiv e-prints* (2020), arXiv–2003.

[92]  Kun Li, Jinsong Zhang, Yebin Liu, Yu-Kun Lai, and Qionghai Dai. 2020. PoNA: Pose-guided non-local attention for human pose transfer. *IEEE Transactions on Image Processing* 29 (2020), 9584–9599.

[93]  Lincheng Li, Suzhen Wang, Zhimeng Zhang, Yu Ding, Yixing Zheng, Xin Yu, and Changjie Fan. 2021. Write-a-speaker: Text-based Emotional and Rhythmic Talking-head Generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 1911–1920.

[94]  Yining Li, Chen Huang, and Chen Change Loy. 2019. Dense intrinsic appearance flow for human pose transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3693–3702.

[95]  Dong Liang, Rui Wang, Xiaowei Tian, and Cong Zou. 2019. PCGAN: Partition-Controlled Human Image Generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 8698–8705.

[96]  Xiaodan Liang, Chunyan Xu, Xiaohui Shen, Jianchao Yang, Si Liu, Jinhui Tang, Liang Lin, and Shuicheng Yan. 2015. Human parsing with contextualized convolutional neural network. In *Proceedings of the IEEE international conference on computer vision*. 1386–1394.

[97]  Ji Liu, Heshan Liu, Mang-Tik Chiu, Yu-Wing Tai, and Chi-Keung Tang. 2020. Pose-Guided High-Resolution Appearance Transfer via Progressive Training. *arXiv e-prints* (2020), arXiv–2008.

[98]  Jinxian Liu, Bingbing Ni, Yichao Yan, Peng Zhou, Shuo Cheng, and Jianguo Hu. 2018. Pose transferrable person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4099–4108.

[99]  Kang Liu and Joern Ostermann. 2011. Realistic facial expression synthesis for an image-based talking head. In *2011 IEEE International Conference on Multimedia and Expo*. IEEE, 1–6.

[100]  Kuan-Hsien Liu, Ting-Yen Chen, and Chu-Song Chen. 2016. Mvc: A dataset for view-invariant clothing retrieval and attribute prediction. In *Proceedings of the 2016 ACM on International Conference on Multimedia Retrieval*. 313–316.

[101]  Lingjie Liu, Weipeng Xu, Michael Zollhoefer, Hyeongwoo Kim, Florian Bernard, Marc Habermann, Wenping Wang, and Christian Theobalt. 2019. Neural rendering and reenactment of human actor videos. *ACM Transactions on Graphics (TOG)* 38, 5 (2019), 1–14.

[102]  Meichen Liu, Kejun Wang, Juihang Ji, and Shuzhi Sam Ge. 2020. Person image generation with semantic attention network for person re-identification. *arXiv e-prints* (2020), arXiv–2008.

[103]  Meichen Liu, Kejun Wang, Ruihang Ji, Shuzhi Sam Ge, and Jing Chen. 2021. Pose transfer generation with semantic parsing attention network for person re-identification. *Knowledge-Based Systems* 223 (2021), 107024.

[104]  Meichen Liu, Xin Yan, Chenhui Wang, and Kejun Wang. 2020. Segmentation mask-guided person image generation. *Applied Intelligence* (2020), 1–16.

[105]  Ming-Yu Liu, Xun Huang, Jiahui Yu, Ting-Chun Wang, and Arun Mallya. 2021. Generative adversarial networks for image and video synthesis: Algorithms and applications. *Proc. IEEE* 109, 5 (2021), 839–862.

[106]  Wen Liu, Zhixin Piao, Jie Min, Wenhan Luo, Lin Ma, and Shenghua Gao. 2019. Liquid warping GAN: A unified framework for human motion imitation, appearance transfer and novel view synthesis. In *Proceedings of the IEEE International Conference on Computer Vision*. 5904–5913.

[107]  Yu Liu, Wei Chen, Li Liu, and Michael S Lew. 2019. SwapGAN: A Multistage Generative Approach for Person-to-Person Fashion Style Transfer. *IEEE Transactions on Multimedia* 21, 9 (2019), 2209–2222.

[108]  Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, and Xiaoou Tang. 2016. Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1096–1104.

[109]  Ziwei Liu, Sijie Yan, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2016. Fashion landmark detection in the wild. In *European Conference on Computer Vision*. Springer, 229–245.

