

# Size-Variable Virtual Try-On with Physical Clothes Size

Yohei Yamashita Chihiro Nakatani Norimichi Ukita  
Toyota Technological Institute  
{sd23501,ukita}@toyota-ti.ac.jp

## Abstract

This paper addresses a new virtual try-on problem of fitting any size of clothes to a reference person in the image domain. While previous image-based virtual try-on methods can produce highly natural try-on images, these methods fit the clothes on the person without considering the relative relationship between the physical sizes of the clothes and the person. Different from these methods, our method achieves size-variable virtual try-on in which the image size of the try-on clothes is changed depending on this relative relationship of the physical sizes. To relieve the difficulty in maintaining the physical size of the clothes while synthesizing the high-fidelity image of the whole clothes, our proposed method focuses on the residual between the silhouettes of the clothes in the reference and try-on images. We also develop a size-variable virtual try-on dataset consisting of 1,524 images provided by 26 subjects. Furthermore, we propose an evaluation metric for size-variable virtual try-on. Quantitative and qualitative experimental results show that our method can achieve size-variable virtual try-on better than general virtual try-on methods.

## 1. Introduction

The apparel industry is rapidly shifting to e-commerce, and its market size is expanding. However, customers cannot try on clothes online. The alternative is virtual try-on [10, 17, 35, 44]. We can try on any clothes in images. The quality of virtual try-on is improved by the success of image synthesis tasks such as image harmonization [8, 15, 33], image inpainting [28, 32, 40], and image editing [7, 22, 25].

In the 3D virtual try-on task [3, 31, 43], a 3D avatar [4, 29] generated from a user tries on any clothes by rendering the texture of the clothes on the avatar. However, generating a realistic 3D avatar for each user takes much time and cost.

In contrast to the 3D virtual try-on task, the 2D virtual try-on task only requires the images of the user and clothes. This image-based virtual try-on task can be categorized into (i) warping-based virtual try-on and (ii) Image Style Transfer-based (IST-based) virtual try-on.

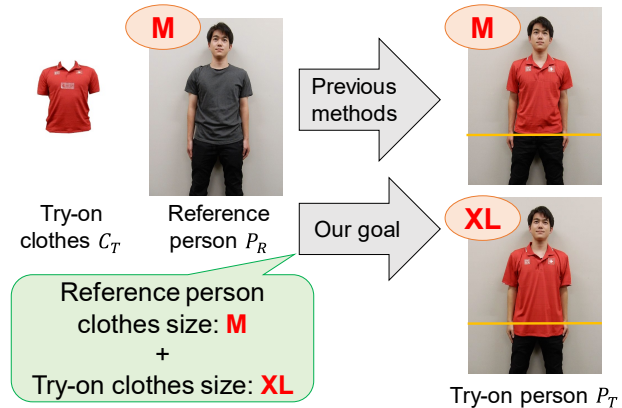


Figure 1. Comparison between previous methods and our method. (i) Previous methods [16, 39] only generate the try-on image in which the clothes size is the same as the one of the reference person. (ii) A user can change the try-on cloth size in our method.

Warping-based methods [16, 39] consist of two steps. First, the silhouette of the clothes is estimated as a mask in the images of the try-on clothes and the reference person (i.e., a user). The mask of the try-on clothes is warped to fit with that of the reference person's clothes without considering the physical size of the try-on clothes (e.g., small, medium, and large), as shown in the upper row of Fig. 1.

While the warping-based methods cannot control the size of the try-on clothes, IST-based methods [26] potentially allow us to control the clothes size by adjusting the weights of conditioning for image style transfer. However, it is not easy to appropriately control the clothes size (so that the try-on clothes fits with its physical size in the image) because the relationship between the conditioning weights and the physical size is unknown and highly complex.

In summary, the previous warping- and IST-based virtual try-on methods cannot reflect the physical size of try-on clothes in the try-on image. To address this issue, we propose size-variable virtual try-on with the physical size, as shown in the lower row of Fig. 1. Our contributions are as follows:

1. **Size-variable virtual try-on:** This paper defines a new problem, namely size-variable virtual try-on. Its goal is to fit the try-on clothes with the human body image by taking into account their physical sizes given as a user’s preference.
2. **Size-variable mask deformation network:** Size-variable virtual try-on is divided into two sub-tasks, size-variable mask deformation and texture rendering within the deformed mask. This paper focuses on the former (MDN: Mask Deformation Network in Fig. 2), while the latter is done with existing methods (TPS: Thin Plate Spline and CFN: Content Fusion Network in Fig. 2). Our MDN maintains the whole silhouette of the try-on clothes while adjusting its image size in accordance with its physical size given by a user. Our method achieves this silhouette maintenance and size adjustment by focusing on the residual between the cloth silhouettes in the reference and try-on images.
3. **Size evaluation:** Size-variable virtual try-on is a new problem, so we propose a new evaluation metric, namely the Size Evaluation Metric (SEM). SEM evaluates the size differences of hem and sleeve areas that are important for size-variable virtual try-on.
4. **Size-variable virtual try-on dataset:** We also develop a new dataset for size-variable virtual try-on.

## 2. Related Work

### 2.1. Warping-based Virtual Try-On

Figure 2 shows the two-stage pipeline of general warping-based methods [16, 38, 41]. (i) The mask of the try-on clothes ( $C_T$ ) is deformed to fit with the reference person image ( $P_R$ ). To make this deformation easier, the segmentation image and the person key-points are extracted as the auxiliary images from  $P_R$  and fed into MDN with  $C_T$ . (ii)  $C_T$  is warped to fit with the deformed mask ( $M_D$ ) by TPS [9], and then the warped try-on clothes ( $C_W$ ) and  $P_R$  are fused to produce the try-on person image ( $P_T$ ) by CFN.

However, TPS sometimes causes a large erroneous deformation on  $C_W$ . In [11, 39], this problem is relieved by segmenting  $P_T$  to the generated and original pixels so that the pixel values in the original pixels are copied from  $P_R$ . While these methods [11, 39] can preserve the quality of the original pixels, the quality of the generated pixels in  $P_T$  is degraded if these methods are applied to high-resolution images. VITON-HD [6] iteratively updates the segmentation image and increases the resolution of  $P_T$  for high-resolution virtual try-on. While such segmentation-based methods can be affected by erroneous segments, knowledge distillation-based methods [12, 18, 21] can generate  $P_T$  without segmentation.

