

# PISE: Person Image Synthesis and Editing with Decoupled GAN

Jinsong Zhang<sup>1</sup>, Kun Li<sup>1\*</sup>, Yu-Kun Lai<sup>2</sup>, Jingyu Yang<sup>1</sup>  
<sup>1</sup>Tianjin University, China <sup>2</sup> Cardiff University, United Kingdom

## Abstract

Person image synthesis, e.g., pose transfer, is a challenging problem due to large variation and occlusion. Existing methods have difficulties predicting reasonable invisible regions and fail to decouple the shape and style of clothing, which limits their applications on person image editing. In this paper, we propose PISE, a novel two-stage generative model for Person Image Synthesis and Editing, which is able to generate realistic person images with desired poses, textures, or semantic layouts. For human pose transfer, we first synthesize a human parsing map aligned with the target pose to represent the shape of clothing by a parsing generator, and then generate the final image by an image generator. To decouple the shape and style of clothing, we propose joint global and local per-region encoding and normalization to predict the reasonable style of clothing for invisible regions. We also propose spatial-aware normalization to retain the spatial context relationship in the source image. The results of qualitative and quantitative experiments demonstrate the superiority of our model on human pose transfer. Besides, the results of texture transfer and region editing show that our model can be applied to person image editing. The code is available for research purposes at <https://github.com/Zhangjinso/PISE>.

## 1. Introduction

Person image synthesis is a challenging problem in computer vision and computer graphics, which has great application potentials in image editing, video generation, virtual try-on, etc. Human pose transfer [16, 20, 23, 24, 32], i.e., synthesizing a new image for the same person in a target pose, is an active topic in person image synthesis.

Recently, Generative Adversarial Networks (GANs) [4] achieve great success in human pose transfer. Many methods directly learn the mapping from the source image and pose to the target image using neural networks [32, 23, 24, 12]. Most of these methods utilize a two-branch (pose branch and image branch) framework to transfer the feature of the source image from the source pose to the target pose. However, by taking keypoints as the pose representation,

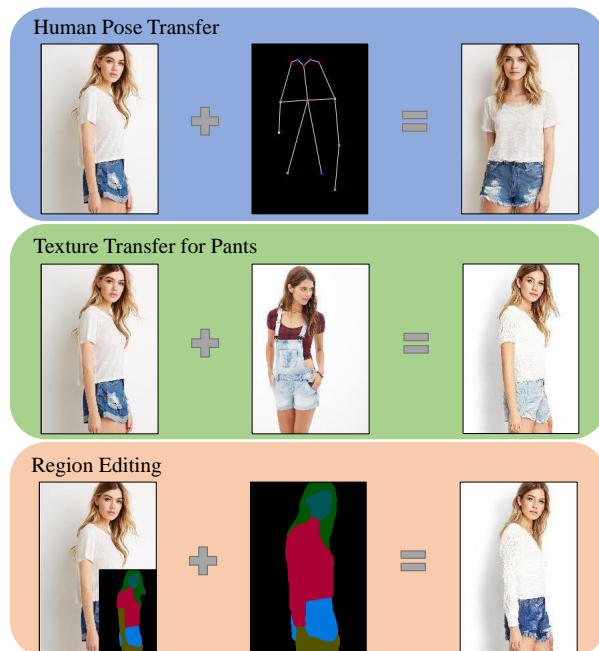


Figure 1. Our model PISE allows to transfer new pose or texture to a single person image, and also enables region editing.

it is difficult to predict a sharp and reasonable image with sparse correspondences when the source pose and the target pose have large differences. To deal with this problem, flow-based methods [13, 20] estimate an appearance flow to obtain denser correspondences, which is used to warp the source image or its feature to align with the target pose. The final image delivered by refining the warped image or decoding the warped image feature is a rearrangement of the source image elements. Thus, the generated image can preserve details of the source image, but the invisible region due to occlusion is not satisfactorily recovered. To predict invisible regions, some methods [16, 29] introduce human parsing maps to human pose transfer. The human parsing map provides semantic correspondence to synthesize the final image and enables applications of person image editing. However, these methods cannot disentangle the shape and style information (e.g., the category and texture of clothing) and fail to preserve spatial context relationships. For

\*Corresponding author

flexible and detailed editing, it is better to disentangle the shape and style. Meanwhile, preserving spatial information of the source image and reasonable prediction of invisible regions are also important for producing the desired output for human pose transfer.

The aforementioned methods encounter three challenges to synthesize satisfactory images: 1) the coupling of the shape and style of clothing, 2) potential uncertainties in invisible regions, and 3) loss of spatial context relationships.

To address these problems, we propose a novel Decoupled GAN for person image synthesis and editing. Instead of directly learning a mapping from the source to the target, we take the human parsing map as the intermediate result to provide semantic guidance to predict a reasonable shape of clothing. We propose joint global and local per-region encoding and normalization to control the texture style on a semantic region basis, better utilizing information for both visible and invisible regions. Specifically, for the region visible in the source image, we use the local feature of the corresponding region to predict the style of clothing. For the region invisible in the source image but visible in the target image, we obtain the global feature of the source image to predict the reasonable style of clothing, which can well deal with the generation of invisible regions. Benefiting from the human parsing map and per-region texture control, the shape and the style of clothing are disentangled for more flexible editing. Besides, to preserve the spatial context relationship, we propose a novel spatial-aware normalization to transfer spatial information of the source image to the generated image. After per-region normalization and spatial-aware normalization, the generated target feature passes through a decoder to output the final image. Figure 1 shows some applications of our model.