[110]  Steven R Livingstone and Frank A Russo. 2018. The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English. *PloS one* 13, 5 (2018), e0196391.

[111]  Liqian Ma, Xu Jia, Qianru Sun, Bernt Schiele, Tinne Tuytelaars, and Luc Van Gool. 2017. Pose guided person image generation. In *Advances in Neural Information Processing Systems*. 406–416.

[112]  Liqian Ma, Qianru Sun, Stamatios Georgoulis, Luc Van Gool, Bernt Schiele, and Mario Fritz. 2018. Disentangled person image generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 99–108.

[113] Yifang Men, Yiming Mao, Yuning Jiang, Wei-Ying Ma, and Zhouhui Lian. 2020. Controllable person image synthesis with attribute-decomposed gan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5084–5093.

[114] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2020. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*. Springer, 405–421.

[115] Matiur Rahman Minar and Heejune Ahn. 2020. CloTH-VTON: Clothing Three-dimensional reconstruction for Hybrid image-based Virtual Try-ON. In *Proceedings of the Asian Conference on Computer Vision*.

[116] Matiur Rahman Minar, Thai Thanh Tuan, Heejune Ahn, Paul Rosin, and Yu-Kun Lai. 2020. 3D Reconstruction of Clothes using a Human Body Model and its Application to Image-based Virtual Try-On. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*.

[117] Matiur Rahman Minar, Thai Thanh Tuan, Heejune Ahn, Paul Rosin, and Yu-Kun Lai. 2020. CP-VTON+: Clothing Shape and Texture Preserving Image-Based Virtual Try-On. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*.

[118] Mehdi Mirza and Simon Osindero. 2014. Conditional Generative Adversarial Nets. *arXiv e-prints* (2014), arXiv–1411.

[119] Koki Nagano, Jaewoo Seo, Jun Xing, Lingyu Wei, Zimo Li, Shunsuke Saito, Aviral Agarwal, Jens Fursund, and Hao Li. 2018. paGAN: real-time avatars using dynamic textures. *ACM Transactions on Graphics (TOG)* 37, 6 (2018), 1–12.

[120] Arsha Nagrani, Joon Son Chung, and Andrew Zisserman. 2017. VoxCeleb: a large-scale speaker identification dataset. *arXiv e-prints* (2017), arXiv–1706.

[121] Assaf Neuberger, Eran Borenstein, Bar Hilleli, Eduard Oks, and Sharon Alpert. 2020. Image Based Virtual Try-On Network From Unpaired Data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5184–5193.

[122] Natalia Neverova, Riza Alp Guler, and Iasonas Kokkinos. 2018. Dense pose transfer. In *Proceedings of the European conference on computer vision (ECCV)*. 123–138.

[123] Yuval Nirkin, Yosi Keller, and Tal Hassner. 2019. FSGAN: Subject agnostic face swapping and reenactment. In *Proceedings of the IEEE international conference on computer vision*. 7184–7193.

[124] Kyle Olszewski, Zimo Li, Chao Yang, Yi Zhou, Ronald Yu, Zeng Huang, Sitao Xiang, Shunsuke Saito, Pushmeet Kohli, and Hao Li. 2017. Realistic dynamic facial textures from a single image using gans. In *Proceedings of the IEEE International Conference on Computer Vision*. 5429–5438.

[125] Xi Ouyang, Yu Cheng, Yifan Jiang, Chun-Liang Li, and Pan Zhou. 2018. Pedestrian-Synthesis-GAN: Generating Pedestrian Data in Real Scene and Beyond. *arXiv e-prints* (2018), arXiv–1804.

[126] Ivan Perov, Daiheng Gao, Nikolay Chervoniy, Kunlin Liu, Sugasa Marangonda, Chris Umé, Mr Dpfks, Carl Shift Facenheim, Luis RP, Jian Jiang, et al. 2020. DeepFaceLab: Integrated, flexible and extensible face-swapping framework. *arXiv e-prints* (2020), arXiv–2005.

[127] Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. 2019. MELD: A Multimodal Multi-Party Dataset for Emotion Recognition in Conversations. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 527–536.

[128] KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, and CV Jawahar. 2020. A Lip Sync Expert Is All You Need for Speech to Lip Generation In The Wild. In *Proceedings of the 28th ACM International Conference on Multimedia*. 484–492.