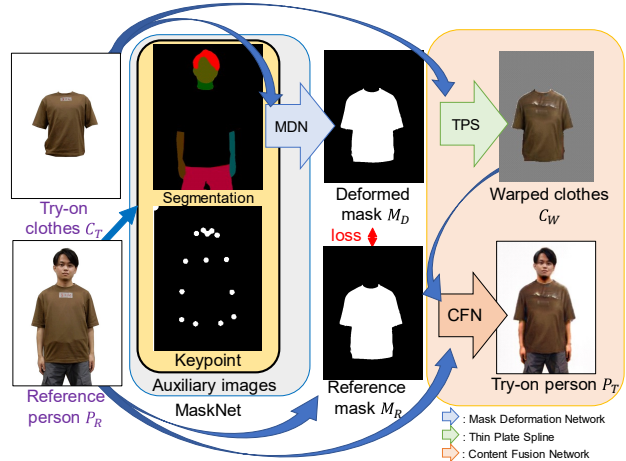


Figure 2. Overview of warping-based virtual try-on methods. The try-on mask is estimated from the auxiliary and clothes images. The clothes are warped to fit the mask and integrated with the reference person image to generate the try-on image.

In all of these methods, the physical size of the try-on clothes (e.g., hem and sleeve) is not explicitly addressed. Unlike these methods, our method estimates a size-variable try-on mask according to the physical size of the clothes.

### 2.2. Image Style Transfer-based Virtual Try-On

Image Style Transfer (IST) such as StyleGAN [1, 2, 23, 24] can generate images from the learned disentangled latent space. The disentangled latent space allows us to edit the specific regions of the generated image (e.g., hair length and eye color). StyleGAN is extended to virtual try-on in TryOnGAN [26]. TryOnGAN generates  $P_T$  by fusing disentangled style vectors representing the attributes of a person in an image and clothes in another image. However, TryOnGAN optimizes the style vectors of the clothes so that the try-on clothes in  $C_T$  fit with the body shape in  $P_R$  without taking into account the physical size of the clothes.

IST-based methods such as TryOnGAN can be extended to change the size of try-on clothes in  $P_T$  by changing noise given to the disentangled style vectors. However, the relationship between the image and physical sizes of try-on clothes is unknown. Different from such IST-based methods, our method estimates the try-on mask from the physical size of try-on clothes for size-variable virtual try-on.

## 3. Size-Variable Virtual Try-On

For size-variable virtual try-on, we collected a new dataset introduced in Sec. 3.1. Our proposed size-variable MDN is described in Sec. 3.2. Furthermore, we propose a new evaluation metric for size-variable virtual try-on (Sec. 3.3).

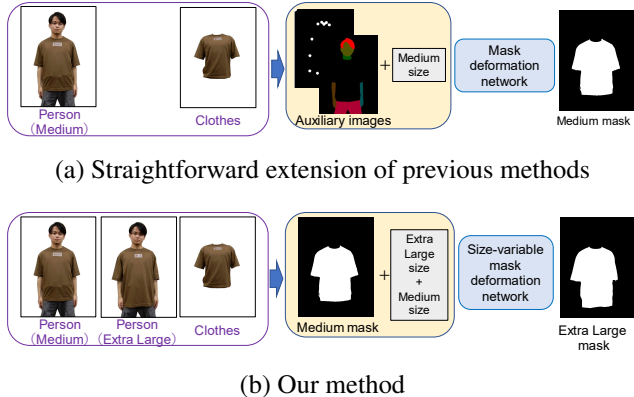


Figure 3. Two different approaches of mask deformation. (a) Extension of previous methods. The size of the try-on clothes is used for training MDN with auxiliary images in previous methods. (b) Our method. The paired person images in which each person wears different sizes of the same clothes are used in training.



Figure 4. Sample images of 14 clothes.

For size-variable virtual try-on, a straightforward scheme is to provide the physical size of try-on clothes as auxiliary cues to a previous virtual try-on method, as shown in Fig. 3 (a). However, such a straightforward scheme is not effective because only the size of clothes (as numerical parameters, which are indicated by “Parameters (Try-on size)” in Fig. 3 (a)) is not enough informative to generate the size-aware clothes mask. On the other hand, our method generates the size-aware clothes mask by adding the clothes image of the target size (which is indicated by “Person (Try-on size)” in Fig. 3 (b)) as well as the sizes of reference clothes and try-on clothes (as numerical parameters, which are denoted by “Parameters (Reference size)” and “Parameters (Try-on size),” respectively) in training. This training is achieved by conditioning our MDN by the preferred size of the clothes so that the MDN output coincides with the mask of the clothes of the preferred size.

### 3.1. Size-variable Virtual Try-on Dataset

**Motivation.** While the Zalando Dataset [16] has been widely used for the virtual try-on task, this dataset only contains the pairs of clothes and person images without the

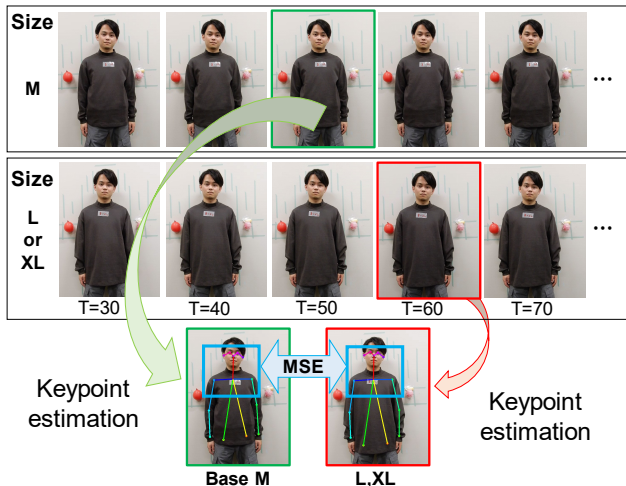


Figure 5. Posture matching for collecting image pairs, in each of which each subject’s postures are similar. This matching is done with the Mean Squared Error (MSE) computed between the sets of several body key-points.

clothes size. This disadvantage motivates us to collect a new dataset that contains the clothes size that can be used for size-variable virtual try-on.