The main contributions of this work are summarized as follows:

- We propose a two-stage model with per-region control to decouple the shape and style of clothing. Experimental results on human pose transfer, texture transfer, and region editing show the flexibility and superior performance of our person image synthesis and editing method.
- We propose joint global and local per-region encoding and normalization to predict the reasonable style of clothing for invisible regions, and preserve the original style of clothing in the target image.
- We propose a spatial-aware normalization to retain the spatial context relationship in the source image, and transfer it by modulating the scale and bias of the generated image feature.

## 2. Related Work

### 2.1. Image Synthesis

In recent years, Generative Adversarial Networks (GANs) [4] have made great success in image synthesis [28, 27, 21]. Isola *et al.* [9] first introduced conditional GANs [17] to solve the image-to-image generation task, which was extended to high-resolution image synthesis [25]. Zhu *et al.* proposed an unsupervised method with cycle consistency to transfer images between two domains without paired data. StyleGAN [11] used adaptive instance normalization (AdaIN) [8] to achieve scale-specific control of image synthesis. SPADE [19] adopted spatially-adaptive normalization to synthesize new images given semantic input by modulating the activations in normalization layers. Zhu *et al.* [33] leveraged group convolution and designed a Group Decreasing Network to alleviate memory access cost problem. SEAN [31] improved SPADE by proposing per-region encoding to control the style of individual regions in the generated images. However, these methods have limited editable capacity in human pose transfer due to sparse correspondence of keypoints and large variation in pose and texture. In this paper, we propose a two-stage model to obtain semantic guidance to achieve more controllable image synthesis.

### 2.2. Human Pose Transfer

Human pose transfer is a highly active topic in computer vision and computer graphics. PG<sup>2</sup> [15] firstly introduced this problem and utilized a coarse-to-fine framework to alleviate the challenging generation problem. It concatenated the source image, the source pose, and the target pose as inputs to learn the target image, which leads to feature misalignment. Some methods used a two-branch framework with the image branch and pose branch to deal with the misalignment between the source and target images. Zhu *et al.* [32] proposed to transfer image information from the source pose to the target pose progressively with a local attention mechanism. Tang *et al.* [24] modeled appearance information and shape information with two novel blocks. Li *et al.* [12] designed pre-posed image-guided pose feature update and post-posed pose guided image feature update to better utilize the pose and image features. These methods used keypoints as their pose representation and focused on transferring image information with the guidance of pose information. Therefore, only sparse correspondence between the source image and the target image can be obtained, which is difficult to transfer image information from the source pose to the target pose. To provide semantic correspondence between the source image and the target image, some methods introduced the human parsing map as semantic guidance. Dong *et al.* [2] used a two-stage model, which first synthesized a target semantic segmentation map

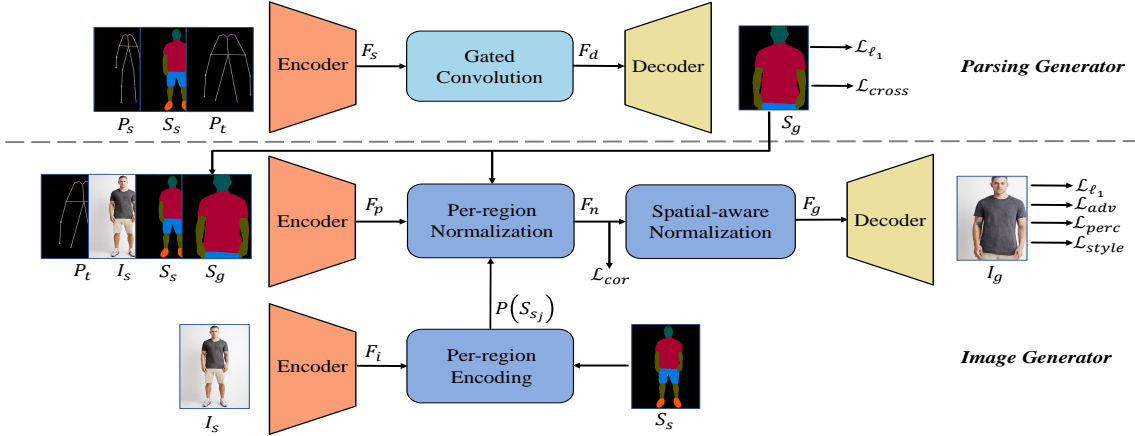


Figure 2. Overview of our model.

and then rendered textures from the original image using a soft-gated warping block. Han *et al.* [26] first synthesized a target human parsing map, then estimated a dense flow to warp source image feature and finally refined the image using a U-Net-based [18] network. These flow-based methods can preserve some details, but cannot cope well with large occlusion, which limits their application on person image synthesis. Men *et al.* [16] used human parsing maps for attribute-controllable person image synthesis. Zhang *et al.* [29] introduced gated convolution to learn a dynamic feature selection mechanism and adaptively deform the image layer by layer. However, these methods fail to disentangle the shape and style information and cannot edit the image flexibly. Our model generates more reasonable and realistic person images with the consistencies of both shape and style. Besides, our model can edit the image more flexibly by disentangling the shape and style information.