[129] Albert Pumarola, Antonio Agudo, Alberto Sanfeliu, and Francesc Moreno-Noguer. 2018. Unsupervised person image synthesis in arbitrary poses. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8620–8628.

[130] Xuelin Qian, Yanwei Fu, Tao Xiang, Wenxuan Wang, Jie Qiu, Yang Wu, Yu-Gang Jiang, and Xiangyang Xue. 2018. Pose-normalized image generation for person re-identification. In *Proceedings of the European conference on computer vision (ECCV)*. 650–667.

[131] Alec Radford, Luke Metz, and Soumith Chintala. 2015. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434* (2015).

[132] Amir Hossein Raffiee and Michael Sollami. 2021. Garmentgan: Photo-realistic adversarial fashion transfer. In *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 3923–3930.

[133] Amit Raj, Patsorn Sangkloy, Huiwen Chang, James Hays, Duygu Ceylan, and Jingwan Lu. 2018. Swapnet: Image based garment transfer. In *European Conference on Computer Vision*. Springer, 679–695.

[134] Bin Ren, Hao Tang, Fanyang Meng, Runwei Ding, Ling Shao, Philip HS Torr, and Nicu Sebe. 2021. Cloth Interactive Transformer for Virtual Try-On. *arXiv e-prints* (2021), arXiv–2104.

[135] Yurui Ren, Ge Li, Shan Liu, and Thomas H Li. 2020. Deep Spatial Transformation for Pose-Guided Person Image Generation and Animation. *IEEE Transactions on Image Processing* 29 (2020), 8622–8635.

[136] Yurui Ren, Xiaoming Yu, Junming Chen, Thomas H Li, and Ge Li. 2020. Deep image spatial transformation for person image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7690–7699.

[137] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. 2019. Faceforensics++: Learning to detect manipulated facial images. In *Proceedings of the IEEE International Conference on Computer Vision*. 1–11.

[138] Gemma Rotger, Felipe Lumbreras, Francese Moreno-Noguer, and Antonio Agudo. 2018. 2D-to-3D facial expression transfer. In *2018 24th International Conference on Pattern Recognition (ICPR)*. IEEE, 2008–2013.

[139] Debapriya Roy, Sanchayan Santra, and Bhabatosh Chanda. 2020. LGVTON: A Landmark Guided Approach to Virtual Try-On. *arXiv e-prints* (2020), arXiv–2004.

[140] Najmeh Sadoughi and Carlos Busso. 2019. Speech-driven expressive talking lips with conditional sequential generative adversarial networks. *IEEE Transactions on Affective Computing* (2019).

[141] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. 2019. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2304–2314.

[142] Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. 2020. Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 84–93.

[143] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. 2016. Improved techniques for training gans. In *Advances in neural information processing systems*. 2234–2242.

[144] Chenyang Si, Wei Wang, Liang Wang, and Tieniu Tan. 2018. Multistage adversarial losses for pose-based human image synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 118–126.

[145] Aliaksandr Siarohin, Stéphane Lathuilière, Enver Sangineto, and Nicu Sebe. 2019. Attention-based Fusion for Multi-source Human Image Generation. (2019).

[146] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. 2019. Animating arbitrary objects via deep motion transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2377–2386.

[147] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. 2019. First order motion model for image animation. *Advances in Neural Information Processing Systems* 32 (2019), 7137–7147.

[148] Aliaksandr Siarohin, Enver Sangineto, Stéphane Lathuilière, and Nicu Sebe. 2018. Deformable gans for pose-based human image generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3408–3416.

[149] Joon Son Chung, Arsha Nagrani, and Andrew Zisserman. 2018. VoxCeleb2: deep speaker recognition. *arXiv e-prints* (2018), arXiv–1806.

[150] Dan Song, Tianbao Li, Zhendong Mao, and An-An Liu. 2019. SP-VITON: shape-preserving image-based virtual try-on network. *Multimedia Tools and Applications* (2019), 1–13.

[151] Linsen Song, Wayne Wu, Chaoyou Fu, Chen Qian, Chen Change Loy, and Ran He. 2021. Everything's Talkin': Pareidolia Face Reenactment. *arXiv e-prints* (2021), arXiv–2104.

[152] Linsen Song, Wayne Wu, Chen Qian, Ran He, and Chen Change Loy. 2022. Everybody's talkin': Let me talk as you want. *IEEE Transactions on Information Forensics and Security* (2022).