**Overview.** In our dataset, each data is a pair of images in which observed clothes are the same except for physical size (e.g., “Person (Reference size)” and “Person (Try-on size)” in Fig. 3 (b)). Each image is annotated with the physical clothes size. The size parameters of each clothes are “Body length back,” “Sleeve length,” “Shoulder width,” “Body width,” and “Neck size.” A set of these size parameters is represented as a 5D vector. Our dataset is generated from 1,524 images of 26 subjects with 14 types of clothes shown in Fig. 4. From the 1,524 images, 3,746 image pairs are collected so that the person’s postures are almost the same in each image pair. All the paired data are split into 3,121 training, 529 validation, and 96 test data. In the 96 test data, there are 24 new clothes data, in which subjects wear new try-on clothes that are not included in the training data, and 72 new person data, in which subjects are excluded from those in the training data.

The distance from a camera to a subject was almost fixed in our dataset images. In real application scenarios, on the other hand, this distance may differ depending on the image capturing condition. This gap can be suppressed by rescaling/normalizing an input image in inference according to the ratio between the pixel and physical sizes (i.e., heights) of a user observed in the input image because the ratio in the dataset image is known.

**Dataset collection.** While it is required to spatially align the human postures between the two images in each image pair for our proposed method, this human posture align-

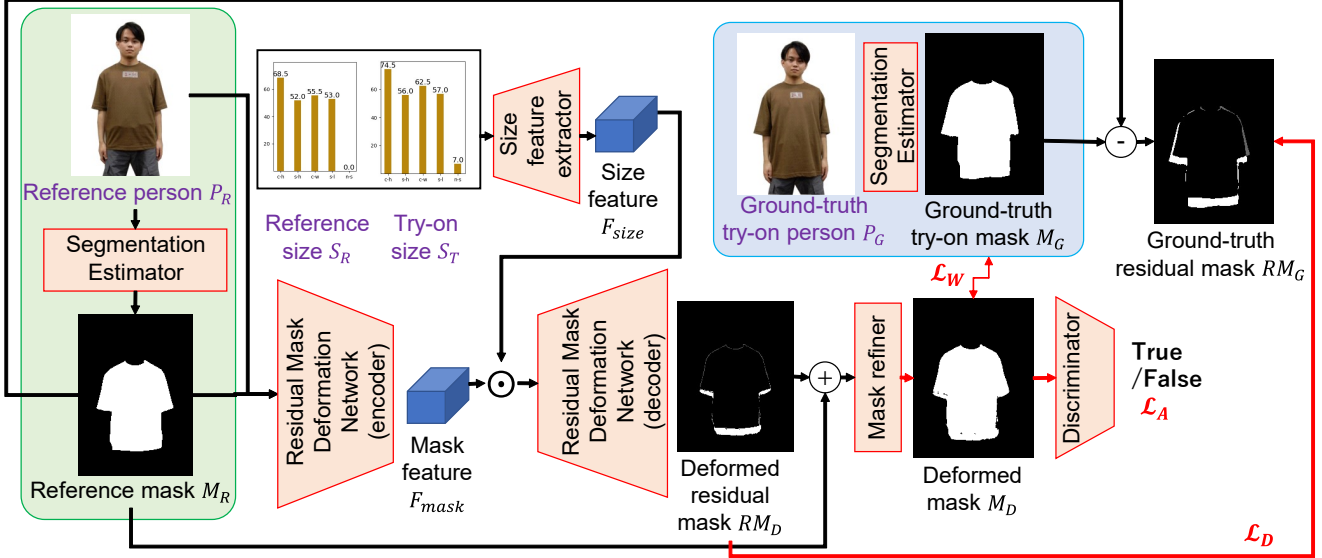


Figure 6. Overview of our size-variable mask deformation network. The final output mask ( $M_D$ ) is estimated from  $P_R$  and the sizes of the reference and try-on clothes ( $S_R$  and  $S_T$ , respectively). To focus on the small difference between the reference mask  $M_R$  and try-on mask  $M_G$ , the residual between these two masks is estimated as the intermediate output ( $RM_D$ ).  $RM_D$  is fed into the Mask Refiner (MR).

ment is not easy because it is difficult for a subject to be in the same posture in two shots. Therefore, we propose posture matching to collect the paired images in which the postures of a person are almost the same, as shown in Fig. 5. This posture matching consists of the following four steps. (i) Videos of each subject wearing different-sized same-type clothes are captured. The subject is requested to be in the same posture between the videos of the different-sized clothes. In each video pair (e.g., upper and lower videos in Fig. 5), the following steps (ii), (iii), and (iv) are done. (ii) The body key-points of the subject are estimated in all frames in the pair videos. In each of all possible frame pairs between the pair videos, the posture similarity between the sets of the key-points above the shoulders (which are within the blue rectangles in Fig. 5) is evaluated. This is because the key-points below the shoulders tend to differ even if the subject tries to be in the same posture. This posture similarity is measured as the Mean Squared Error (MSE) between the set of the key-points. (iii) The frames in which the MSE is smallest are selected as a matched image pair, as enclosed by the green and red rectangles in Fig. 5. We call the dataset collected by the above protocol “BaseDataset.” (iv) For further reducing the spatial displacement between the pair images in the BaseDataset, one of the pair images is warped in order to spatially align the sets of the key-points between the two images by projective transformation. This dataset is called “ProjDataset.”

### 3.2. Size-variable Mask Deformation Network

The detail of our proposed size-variable mask deformation network, which is “MDN” in Fig. 2, is shown in Fig. 6. For training this network, the pair images ( $P_R$  and  $P_G$ ) annotated with the sizes of clothes in these images ( $S_R$  and  $S_T$ ) are given, as described in Sec. 3.1. These inputs are fed into the network to estimate the deformed mask  $M_D$  used for warping-based virtual try-on, shown in Fig. 2.