### 3. Method

As shown in Figure 2, our approach consists of two generators: a parsing generator and an image generator. The inputs of the parsing generator are the source pose  $P_s$ , the target pose  $P_t$ , and the source parsing map  $S_s$ . The parsing generator estimates a human parsing map  $S_g$  aligned with the target pose  $P_t$ . This allows the image editing of the shape of the final image to be controllable. The pose representation includes 18 human keypoints extracted by Human Pose Estimator (HPE) [1], which has 18 channels and encodes the locations of 18 joints of a human body. The source parsing map  $S_s$  is extracted by Part Grouping Network (PGN) [3] from the source image  $I_s$ . To clean incorrect labels (*e.g.*, left leg and right leg) and reduce the number of categories, we re-organize the map from 21 categories to 8 categories: hair, upper clothes, dress, pants, face,

upper skin, leg, and background. The image generator synthesizes a high-quality image  $I_g$  of the reposed person conditioned on the human parsing map  $S_g$ . The inputs of the image generator are the source image  $I_s$ , the source parsing map  $S_s$ , the generated parsing map  $S_g$  and the target pose  $P_t$ . To provide detailed control of the styles in individual regions, we propose joint global and local per-region encoding and normalization to decouple the style and shape. With the generated parsing map  $S_g$  and per-region style control, we decouple the shape and style of clothing to facilitate image editing tasks. Besides, we propose spatial-aware normalization to retain the spatial context relationship. Note that the model can deal with other tasks, *e.g.*, texture transfer and region editing.

#### 3.1. Parsing Generator

The parsing generator is responsible for generating the human parsing map aligned with the target pose while keeping the clothing style and body shape of the person in the source image. The inputs, the source pose  $P_s$ , the target pose  $P_t$ , and the source parsing  $S_s$ , are first embedded into a latent space by an encoder, which consists of  $M$  down-sampling convolutional layers ( $M = 4$  in our case). Inspired by PINet [29], instead of applying residual blocks [5] like previous methods [2, 26], we use gated convolution to deform the feature  $F_s$  of source human parsing map  $I_s$  from the source pose to the target pose to avoid the drawback of vanilla convolution that treats all the pixels as valid information. The formulation of gated convolution is

$$\begin{aligned}
 O_{x,y} &= \phi \left( \sum_{i=-k_h}^{k_h} \sum_{j=-k_w}^{k_w} u_{k_h+i, k_w+j} \cdot I_{y+i, x+j} \right) \cdot \\
 &\quad \sigma \left( \sum_{i=-k_h}^{k_h} \sum_{j=-k_w}^{k_w} v_{k_h+i, k_w+j} \cdot I_{y+i, x+j} \right),
 \end{aligned} \tag{1}$$

where  $\cdot$  denotes the element-wise multiplication of two feature maps.  $I_{x,y}$  and  $O_{x,y}$  are the input and output at position  $(x, y)$ ,  $k_h = (k_{sh} - 1)/2$  and  $k_w = (k_{sw} - 1)/2$ .  $k_{sh}$  and  $k_{sw}$  are the kernel sizes (e.g.,  $3 \times 3$ ).  $\sigma$  denotes sigmoid function which ensures the output gating values are between 0 and 1.  $\phi$  denotes the activation function.  $u$  and  $v$  are two different convolutional filters.

Gated convolution can learn a dynamic selection mechanism for each spatial location, which is suitable for unaligned generation tasks [28, 27, 29]. Finally, with the deformed parsing feature  $F_d$ , the generated human parsing map  $S_g$  is delivered by a decoder with standard configuration. The detailed network design can be found in the supplementary material.

### 3.2. Image Generator

The image generator aims at transferring the textures of individual regions in the source image to the generated parsing map. From another point of view, this can be seen as a semantic map-to-image translation problem conditioned on the source image. Inspired by SEAN [31], we first extract per-region styles of the source image  $I_s$  with the source human parsing map  $P_s$ , which are then transferred using normalization techniques. Because some visible regions in the target image are invisible in the source image due to large variation in pose, the number of regions in our generated parsing map is different from that of the source image. SEAN sets the styles of invisible regions to zero, which ignores the implied relationship between visible and invisible regions of person images (e.g., a man in a coat is more likely to wear trousers than shorts). Therefore, instead of using only local region-wise average pooling, we propose *joint global and local per-region average pooling* to extract the style of the region in the source image. We then concatenate the source image  $I_s$ , the source parsing map  $S_s$ , the generated parsing map  $S_g$  and the target pose  $P_t$  in depth (channel) dimension and extract its feature  $F_p$ . Finally, similar to existing normalization techniques, we control the per-region style of  $F_p$  by modulating its scale and bias. However, previous normalization techniques lose the spatial information of the source image. To solve this problem, we propose a *spatial-aware normalization* method to preserve the spatial context relationship of the source image. After spatial-aware normalization, the desired person image  $I_g$  is delivered by a decoder.