[153] Sijie Song, Wei Zhang, Jiaying Liu, and Tao Mei. 2019. Unsupervised person image generation with semantic parsing transformation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2357–2366.

[154] Yang Song, Jingwen Zhu, Dawei Li, Andy Wang, and Hairong Qi. 2019. Talking face generation by conditional recurrent adversarial network. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*. 919–925.

[155] Feng Sun, Jiaming Guo, Zhuo Su, and Chengying Gao. 2019. Image-Based Virtual Try-on Network with Structural Coherence. In *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 519–523.

[156] Wei Sun, Jawadul H Bappy, Shanglin Yang, Yi Xu, Tianfu Wu, and Hui Zhou. 2019. Pose Guided Fashion Image Synthesis Using Deep Generative Model. *arXiv e-prints* (2019), arXiv–1906.

[157] Yang-Tian Sun, Hao-Zhi Huang, Xuan Wang, Yu-Kun Lai, Wei Liu, and Lin Gao. 2022. Robust pose transfer with dynamic details using neural video rendering. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2022).

[158] Supasorn Suwajanakorn, Steven M Seitz, and Ira Kemelmacher-Shlizerman. 2017. Synthesizing obama: learning lip sync from audio. *ACM Transactions on Graphics (TOG)* 36, 4 (2017), 1–13.

[159] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2818–2826.

[160] Hao Tang, Song Bai, Li Zhang, Philip HS Torr, and Nicu Sebe. 2020. Xinggan for person image generation. In *European Conference on Computer Vision*. Springer, 717–734.

[161] Hao Tang, Dan Xu, Gaowen Liu, Wei Wang, Nicu Sebe, and Yan Yan. 2019. Cycle in cycle generative adversarial networks for keypoint-guided image generation. In *Proceedings of the 27th ACM International Conference on Multimedia*. 2052–2060.

[162] Ayush Tewari, Ohad Fried, Justus Thies, Vincent Sitzmann, Stephen Lombardi, Kalyan Sunkavalli, Ricardo Martin-Brualla, Tomas Simon, Jason Saragih, Matthias Nießner, et al. 2020. State of the art on neural rendering. In *Computer Graphics Forum*, Vol. 39. Wiley Online Library, 701–727.

[163] Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. 2020. Neural voice puppetry: Audio-driven facial reenactment. In *European Conference on Computer Vision*. Springer, 716–731.

[164] Justus Thies, Michael Zollhöfer, and Matthias Nießner. 2019. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)* 38, 4 (2019), 1–12.

[165] Justus Thies, Michael Zollhöfer, Matthias Nießner, Levi Valgaerts, Marc Stamminger, and Christian Theobalt. 2015. Real-time expression transfer for facial reenactment. *ACM Transactions on Graphics (TOG)* 34, 6 (2015), 1–14.

[166] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. 2016. Face2face: Real-time face capture and reenactment of rgb videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2387–2395.

[167] Ruben Tolosana, Ruben Vera-Rodriguez, Julian Fierrez, Aythami Morales, and Javier Ortega-Garcia. 2020. Deepfakes and beyond: A survey of face manipulation and fake detection. *Information Fusion* 64 (2020), 131–148.

[168] Soumya Tripathy, Juho Kannala, and Esa Rahtu. 2020. Icface: Interpretable and controllable face reenactment using gans. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*. 3385–3394.

[169] Soumya Tripathy, Juho Kannala, and Esa Rahtu. 2021. Facegan: Facial attribute controllable reenactment gan. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 1329–1338.

[170] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. 2018. Mocogan: Decomposing motion and content for video generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1526–1535.

[171] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.

[172] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2017. Graph attention networks. *stat* 1050 (2017), 20.

[173] Luisa Verdoliva. 2020. Media forensics and deepfakes: an overview. *IEEE Journal of Selected Topics in Signal Processing* 14, 5 (2020), 910–932.

[174] Ruben Villegas, Jimei Yang, Yuliang Zou, Sungryull Sohn, Xunyu Lin, and Honglak Lee. 2017. Learning to generate long-term future via hierarchical prediction. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 3560–3569.

[175] Daniel Vlasic, Matthew Brand, Hanspeter Pfister, and Jovan Popovic. 2006. Face transfer with multilinear models. In *ACM SIGGRAPH 2006 Courses*. 24–es.