**Architecture.** The image of reference person  $P_R \in \mathbb{R}^{3 \times H \times W}$  is fed into a segmentation estimator (SE) to obtain its reference mask  $M_R \in \mathbb{R}^{1 \times H \times W}$ .  $H$  and  $W$  denote the height and width of the image, respectively.  $P_R$  and  $M_R$  are fed into the Residual Mask Deformation Network (RMDN) to extract their feature map  $F_{mask} \in \mathbb{R}^{C \times H' \times W'}$  where  $C$ ,  $H'$ , and  $W'$  denote the dimension of  $F_{mask}$ . The size parameters of clothes in  $P_R$  and  $P_G$  (i.e.,  $S_R \in \mathbb{R}^S$  and  $S_T \in \mathbb{R}^S$  where  $S$  denotes the number of parameters representing the size of clothes) are fed into the Size Feature Extractor (SFE) to extract their feature map ( $F_{size} \in \mathbb{R}^{C \times H' \times W'}$ ). Then,  $F_{mask}$  and  $F_{size}$  are fused by element-wise multiplication. The fused feature map is decoded to estimate the residual mask  $RM_D \in \mathbb{R}^{1 \times H \times W}$ .  $RM_D$  and  $M_R$  are elementwise added to obtain the whole mask. Finally, the whole mask is fed into the Mask Refiner (MR) consisting of two convolutional layers to obtain the deformed mask  $M_D \in \mathbb{R}^{1 \times H \times W}$ . Our RMDN allows us to focus on the minor difference between  $M_R$  and  $M_G$ .

**Ground-truth Mask Generation.** The ground-truth try-on mask  $M_G \in \mathbb{R}^{1 \times H \times W}$  and the ground-truth residual mask

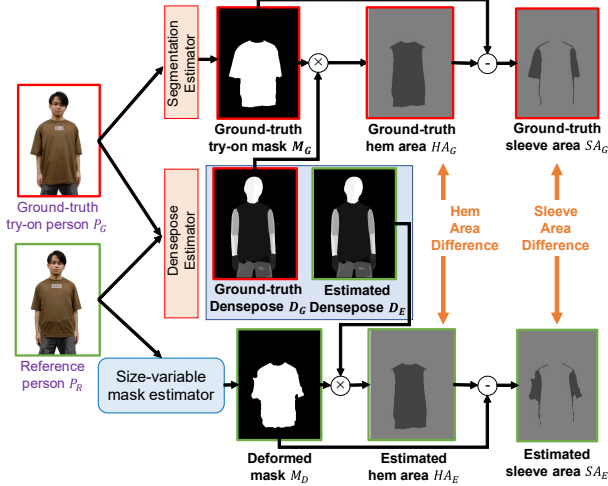


Figure 7. Size Evaluation Metric (SEM).

$RM_G \in \mathbb{R}^{1 \times H \times W}$  are generated from the ground-truth try-on person image  $P_G$ . As with  $M_R$ ,  $M_G$  is generated by SE.  $RM_G$  is generated from the elementwise subtraction between  $M_R$  and  $M_G$ . Both  $M_G$  and  $RM_G$  are used to train the whole network by the following loss functions.

**Loss functions.** The whole network is trained on the following loss functions:

$$\mathcal{L} = \lambda_W \mathcal{L}_W + \lambda_D \mathcal{L}_D + \lambda_A \mathcal{L}_A, \quad (1)$$

where  $\mathcal{L}_W$ ,  $\mathcal{L}_D$ , and  $\mathcal{L}_A$  denote the Weighted Binary Cross Entropy loss, the Dice loss [30], and the Adversarial loss [20], respectively.  $\lambda_W$ ,  $\lambda_D$ , and  $\lambda_A$  are their weights.  $\mathcal{L}_W$  is computed with  $M_D$  and  $M_G$ .  $\mathcal{L}_D$  is computed with  $RM_D$  and  $RM_G$  to focus on the minor difference between  $M_R$  and  $M_G$ .  $\mathcal{L}_A$  is computed with  $M_D$  and  $M_G$  so that  $M_D$  as a fake data gets close to  $M_G$  as a true data.  $\mathcal{L}_A$  can improve the reality of the boundary of  $M_D$ .

**Inference.** In inference,  $P_R$ ,  $S_R$ , and  $S_T$  are given and fed into our size-variable mask deformation network. Its output  $M_D$  is used to generate a warped clothes image  $C_W$  and a try-on image  $P_T$  by a general warping-based virtual try-on method, as described in Sec. 2.1 and shown in Fig. 2.

### 3.3. Size Evaluation Metric

While each of  $M_D$  and  $M_G$  is estimated as a heatmap image, the pixelwise difference between these two masks cannot be directly used for evaluating how much  $M_D$  looks like  $M_G$  because  $M_D$  and  $M_G$  are misaligned, as mentioned in Sec. 3.1. Furthermore, evaluation using all pixels cannot focus on the changes of areas important for size-variable virtual try-on (e.g., sleeves and hems). To resolve these problems, we propose the Size Evaluation Metric (SEM), as shown in Fig. 7.

For SEM, the torso hem and sleeves are considered to be areas important for size-variable virtual try-on. These areas are identified based on human body-part segments. These segments are detected by Densepose [14].  $D_R$  and  $D_G$  denote the Densepose heatmaps estimated from  $P_R$  and  $P_G$ , respectively. In the ground-truth image, its torso area ( $T_G$ ) is identified to be the pixelwise multiplication between  $M_G$  and the sum of the torso and upper-legs segments in  $D_G$ . In the same manner, the torso area in the reference image ( $T_D$ ) is identified with  $M_D$  and  $D_R$ . The sleeve areas are identified to be the pixelwise subtraction between the clothes mask and the torso area. The sleeve areas are denoted by  $S_G$  and  $S_D$ . The differences between “ $T_G$  and  $T_D$ ” and “ $S_G$  and  $S_D$ ” are calculated as follows:

$$T_- = \frac{1}{HW} \left| \sum_{i=1}^H \sum_{j=1}^W T_G(i, j) - \sum_{i=1}^H \sum_{j=1}^W T_D(i, j) \right| \quad (2)$$

$$S_- = \frac{1}{HW} \left| \sum_{i=1}^H \sum_{j=1}^W S_G(i, j) - \sum_{i=1}^H \sum_{j=1}^W S_D(i, j) \right| \quad (3)$$

The balance between  $T_-$  and  $S_-$  is quantified by their harmonic mean as follows:

$$SEM = \frac{2T_- S_-}{T_- + S_-} \quad (4)$$

The above SEM score gets smaller as the size-variable mask estimation works better.

## 4. Experiments

### 4.1. Implementation Details

We employ Graphonomy [13, 27] as SE. For estimating Densepose and key-points of each human body, Güler *et al.* [14] and Cao *et al.* [5] are used in our experiments, respectively. As for a warping-based virtual try-on method that accepts  $M_D$ , we used TPS and CFN in the pre-trained ACGPN [39]. All these components for our network are modularized so that they can be replaced with the SOTA methods without difficulty. SFE consists of three full-connection layers, each of which has ReLe activation. RMDN is an encoder-decoder network. The encoder and decode consist of four and five convolutional layers, respectively.