#### 3.2.1 Per-region encoding

Given the source image  $I_s$ , we first extract its feature map using a bottleneck structure with 4 down-sampling convolution layers and 2 transposed convolution layers. The output of encoder  $F_i$  with 256 channels contains the spatial context relationship and style of the source image. Intuitively, the styles of individual regions are independent of their shapes. With the source parsing map  $S_s$ , we disen-

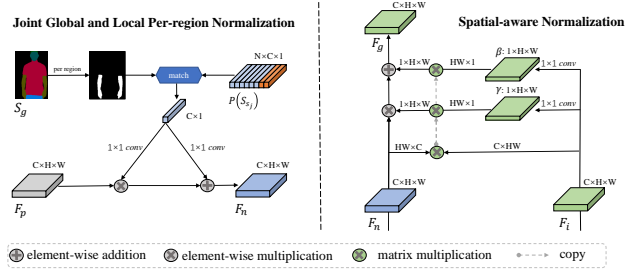


Figure 3. Details of per-region normalization and spatial-aware normalization in our model.

tangle the shape and style by extracting style information using average pooling to remove the shape information. To preserve the style information visible in the source image and predict a reasonable style for invisible regions conditioned on the source image, we propose *joint global and local per-region average pooling* to extract the style of the source image, which is formulated as

$$P(S_{s_j}) = \begin{cases} \text{avg}_{w,h}(F_i \cdot S_{s_j}), & \sum S_{s_j} > 0 \\ \text{avg}_{w,h}(F_i), & \sum S_{s_j} \leq 0 \end{cases}, \quad (2)$$

where  $\text{avg}(\cdot)$  denotes average pooling in the spatial dimension.  $w$  and  $h$  are the width and height of the feature map  $F_i$ .  $j$  indicates the category index, and  $S_{s_j}$  is the source semantic map w.r.t. category  $j$ . The first case considers local average pooling where category  $j$  appears in the source image, and the second case is global average pooling for unseen categories. The size of output  $P(S_{s_j})$  after average pooling is  $256 \times N$ , where  $N$  is the number of label categories (8 in our case).

#### 3.2.2 Per-region normalization

The newly generated feature  $F_p$  from a basic encoder contains the information of the inputs and aligns with the generated parsing map  $S_g$ . With the per-region style code  $P(S_{s_j})$  w.r.t. label category  $j$ , we can transfer the style to the newly generated feature  $F_p$  by modulating its scale and bias. As shown in Figure 3, for each region  $S_{g_j}$  in the generated parsing map  $S_g$ , we find the corresponding region style in  $P(S_{s_j})$ . Thanks to our joint global and local per-region encoding, the style codes of all the regions can be found in  $P(S_{s_j})$ . Then we use two fully connected layers to predict the scale and bias for  $F_p$ , which are then applied for per-region normalization. The generated feature  $F_p$  is updated to  $F_n$  that contains the per-region styles of the source image.

#### 3.2.3 Spatial-aware normalization

As illustrated in Section 3.2.1, the style included in  $P(S_{s_j})$  loses the spatial context relationship of the source image



due to average pooling. It is difficult to make the generated images preserve the details (especially spatial details) in the source image. As shown in Figure 3, to preserve the spatial context relationship in the source image, in addition to extracting styles by average pooling in the spatial dimension, we also extract spatial scale ( $\gamma$ ) and bias ( $\beta$ ) from the source image code  $F_i$  using  $1 \times 1$  convolution layers. We retain spatial context relationships using  $\gamma$  and  $\beta$ . However, due to the misalignment between the source image and the target image, how to transfer  $\gamma$  and  $\beta$  to the generated image is a challenging problem. To tackle this problem, we try to compute the correspondence between the feature after per-region normalization  $F_n$  and the feature of the source image  $I_s$ . We first use a correspondence loss to constrain the similarity between  $F_n$  and the pre-trained VGG-19 [22] feature of the target image. With the correspondence loss,  $F_n$  can be more aligned with the feature of the target image in the latent space, which constrains them to be in the same domain. Then, we compute the correspondence between  $F_n$  and the VGG-19 feature of the source image  $I_s$ . We use the correspondence layer [6] to compute a correlation matrix

$$\mathcal{M}(p_1, p_2) = \frac{F_n(p_1)^T \phi_i(I_s)(p_2)}{\|F_n(p_1)\| \|\phi_i(I_s)(p_2)\|}, \quad (3)$$

where  $\phi_i$  denotes the activation map of the  $i$ -th layer of the VGG-19 network. We use *conv3\_1* in our experiments.  $F_n(p_1)$  and  $F_n(p_2)$  denote the channel-wise centralized features of  $F_n$  at the positions  $p_1$  and  $p_2$ , respectively.

The  $\gamma$  and  $\beta$  that denote spatial context relationships at each position can be transformed from the source image to the target image by multiplying them with the correlation matrix  $\mathcal{M}$ . Then,  $F_n$  is further updated to  $F_g$  by modulating the scale and bias. Therefore, after passing through the spatial-aware normalization module,  $F_g$  retains both the style and spatial relationships of the source image. Finally, The decoder outputs the final image  $I_g$ .

### 3.3. Loss Functions

We first train the parsing generator and the image generator respectively and then end-to-end train our full model. In the following, we will describe the loss function of two generators in details.

**Parsing Generator Loss.** The full loss function of the parsing generator can be formulated as:

$$\mathcal{L}_{parsing} = \mathcal{L}_{cross} + \lambda_{pl} \mathcal{L}_{\ell_1}. \quad (4)$$

where  $\lambda_{pl}$  is the coefficient of the  $\ell_1$  item.