[176] Antonín Vobecký, Michal Uřičář, David Hurych, and Radoslav Škoviera. 2019. Advanced Pedestrian Dataset Augmentation for Autonomous Driving. In *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*. IEEE, 2367–2372.

[177] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. 2019. End-to-End Speech-Driven Realistic Facial Animation with Temporal GANs. In *CVPR Workshops*. 37–40.

[178] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. 2019. Realistic speech-driven facial animation with gans. *International Journal of Computer Vision* (2019), 1–16.

[179] Jacob Walker, Kenneth Marino, Abhinav Gupta, and Martial Hebert. 2017. The pose knows: Video forecasting by generating pose futures. In *Proceedings of the IEEE international conference on computer vision*. 3332–3341.

[180] Vincent Wan, Robert Anderson, Art Blokland, Norbert Braunschweiler, Langzhou Chen, BalaKrishna Kolluru, Javier Latorre, Ranniery Maia, Björn Stenger, Kayoko Yanagisawa, et al. 2013. Photo-realistic expressive text to talking head synthesis.. In *INTERSPEECH*. 2667–2669.

[181] Bochao Wang, Huabin Zheng, Xiaodan Liang, Yimin Chen, Liang Lin, and Meng Yang. 2018. Toward characteristic-preserving image-based virtual try-on network. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 589–604.

[182] Jiahang Wang, Tong Sha, Wei Zhang, Zhoujun Li, and Tao Mei. 2020. Down to the last detail: Virtual try-on with fine-grained details. In *Proceedings of the 28th ACM International Conference on Multimedia*. 466–474.

[183] Kaisiyuan Wang, Qianyi Wu, Linsen Song, Zhuoqian Yang, Wayne Wu, Chen Qian, Ran He, Yu Qiao, and Chen Change Loy. 2020. Mead: A large-scale audio-visual dataset for emotional talking-face generation. In *European Conference on Computer Vision*. Springer, 700–717.

[184] Ting-Chun Wang, Ming-Yu Liu, Andrew Tao, Guilin Liu, Jan Kautz, and Bryan Catanzaro. 2019. Few-shot video-to-video synthesis. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. 5013–5024.

[185] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. 2018. Video-to-Video Synthesis. *Advances in Neural Information Processing Systems* 31 (2018).

[186] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu. 2021. One-shot free-view neural talking-head synthesis for video conferencing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10039–10049.

[187] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* 13, 4 (2004), 600–612.

[188] Olivia Wiles, A Sophia Koepke, and Andrew Zisserman. 2018. X2face: A network for controlling face generation using images, audio, and pose codes. In *Proceedings of the European conference on computer vision (ECCV)*. 670–686.

[189] Kun Wu, Chengxiang Yin, Zhengping Che, Bo Jiang, Jian Tang, Zheng Guan, and Gangyi Ding. 2021. Human Pose Transfer with Disentangled Feature Consistency. *arXiv e-prints* (2021), arXiv–2107.

[190] Wayne Wu, Yunxuan Zhang, Cheng Li, Chen Qian, and Chen Change Loy. 2018. Reenactgan: Learning to reenact faces via boundary transfer. In *Proceedings of the European conference on computer vision (ECCV)*. 603–619.

[191] Zhonghua Wu, Guosheng Lin, Qingyi Tao, and Jianfei Cai. 2019. M2e-try on net: Fashion from model to everyone. In *Proceedings of the 27th ACM International Conference on Multimedia*. 293–301.

[192] Zhenyu Xie, Jianhuang Lai, and Xiaohua Xie. 2020. LG-VTON: Fashion Landmark Meets Image-Based Virtual Try-On. In *Chinese Conference on Pattern Recognition and Computer Vision (PRCV)*. Springer, 286–297.

[193] Chengming Xu, Yanwei Fu, Chao Wen, Ye Pan, Yu-Gang Jiang, and Xiangyang Xue. 2020. Pose-Guided Person Image Synthesis in the Non-Iconic Views. *IEEE Transactions on Image Processing* 29 (2020), 9060–9072.

[194] Yichao Yan, Jingwei Xu, Bingbing Ni, Wendong Zhang, and Xiaokang Yang. 2017. Skeleton-aided articulated motion generation. In *Proceedings of the 25th ACM international conference on Multimedia*. 199–207.