### 4.2. Evaluation Metrics and Dataset

**Evaluation metrics.**  $M_D$  estimated by our size-variable mask deformation network is evaluated with SEM proposed in Sec. 3.3. Furthermore, the try-on image  $P_T$  is also evaluated with Learned Perceptual Image Patch Similarity (LPIPS) [42] and Fréchet Inception Distance (FID) [19], both of which are widely used in the field of image generation to evaluate the perceptual image quality.

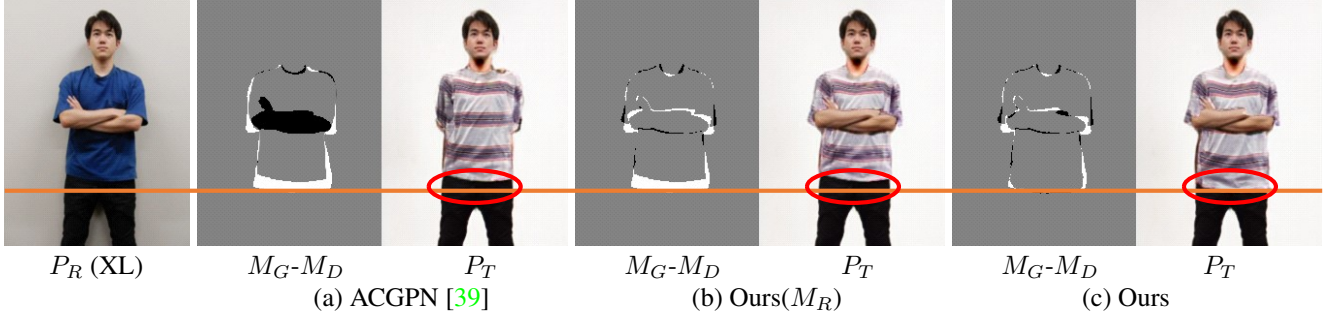


Figure 8. Visual comparison in comparative experiments. In this example, M-size and XL-size clothes are used as clothes in the reference person and try-on clothes images. In (a) and (b), the torso hem and sleeves are not changed from  $P_R$  to  $P_T$ . In (c), on the other hand, the torso hem and sleeves in  $P_T$  are extended in accordance with the physical size of the try-on clothes.

Table 1. Comparison of BaseDataset and ProjDataset.

Method	SEM( $\times 10^2$ ) $\downarrow$	LPIPS	FID
BaseDataset	0.49	0.44	16.83
ProjDataset (Ours)	0.42	0.44	16.25

Table 2. Quantitative comparison. The best and second-best results in each column are colored in red and blue. While StyleGAN is an IST-based method, other methods are warping-based methods.

Method	SEM( $\times 10^2$ ) $\downarrow$	LPIPS $\downarrow$	FID $\downarrow$
StyleGAN [23]	1.40	0.43	84.66
ACGPN [39]	1.03	0.45	35.10
Ours( $M_R$ )	0.73	0.44	15.77
Ours	0.42	0.44	16.25

**Dataset.** Our size-variable virtual try-on dataset, which is proposed in Sec. 3.1, is used. For training in all experiments, person images  $P_R$  and  $P_G$  are augmented by flip and rotation. The two types of datasets (i.e., BaseDataset and ProjDataset) are compared to validate the effectiveness of our proposed dataset generation method. For this validation, the performance of virtual try-on is considered to be the measure of the dataset quality. That is,  $M_D$  and  $P_T$  as the outputs of our method are evaluated by “SEM” and “LPIPS and FID,” respectively.

The results are shown in Table. 1. Since we can see that ProjDataset is better, ProjDataset is used in all experiments in what follows.

### 4.3. Comparative Experiments

Note that, since all existing virtual try-on methods have no function for changing the cloth size, it is impossible to show fair comparative experiments in terms of the performance on the virtual try-on task. However, to validate the necessity of the function for changing the cloth size,

our method is compared with ACGPN [39] as a general virtual try-on method that cannot change the cloth size. ACGPN is selected because (i) it is one of the SoTA methods for warping-based virtual try-on and (ii) since ACGPN and our method have the same networks of TPS and CFN, the difference between these two methods is how to produce the mask  $M_D$ . Therefore, a comparison between ACGPN and our method clearly validates the effectiveness of our proposed size-variable MDN. In addition to ACGPN, the effectiveness of our size-variable MDN is verified by using  $M_R$  instead of  $M_D$  deformed by the size-variable MDN as the input for TPS and CFN. This method is called Ours( $M_R$ ). While all the above methods are warping-based try-on methods, StyleGAN [23] as an IST-based method is also evaluated so that style parameters were manually optimized for changing the clothes size.

Quantitative results are shown in Table 2. Our method outperforms the other methods on SEM. The SEM scores of the previous methods [23, 39] are inferior to Ours because the physical clothes size is not directly given to those methods. This result demonstrates that our size-variable MDN can deform the mask according to the physical size of the clothes. This is the biggest goal of this work.

The visual quality of the generated virtual try-on images is also important. Ours and Ours( $M_R$ ) are the best on LPIPS, and Ours is the second-best on FID, while the gap between Ours( $M_R$ ) and Ours is not large (i.e., 15.77 vs. 16.25). The superiority of Ours( $M_R$ ) may be because the reality of  $M_D$  generated by our size-variable MDN cannot reach that of the real clothes silhouette (i.e.,  $M_R$ ) even if our size-variable MDN is optimized by the adversarial loss. FID of StyleGAN is much inferior to the other methods, probably because of the difficulty in hand-tuning of style parameters for changing the clothes size

Visual results are shown in Fig. 8. The orange auxiliary line is located along the torso hem of the clothes in  $P_R$ . Since the clothes sizes of the reference person image  $P_R$  and the try-on clothes are M-size and XL-size, respectively,