In order to generate reasonable human parsing maps we apply  $\ell_1$  distance loss between the generated parsing map and the target parsing map:

$$\mathcal{L}_{\ell_1} = \|S_g - S_t\|_1. \quad (5)$$

Besides, to generate high-quality human parsing maps, we also use cross-entropy loss  $\mathcal{L}_{cross}$ , which is defined as:

$$\mathcal{L}_{cross} = -\frac{1}{N} \sum_{i=0}^{N-1} S_{t_i} \log(\text{Softmax}(S_{g_i})), \quad (6)$$

where  $N$  is the number of categories of labels ( $N = 8$  in our case).

**Image Generator Loss.** The full loss used for training the image generator consists of correspondence loss, reconstruction  $\ell_1$  loss, perceptual loss, style loss, and adversarial loss, defined as:

$$\mathcal{L}_{image} = \lambda_c \mathcal{L}_{cor} + \lambda_\ell \mathcal{L}_\ell + \lambda_p \mathcal{L}_{per} + \lambda_s \mathcal{L}_{style} + \lambda_a \mathcal{L}_{adv}, \quad (7)$$

where  $\lambda_c$ ,  $\lambda_\ell$ ,  $\lambda_p$ ,  $\lambda_s$  and  $\lambda_a$  are weights that balance contributions of individual loss terms.

The correspondence loss is used to constrain the generated feature  $F_n$  and the pre-trained VGG-19 [22] feature of the target image to align in the same domain, which is defined as:

$$\mathcal{L}_{cor} = \|F_n - \phi_i(I_t)\|_2. \quad (8)$$

In practice, we use the feature from *conv3\_1* of VGG-19 to compute the correspondence loss.

The reconstruction  $\ell_1$  loss is used to encourage the generated image  $I_g$  to be similar with the target image  $I_t$  at the pixel level, which is defined as:

$$\mathcal{L}_{\ell_1} = \|I_g - I_t\|_1. \quad (9)$$

We also adopt perceptual loss and style loss [10] to generate more realistic images. The perceptual loss calculates the  $\ell_1$  distance between activation maps of the pre-trained VGG-19 network, which can be written as:

$$\mathcal{L}_{per} = \sum_i \|\phi_i(I_g) - \phi_i(I_t)\|_1. \quad (10)$$

We adopt the feature of [*relu1\_1*, *relu2\_1*, *relu3\_1*, *relu4\_1*, *relu5\_1*] with the same weight. The style loss measures the statistical difference of the activation maps between the generated image  $I_g$  and the target image  $I_t$ , which is formulated as:

$$\mathcal{L}_{style} = \sum_j \|G_j^\phi(I_t) - G_j^\phi(I_g)\|_1. \quad (11)$$

In practice, we use the feature of [*relu2\_2*, *relu3\_4*, *relu4\_4*, *relu5\_2*] with the same weight. The adversarial loss with discriminator  $D$  is employed to penalize the distribution difference between generated (fake) images  $I_g$  and target (real) images  $I_t$ , which is written as:

$$\mathcal{L}_{adv} = \mathbb{E}[\log(1 - D(I_g))] + \mathbb{E}[\log D(I_t)]. \quad (12)$$

## 4. Experimental Results

**Dataset.** We conduct our experiment on DeepFashion In-shop Clothes Retrieval Benchmark [14], which contains 52712 images with the resolution of  $256 \times 256$ . The images vary in terms of poses and appearances. We split the



Figure 4. Our results of person image synthesis in arbitrary poses.