[195] Ceyuan Yang, Zhe Wang, Xinge Zhu, Chen Huang, Jianping Shi, and Dahua Lin. 2018. Pose guided human video generation. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 201–216.

[196] Han Yang, Ruimao Zhang, Xiaobao Guo, Wei Liu, Wangmeng Zuo, and Ping Luo. 2020. Towards Photo-Realistic Virtual Try-On by Adaptively Generating-Preserving Image Content. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7850–7859.

[197] Lingbo Yang, Pan Wang, Chang Liu, Zhanning Gao, Peiran Ren, Xinfeng Zhang, Shanshe Wang, Siwei Ma, Xiansheng Hua, and Wen Gao. 2021. Towards fine-grained human pose transfer with detail replenishing network. *IEEE Transactions on Image Processing* 30 (2021), 2422–2435.

[198] Lingbo Yang, Pan Wang, Xinfeng Zhang, Shanshe Wang, Zhanning Gao, Peiran Ren, Xuansong Xie, Siwei Ma, and Wen Gao. 2020. Region-adaptive texture enhancement for detailed person image synthesis. In *2020 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 1–6.

[199] Lingbo Yang, Zhenghui Zhao, Shiqi Wang, Shanshe Wang, Siwei Ma, and Wen Gao. 2019. Disentangled Human Action Video Generation via Decoupled Learning. In *2019 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*. IEEE, 495–500.

[200] Ran Yi, Zipeng Ye, Juyong Zhang, Hujun Bao, and Yong-Jin Liu. 2020. Audio-driven Talking Face Video Generation with Learning-based Personalized Head Pose. *arXiv e-prints* (2020), arXiv–2002.

[201] Gokhan Yildirim, Nikolay Jetchev, Roland Vollgraf, and Urs Bergmann. 2019. Generating high-resolution fashion model images wearing custom outfits. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 0–0.

[202] Weidong Yin, Ziwei Liu, and Leonid Sigal. 2021. Person-in-context synthesis with compositional structural space. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 2827–2836.

[203] Donggeun Yoo, Namil Kim, Sunggyun Park, Anthony S Paek, and In So Kweon. 2016. Pixel-level domain transfer. In *European Conference on Computer Vision*. Springer, 517–532.

[204] Jae Shin Yoon, Lingjie Liu, Vladislav Golyanik, Kripasindhu Sarkar, Hyun Soo Park, and Christian Theobalt. 2021. Pose-Guided Human Animation from a Single Image in the Wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15039–15048.

[205] Lingyun Yu, Jun Yu, and Qiang Ling. 2019. Mining audio, text and visual information for talking face generation. In *2019 IEEE International Conference on Data Mining (ICDM)*. IEEE, 787–795.

[206] Li Yu, Yueqi Zhong, and Xin Wang. 2019. Inpainting-Based Virtual Try-on Network for Selective Garment Transfer. *IEEE Access* 7 (2019), 134125–134136.

[207] Ruiyun Yu, Xiaoqi Wang, and Xiaohui Xie. 2019. Vtnfp: An image-based virtual try-on network with body and clothing feature preservation. In *Proceedings of the IEEE International Conference on Computer Vision*. 10511–10520.

[208] Polina Zablotskaia, Aliaksandr Siarohin, Bo Zhao, and Leonid Sigal. 2019. DwNet: Dense warp-based network for pose-guided human video generation. *arXiv e-prints* (2019), arXiv–1910.

[209] Egor Zakharov, Aleksei Ivakhnenko, Aliaksandra Shysheya, and Victor Lempitsky. 2020. Fast Bi-layer Neural Synthesis of One-Shot Realistic Head Avatars. In *European Conference on Computer Vision*. Springer, 524–540.

[210] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. 2019. Few-shot adversarial learning of realistic neural talking head models. In *Proceedings of the IEEE International Conference on Computer Vision*. 9459–9468.

[211] Mihai Zanfir, Alin-Ionut Popa, Andrei Zanfir, and Cristian Sminchisescu. 2018. Human appearance transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5391–5399.

[212] Xianfang Zeng, Yusu Pan, Mengmeng Wang, Jiangning Zhang, and Yong Liu. 2020. Realistic face reenactment via self-supervised disentangling of identity and pose. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 12757–12764.

[213] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 586–595.