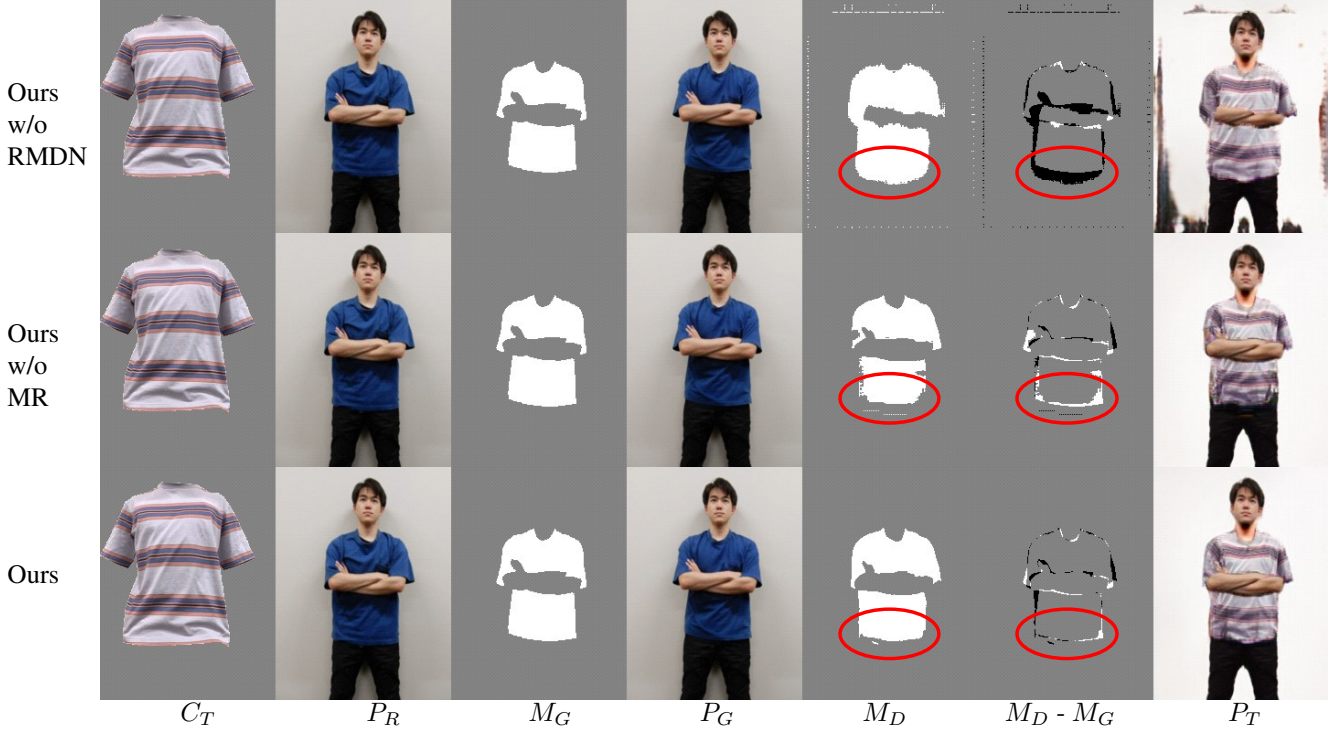


Figure 9. Visual results for validating the effectiveness of RMDN and MR in ablation studies.

Table 3. Ablation studies of our method.

Method	SEM( $\times 10^2$ ) $\downarrow$	LPIPS $\downarrow$	FID $\downarrow$
Ours w/o RMDN	1.21	0.54	67.51
Ours w/o MR	0.78	0.44	19.92
Ours w/o $P_R$	0.44	0.44	15.94
Ours	0.42	0.44	16.25

the torso hem and sleeves should be extended in the try-on image  $P_T$ . With (a) ACGPN [39] and (b) Ours( $M_R$ ), it can be seen that the torso hem and sleeves are not changed. (c) Our method, on the other hand, can extend the torso hem and sleeve in  $P_T$ . This result clearly demonstrates that our method can achieve size-variable virtual try-on in  $P_T$  different from the other methods.

#### 4.4. Ablation Study

The effectiveness of each component in our method is verified by ablating the following three components from our proposed method. We ablate the Residual Mask Deformation Network (RMDN) by directly estimating  $M_D$  from  $P_R$  and  $M_R$  without the residual connection. We also ablate the Mask Refiner (MR) in which  $RM_D$  is added with  $M_R$  to estimate  $M_D$  without the MR. Furthermore,  $P_R$  given to RMDN is also ablated.

The quantitative results are shown in Table 3. Our

method is the best on SEM and LPIPS. The improvements on SEM validate that RMDN, MR, and  $P_R$  certainly improve the virtual try-on quality in terms of the clothes size. In particular, the large gap on SEM between “Ours” and “Ours w/o RMDN” reveals that RMDN successfully estimates  $M_D$  from the clothes sizes. The difference between “Ours” and “Ours w/o MR” shows that the mask refiner allows fine adjustment of  $M_D$ , especially in the boundaries, rather than just elementwise adding  $RM_D$  to  $M_R$ .

As for FID, the difference between “Ours” and “Ours w/o  $P_R$ ” is not significant, while Ours is the second-best. This result is also demonstrated in visual results shown in Fig. 9. The first, second, and third rows show the results obtained by “Ours w/o RMDN,” “Ours w/o MR,” and “Ours,” respectively. Compared with the other methods, Ours can spatially align  $M_D$  with  $P_R$ , as shown in the regions enclosed by the red ellipses in  $M_D$  and  $M_D - M_G$ ; residual pixels in  $M_D - M_G$  are decreased as  $M_D$  becomes better. In “Ours w/o RMDN,” noisy lines appeared in  $M_D$  give a negative impact on  $P_T$ . Such noisy  $M_D$  degrades the perceptual quality of  $P_T$ , as also shown in Table 3.

#### 4.5. Detailed Analysis

**Comparison of loss functions for comparing  $M_D$  with  $M_G$ :** While the Weighted Binary Cross Entropy loss,  $\mathcal{L}_W$ , is employed for comparing  $M_D$  with  $M_G$  as expressed in (1) in our method,  $\mathcal{L}_W$  is replaced by the Binary Cross

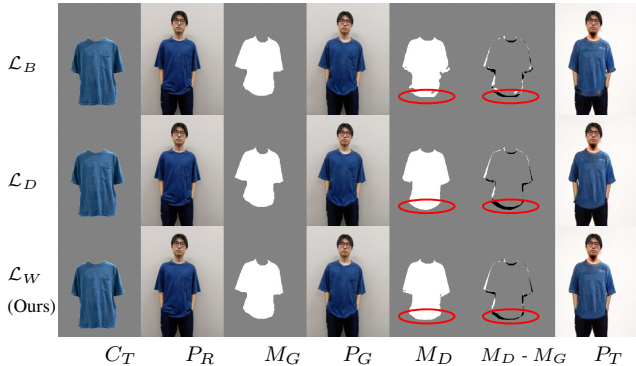


Figure 10. Visual comparison of loss functions for comparing  $M_D$  with  $M_G$ .