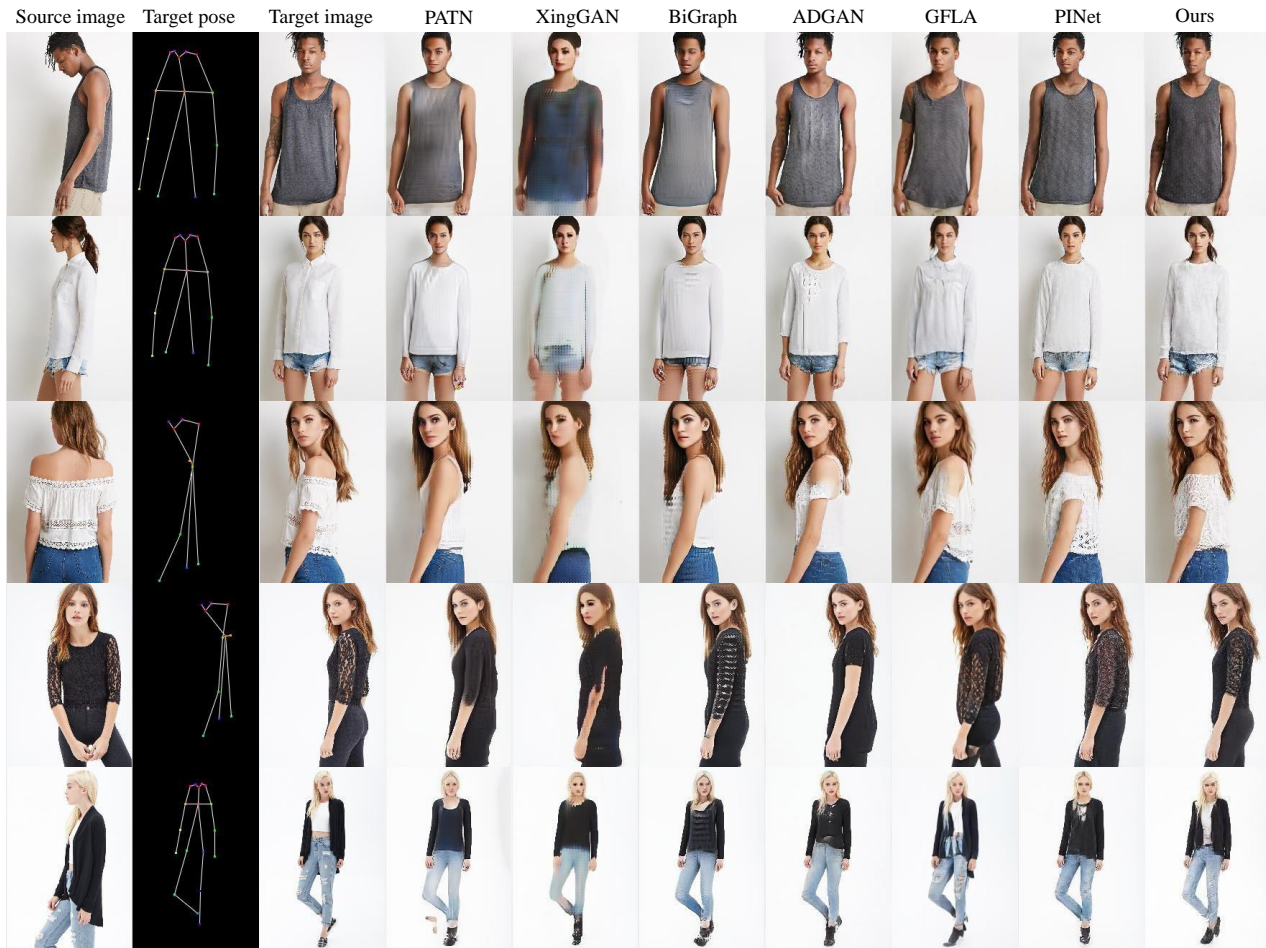


Figure 5. Qualitative comparisons with state-of-the-art methods. From left to right are the results of PATN [32], XingGAN [24], BiGraph [23], ADGAN [16], GFLA [20], PINet [29] and ours, respectively.

data with the same configuration as PATN [32] and collect 101,966 pairs of images for training and 8,750 pairs for testing. Note that the person identities in the test set are different from those in the training set.

**Metrics.** We use Learned Perceptual Image Patch Similarity (LPIPS) [30] to measure the distance between the generated image and the ground-truth image in the perceptual level. Meanwhile, Peak Signal to Noise Ratio (PSNR) is employed to compute the error between the generated image and the ground-truth image in the pixel level. Besides,

we adopt Fréchet Inception Distance (FID) [7] to compute the distance between distributions of the generated images and the ground-truth images, which is used to measure the realism of the generated images.

#### 4.1. Person Image Synthesis in Different Poses

Human pose transfer aims at synthesizing new images for the same source person in different poses. Figure 4 shows some visual results generated by our method. Given a source image and some different poses extracted from im-

Table 1. Quantitative comparison with state-of-the-art methods<sup>1</sup>.

Model	FID ↓	LPIPS ↓	PSNR ↑
PATN [32]	23.70	0.2520	31.16
BiGraph [23]	24.36	0.2428	<b>31.38</b>
XingGAN [24]	41.79	0.2914	31.08
GFLA [20]	<u>14.52</u>	0.2219	31.28
ADGAN [16]	16.00	0.2242	31.30
PINet [29]	15.28	<u>0.2152</u>	<u>31.31</u>
Ours	<b>13.61</b>	<b>0.2059</b>	<b>31.38</b>

Table 2. Quantitative results of ablation study.

Model	FID ↓	LPIPS ↓	PSNR ↑
Global-Enc	15.21	0.2137	31.35
Local-Enc	15.50	0.2138	31.30
w/o SN	<u>14.15</u>	<u>0.2071</u>	<b>31.43</b>
Full	<b>13.61</b>	<b>0.2059</b>	<u>31.38</u>

ages in the test set, our model generates realistic images with fine details.

## 4.2. Comparison with State-of-the-art Methods

We conduct qualitative and quantitative comparisons with several state-of-the-art methods.

**Qualitative Comparison.** We compare the visual results of our method with six state-of-the-art methods: PATN [32], XingGAN [24], BiGraph [23], ADGAN [16], GFLA [20] and PINet [29]. Figure 5 shows some examples. PATN, XingGAN, and BiGraph fail to generate sharp images and cannot keep the consistency of shape and texture due to the use of sparse correspondence extracted from keypoints. The flow-based method, GFLA, preserves the detailed texture in the source image. However, it is difficult to obtain reasonable results for the invisible regions of the source image. ADGAN and PINet succeed in keeping shape consistency and predicting reasonable shapes of clothing, but they lose spatial context relationships. Our model uses spatial-aware normalization to retain spatial context relationships, and hence can generate more natural and sharper images (the second and fourth rows). Besides, our model retains the shape (the first and third rows) and predicts more reasonable results (the second and fifth rows).