[214] Weiyu Zhang, Menglong Zhu, and Konstantinos G Derpanis. 2013. From actemes to action: A strongly-supervised representation for detailed action understanding. In *Proceedings of the IEEE International Conference on Computer Vision*. 2248–2255.

[215] Yulei Zhang, Qingjie Zhao, and You Li. 2019. Multi-view Based Pose Alignment Method for Person Re-identification. In *Chinese Intelligent Automation Conference*. Springer, 439–447.

[216] Bo Zhao, Xiao Wu, Zhi-Qi Cheng, Hao Liu, Zequn Jie, and Jiashi Feng. 2018. Multi-view image generation from a single-view. In *Proceedings of the 26th ACM international conference on Multimedia*. 383–391.

[217] Long Zhao, Xi Peng, Yu Tian, Mubbasir Kapadia, and Dimitris Metaxas. 2018. Learning to forecast and refine residual motion for image-to-video generation. In *Proceedings of the European conference on computer vision (ECCV)*. 387–403.

[218] Wenbin Zhao, Qing Xie, Yanchun Ma, Yongjian Liu, and Shengwu Xiong. 2020. Pose Guided Person Image Generation Based on Pose Skeleton Sequence and 3D Convolution. In *2020 IEEE International Conference on Image Processing (ICIP)*. IEEE, 1561–1565.

[219] Haitian Zheng, Lele Chen, Chenliang Xu, and Jiebo Luo. 2019. Unsupervised Pose Flow Learning for Pose Guided Synthesis. *arXiv e-prints* (2019), arXiv–1909.

[220] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. 2015. Scalable person re-identification: A benchmark. In *Proceedings of the IEEE international conference on computer vision*. 1116–1124.

[221] Na Zheng, Xuemeng Song, Zhaozheng Chen, Linmei Hu, Da Cao, and Liqiang Nie. 2019. Virtually trying on new clothing with arbitrary poses. In *Proceedings of the 27th ACM International Conference on Multimedia*. 266–274.

[222] Zhedong Zheng, Xiaodong Yang, Zhiding Yu, Liang Zheng, Yi Yang, and Jan Kautz. 2019. Joint discriminative and generative learning for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2138–2147.

[223] Zerong Zheng, Tao Yu, Yebin Liu, and Qionghai Dai. 2021. Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction. *IEEE transactions on pattern analysis and machine intelligence* 44, 6 (2021), 3170–3184.

[224] Zhedong Zheng, Liang Zheng, and Yi Yang. 2017. Unlabeled samples generated by gan improve the person re-identification baseline in vitro. In *Proceedings of the IEEE International Conference on Computer Vision*. 3754–3762.

[225] Hang Zhou, Yu Liu, Ziwei Liu, Ping Luo, and Xiaogang Wang. 2019. Talking face generation by adversarially disentangled audio-visual representation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 9299–9306.

[226] Hang Zhou, Yasheng Sun, Wayne Wu, Chen Change Loy, Xiaogang Wang, and Ziwei Liu. 2021. Pose-controllable talking face generation by implicitly modularized audio-visual representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4176–4186.

[227] Mohan Zhou, Yalong Bai, Wei Zhang, Tiejun Zhao, and Tao Mei. 2021. Responsive Listening Head Generation: A Benchmark Dataset and Baseline. *arXiv preprint arXiv:2112.13548* (2021).

[228] Xingran Zhou, Siyu Huang, Bin Li, Yingming Li, Jiachen Li, and Zhongfei Zhang. 2019. Text guided person image synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3663–3672.

[229] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li. 2020. MakeItTalk: speaker-aware talking-head animation. *ACM Transactions on Graphics (TOG)* 39, 6 (2020), 1–15.

[230] Yipin Zhou, Zhaowen Wang, Chen Fang, Trung Bui, and Tamara Berg. 2019. Dance dance generation: Motion transfer for internet videos. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*. 0–0.

[231] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*. 2223–2232.

[232] Shizhan Zhu, Raquel Urtasun, Sanja Fidler, Dahua Lin, and Chen Change Loy. 2017. Be your own prada: Fashion synthesis with structural coherence. In *Proceedings of the IEEE international conference on computer vision*. 1680–1688.

[233] Zhen Zhu, Tengteng Huang, Baoguang Shi, Miao Yu, Bofei Wang, and Xiang Bai. 2019. Progressive pose attention transfer for person image generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2347–2356.