Table 4. Comparison of loss functions for comparing  $M_D$  with  $M_G$ .

Loss function	SEM( $\times 10^2$ ) $\downarrow$	LPIPS $\downarrow$	FID $\downarrow$
$\mathcal{L}_B$	0.46	0.44	16.48
$\mathcal{L}_D$	0.45	0.44	17.27
$\mathcal{L}_W$ (Ours)	0.42	0.44	16.25

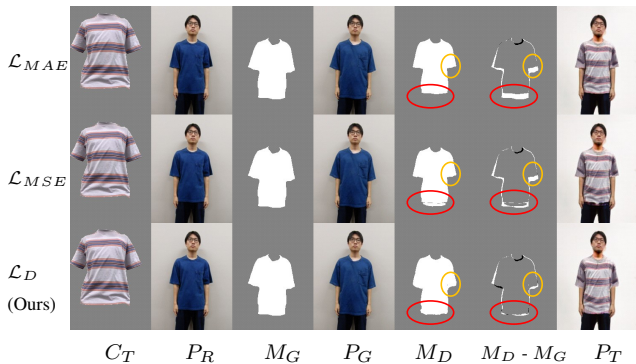


Figure 11. Visual comparison of loss functions for comparing  $RM_D$  with  $RM_G$ .

Entropy loss,  $\mathcal{L}_B$ , or the Dice loss,  $\mathcal{L}_D$  for comparative experiments. The quantitative results are shown in Table 4. Our method with  $\mathcal{L}_W$  achieves the best performance in all metrics. This result shows that  $\mathcal{L}_W$  can properly weight zero and non-zero pixels in  $M_D$  and  $M_G$  for adjusting the clothes size while maintaining the visual reality of the clothes silhouette. The visual comparison is shown in Fig. 10. Compared with the other methods,  $\mathcal{L}_W$  (Ours) successfully estimates  $M_D$ , as shown in the regions enclosed by the red ellipses in  $M_D$  and  $M_D - M_G$ .

**Comparison of loss functions for comparing  $RM_D$  with  $RM_G$ :** The loss function for comparing  $RM_D$  with  $RM_G$ , the Dice loss  $\mathcal{L}_D$  in our method is replaced by alternative loss functions.  $\mathcal{L}_D$  is replaced by the Mean Absolute Error

Table 5. Comparison of loss functions for  $RM_D$  and  $RM_G$ . While  $\mathcal{L}_{Dice}$  is used for  $RM_D$  and  $RM_G$  in our method,  $\mathcal{L}_{MAE}$  and  $\mathcal{L}_{MSE}$  are compared in this experiment.

Loss function	SEM( $\times 10^2$ ) $\downarrow$	LPIPS $\downarrow$	FID $\downarrow$
$\mathcal{L}_{MAE}$	0.72	0.44	15.30
$\mathcal{L}_{MSE}$	0.51	0.44	15.82
$\mathcal{L}_{Dice}$ (Ours)	0.42	0.44	16.25

ror  $\mathcal{L}_{MAE}$  or the Mean Squared Error  $\mathcal{L}_{MSE}$ . As shown in Table 5, our method outperforms the other methods on SEM, as with all the other experiments shown before. The big improvements in the clothes size can also be validated in the visual results, as shown in the red and orange ellipses in Fig. 11. Regarding LPIPS and FID, the LPIPS scores of all the methods are equal, while ours is inferior to the others. However, it is difficult to see any remarkable difference in the visual results, as shown in Fig. 11.

## 5. Conclusion

This paper proposed the size-variable mask deformation network, the size-variable virtual try-on dataset, and the size evaluation metric for size-variable virtual try-on, which is a new problem in this research area. Our proposed mask deformation network can estimate the mask in accordance with the physical size of the try-on clothes. The results of size-variable virtual try-on are evaluated by our size evaluation metric in which the lengths of the torso hem and sleeves are particularly evaluated. Experimental results demonstrate that our method outperforms the other methods quantitatively regarding the performance of size-variable virtual try-on; our method is the best on the size evaluation metric (SEM), with a large margin improvement in all the experiments. In the visual quality evaluated by LPIPS and FID also, our method is comparable with the other methods. Furthermore, various visual results also validate the effectiveness of our proposed method.

Extending our dataset is important future work because its scale is not sufficient yet to validate the performance with more subjects, more poses, and more clothes, while our dataset is the first one that can be used for the size-variable virtual try-on task. While a general pose estimation method [5] is used in our experiments, for key-point estimation under clothes, human pose estimation should be optimized for this purpose [36, 37]. Active learning also benefits efficient and accurate pose estimation in a query video [34]. Another future work is to verify the modularity of our method by replacing a component/sub-network. For example, it is not guaranteed that existing TPS can always generate authentic results for try-on clothes. Since our size-variable mask deformation network is modular, TPS can be replaced by any SoTA texture rendering method.