**Quantitative Comparison.** Table 1 gives the quantitative results on the 8750 test images compared with six state-of-the-art methods: PATN [32], XingGAN [24], BiGraph [23], ADGAN [16], GFLA [20] and PINet [29]. Our results get the best PSNR score, which measures the error in pixel level. Besides, we introduce LPIPS to compute the similarity in perceptual level and FID to measure the realism of

<sup>1</sup>Note that we take the images resized from  $256 \times 176$  to  $256 \times 256$  as the inputs of GFLA.



Figure 6. Qualitative results of ablation study.

the generated images. Our results achieve the best performance in terms of both FID and LPIPS, which indicates that our model not only generates more realistic images but also keeps the best consistency of shape and texture.

## 4.3. Ablation Study

We train several ablation models to prove our hypotheses and validate the effect of our improvements.

**Global Encoding Model (Global-Enc).** The global encoding model replaces per-region styles with the same global feature extracted from the source image using the global average pooling.

**Local Encoding Model (Local-Enc).** The local encoding model adopts per-region encoding. However, for the style of invisible regions, the style code is set to be zero.

**The Model without Spatial-aware Normalization (w/o SN).** This model is designed to measure the contribution of spatial-aware normalization described in Section 3.2.3. We train this model in the same configuration as our full model.

**Full Model (Full).** We use joint global and local per-region encoding and spatial-aware normalization in this model.

Quantitative results on the 8750 test images are shown in Table 2. As shown in Table 2, compared with the global encoding model and local encoding model, the joint global and local per-region encoding and normalization improve the performance by predicting the style of each region, especially for the authenticity of generated images and the similarity to the ground-truth images. Besides, due to the detailed images generated by our full model, our full model has a slightly lower PSNR than the model without spatial-aware normalization in the pixel level, but it gains the best



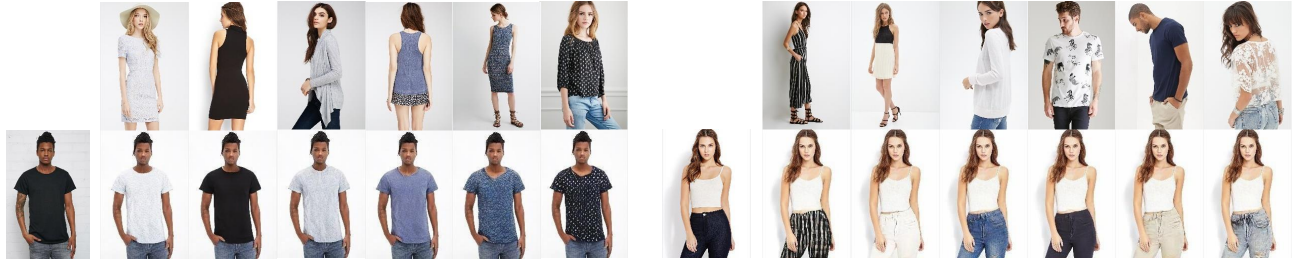


Figure 7. Texture transfer results. The left figure and the right figure show the results of transfer the texture of upper clothes and pants, respectively. In each figure, the first row shows the reference images, and the second row shows the results generated by our method. The first column of each figure shows the source image.



Figure 8. Texture interpolation results.



Figure 9. Region editing results.

performance in similarity in the perceptual level and generates the most realistic images. Figure 6 shows some visual results of ablation study. With joint global and local encoding and normalization, our method can maintain the exact style of visible regions (in the first and second rows) and predict reasonable styles of invisible regions (natural face and shoes in the fifth row). With spatial-aware normalization, our model generates sharper images with the spatial context relationship (in the second, third and fourth rows).

## 5. Applications

Benefiting from our two-stage framework and per-region style encoding, our model decouples the shape and style of clothing and can be applied to various person image editing applications.

**Texture Transfer** The per-region style of our generated images is controllable by the style code  $P(S_{s_j})$ . Therefore, we can change the style of each region by replacing the style code of the corresponding region in  $P(S_{s_j})$ . Specifically, given a reference image, we can extract its style code and transfer the per-region style to our generated image. Figure 7 shows some visual results of texture transfer. We transfer the texture of upper clothes or pants of the reference images to our generated images. Our model generates natural images with detailed texture.

**Texture Interpolation** We control the per-region styles of our generated images by latent code, which can move along the manifold of textures from one style to another style. Taking texture transfer of upper clothes as an example, we interpolate the styles of upper clothes ( $j$ -th in the number of categories) from one image  $I_{s_1}$  to another image  $I_{s_2}$ . The style  $t$  of the upper clothes of the generated image is defined through linear blending as:

$$t = (1 - \alpha)P(S_{s_{1j}}) + \alpha P(S_{s_{2j}}). \quad (13)$$

As shown in Figure 8, the texture of the upper clothes gradually changes from the style of the left reference image to that of the right reference image.

**Region Editing** Since we use a human parsing map as the intermediate result, we can edit the generated images by editing the input parsing of the image generator. Specifically, given a source image and its parsing map, we can flexibly edit the parsing map to automatically synthesize the person image as we desired. As shown in Figure 9, we can change the style of clothing (*e.g.*, T-shirt to waistcoat, pants to dress, and long hair to short hair). Our model generates natural images consistent with the edited parsing map.