## References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *ICCV*, 2019. 2
- [2] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan++: How to edit the embedded images? In *CVPR*, 2020. 2
- [3] Bharat Lal Bhatnagar, Garvita Tiwari, Christian Theobalt, and Gerard Pons-Moll. Multi-garment net: Learning to dress 3d people from images. In *ICCV*, 2019. 1
- [4] Bharat Lal Bhatnagar, Garvita Tiwari, Christian Theobalt, and Gerard Pons-Moll. Multi-garment net: Learning to dress 3d people from images. In *ICCV*, 2019. 1
- [5] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, 2017. 5, 8
- [6] Seunghwan Choi, Sunghyun Park, Minsoo Lee, and Jaegul Choo. VITON-HD: high-resolution virtual try-on via misalignment-aware normalization. *CoRR*, abs/2103.16874, 2021. 2
- [7] Edo Collins, Raja Bala, Bob Price, and Sabine Süssstrunk. Editing in style: Uncovering the local semantics of gans. In *CVPR*, 2020. 1
- [8] Wenyan Cong, Jianfu Zhang, Li Niu, Liu Liu, Zhixin Ling, Weiyuan Li, and Liqing Zhang. Dovenet: Deep image harmonization via domain verification. In *CVPR*, 2020. 1
- [9] Jean Duchon. Splines minimizing rotation-invariant seminorms in sobolev spaces. In *Constructive Theory of Functions of Several Variables*, 1976. 2
- [10] Jun Ehara and Hideo Saito. Texture overlay for virtual clothing based on PCA of silhouettes. In *ISMAR*, 2006. 1
- [11] Chongjian Ge, Yibing Song, Yuying Ge, Han Yang, Wei Liu, and Ping Luo. Disentangled cycle consistency for highly-realistic virtual try-on. In *CVPR*, 2021. 2
- [12] Yuying Ge, Yibing Song, Ruimao Zhang, Chongjian Ge, Wei Liu, and Ping Luo. Parser-free virtual try-on via distilling appearance flows. In *CVPR*, 2021. 2
- [13] Ke Gong, Yiming Gao, Xiaodan Liang, Xiaohui Shen, Meng Wang, and Liang Lin. Graphonomy: Universal human parsing via graph transfer learning. In *CVPR*, 2019. 5
- [14] Riza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation in the wild. In *CVPR*, 2018. 5
- [15] Zonghui Guo, Haiyong Zheng, Yufeng Jiang, Zhaorui Gu, and Bing Zheng. Intrinsic image harmonization. In *CVPR*, 2021. 1
- [16] Xintong Han, Zuxuan Wu, Zhe Wu, Ruichi Yu, and Larry S. Davis. VITON: an image-based virtual try-on network. In *CVPR*, 2018. 1, 2, 3
- [17] Stefan Hauswiesner, Matthias Straka, and Gerhard Reitmayr. Virtual try-on through image-based rendering. *IEEE Trans. Vis. Comput. Graph.*, 19(9):1552–1565, 2013. 1
- [18] Sen He, Yi-Zhe Song, and Tao Xiang. Style-based global appearance flow for virtual try-on. In *CVPR*, 2022. 2
- [19] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, 2017. 5
- [20] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. 5
- [21] Thibaut Issenhuth, Jérémie Mary, and Clément Calauzènes. Do not mask what you do not need to mask: A parser-free virtual try-on. In *ECCV*, 2020. 2
- [22] Youngjoo Jo and Jongyoul Park. SC-FEGAN: face editing generative adversarial network with user’s sketch and color. In *ICCV*, 2019. 1
- [23] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. *CoRR*, abs/1812.04948, 2018. 2, 6
- [24] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *CVPR*, 2020. 2
- [25] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In *CVPR*, 2020. 1
- [26] Kathleen M. Lewis, Srivatsan Varadharajan, and Ira Kemelmacher-Shlizerman. Tryongan: body-aware try-on via layered interpolation. *ACM Trans. Graph.*, 40(4):115:1–115:10, 2021. 1, 2
- [27] Peike Li, Yunqiu Xu, Yunchao Wei, and Yi Yang. Self-correction for human parsing. *CoRR*, abs/1910.09777, 2019. 5
- [28] Wenbo Li, Zhe Lin, Kun Zhou, Lu Qi, Yi Wang, and Jiaya Jia. MAT: mask-aware transformer for large hole image inpainting. In *CVPR*, 2022. 1
- [29] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: a skinned multi-person linear model. *ACM Trans. Graph.*, 34(6):248:1–248:16, 2015. 1
- [30] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *3DV*, 2016. 5
- [31] Aymen Mir, Thiemo Alldieck, and Gerard Pons-Moll. Learning to transfer texture from clothing images to 3d humans. In *CVPR*, 2020. 1
- [32] Kamyar Nazeri, Eric Ng, Tony Joseph, Faisal Z. Qureshi, and Mehran Ebrahimi. Edgeconnect: Structure guided image inpainting using edge prediction. In *ICCV-W*, 2019. 1
- [33] Konstantin Sofiiuk, Polina Popenova, and Anton Konushin. Foreground-aware semantic representations for image harmonization. In *WACV*, 2021. 1
- [34] Hiromu Taketsugu and Norimichi Ukita. Active transfer learning for efficient video-specific human pose estimation. In *WACV*, 2024. 8
- [35] Hiroshi Tanaka and Hideo Saito. Texture overlay onto flexible object with PCA of silhouettes and k-means method for search into database. In *MVA*, 2009. 1
- [36] Norimichi Ukita, Michiro Hirai, and Masatsugu Kidode. Complex volume and pose tracking with probabilistic dynamical models and visual hull constraints. In *ICCV*, 2009. 8

- [37] Norimichi Ukita, Ryosuke Tsuji, and Masatsugu Kidode. Real-time shape analysis of a human body in clothing using time-series part-labeled volumes. In *ECCV*, 2008. 8
- [38] Bochao Wang, Huabin Zheng, Xiaodan Liang, Yimin Chen, Liang Lin, and Meng Yang. Toward characteristic-preserving image-based virtual try-on network. In *ECCV*, 2018. 2
- [39] Han Yang, Ruimao Zhang, Xiaobao Guo, Wei Liu, Wangmeng Zuo, and Ping Luo. Towards photo-realistic virtual try-on by adaptively generating↔preserving image content. In *CVPR*, 2020. 1, 2, 5, 6, 7
- [40] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S. Huang. Free-form image inpainting with gated convolution. In *ICCV*, 2019. 1
- [41] Ruiyun Yu, Xiaoqi Wang, and Xiaohui Xie. VTNFP: an image-based virtual try-on network with body and clothing feature preservation. In *ICCV*, 2019. 2
- [42] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 5
- [43] Fuwei Zhao, Zhenyu Xie, Michael Kampffmeyer, Haoye Dong, Songfang Han, Tianxiang Zheng, Tao Zhang, and Xiaodan Liang. M3D-VTON: A monocular-to-3d virtual try-on network. In *ICCV*, 2021. 1
- [44] Zhenglong Zhou, Bo Shu, Shaojie Zhuo, Xiaoming Deng, Ping Tan, and Stephen Lin. Image-based clothes animation for virtual fitting. In *SIGGRAPH*, 2012. 1