## 6. Conclusion

In this paper, we propose PISE, a novel two-stage generative model for person image synthesis and editing. Our method first synthesizes a human parsing map and then generates the target image by joint global and local encoding and normalization and spatial-aware normalization. Experimental results demonstrate that our model achieves promising results with detailed texture and consistent shape of clothing. Besides, the ablation study also verifies the effectiveness of each designed component. Our model can also be applied in various applications such as texture transfer and region editing, and generates natural images. In the future, we will try to generalize our framework to deal with video generation conditioned on a reference image.

**Acknowledgments.** This work was supported in part by the National Natural Science Foundation of China (61771339) and Tianjin Research Program of Application Foundation and Advanced Technology (18JCYBJC19200).



## References

- [1] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Real-time multi-person 2D pose estimation using part affinity fields. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017.
- [2] Haoye Dong, Xiaodan Liang, Ke Gong, Hanjiang Lai, Jia Zhu, and Jian Yin. Soft-gated warping-gan for pose-guided person image synthesis. In *Adv. Neural Inform. Process. Syst.*, 2018.
- [3] Ke Gong, Xiaodan Liang, Yicheng Li, Yimin Chen, Ming Yang, and Liang Lin. Instance-level human parsing via part grouping network. In *Eur. Conf. Comput. Vis.*, 2018.
- [4] Ian J. Goodfellow, Jean Pouget-Abadie, M. Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *Adv. Neural Inform. Process. Syst.*, 2014.
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2016.
- [6] Mingming He, Dongdong Chen, Jing Liao, Pedro V Sander, and Lu Yuan. Deep exemplar-based colorization. *ACM Trans. Graph.*, 37(4):47, 2018.
- [7] 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 *Adv. Neural Inform. Process. Syst.*, 2017.
- [8] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *IEEE Int. Conf. Comput. Vis.*, 2017.
- [9] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017.
- [10] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Eur. Conf. Comput. Vis.*, 2016.
- [11] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.
- [12] Kun Li, Jinsong Zhang, Yebin Liu, Yu-kun Lai, and Qionghai Dai. PoNA: Pose-guided non-local attention for human pose transfer. *IEEE Trans. Image Process.*, 29:9584–9599, 2020.
- [13] Yining Li, Chen Huang, and Chen Change Loy. Dense intrinsic appearance flow for human pose transfer. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.
- [14] Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, and Xiaoou Tang. DeepFashion: Powering robust clothes recognition and retrieval with rich annotations. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2016.
- [15] Liqian Ma, Xu Jia, Qianru Sun, Bernt Schiele, Tinne Tuytelaars, and Luc Van Gool. Pose guided person image generation. In *Adv. Neural Inform. Process. Syst.*, 2017.
- [16] Yifang Men, Yiming Mao, Yuning Jiang, Wei-Ying Ma, and Zhouhui Lian. Controllable person image synthesis with attribute-decomposed gan. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.
- [17] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [18] Ronneberger Olaf, Fischer Philipp, and Brox Thomas. U-net: Convolutional networks for biomedical image segmentation. *arXiv preprint arXiv:1505.04597*, 2015.
- [19] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.
- [20] Yurui Ren, Ge Li, Shan Liu, and Thomas H Li. Deep spatial transformation for pose-guided person image generation and animation. *IEEE Trans. Image Process.*, 29:8622–8635, 2020.
- [21] Yurui Ren, Xiaoming Yu, Ruonan Zhang, Thomas H. Li, Shan Liu, and Ge Li. Structureflow: Image inpainting via structure-aware appearance flow. In *ICCV*, 2019.
- [22] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *IEEE Int. Conf. Comput. Vis.*, 2015.
- [23] Hao Tang, Song Bai, Philip HS Torr, and Nicu Sebe. Bipartite graph reasoning gans for person image generation. In *Brit. Mach. Vis. Conf.*, 2020.
- [24] Hao Tang, Song Bai, Li Zhang, Philip HS Torr, and Nicu Sebe. XingGAN for person image generation. In *Eur. Conf. Comput. Vis.*, 2020.
- [25] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018.
- [26] Han Xintong, Weulin Huang, Xiaojun Hu, and Scott. Matthew R. ClothFlow: A flow-based model for clothed person generation. In *IEEE Int. Conf. Comput. Vis.*, 2019.
- [27] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Free-form image inpainting with gated convolution. *arXiv preprint arXiv:1806.03589*, 2018.
- [28] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Generative image inpainting with contextual attention. *arXiv preprint arXiv:1801.07892*, 2018.
- [29] Jinsong Zhang, Xingzi Liu, and Kun Li. Human pose transfer by adaptive hierarchical deformation. *Computer Graphics Forum*, 39(7):325–337, 2020.
- [30] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018.
- [31] Peihao Zhu, Rameen Abdal, Yipeng Qin, and Peter Wonka. SEAN: Image synthesis with semantic region-adaptive normalization. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020.
- [32] Zhen Zhu, Tengting Huang, Baoguang Shi, Miao Yu, Bofei Wang, and Xiang Bai. Progressive pose attention transfer for person image generation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2342–2351, 2019.
- [33] Zhen Zhu, Zhiliang Xu, Ansheng You, and Xiang Bai. Semantically multi-modal image synthesis. In *IEEE Conf. Comput. Vis. Pattern Recog.*