# Disentangled Person Image Generation

Liqian Ma[1*]   Qianru Sun[2*]   Stamatios Georgoulis[1]
Luc Van Gool[1,3]   Bernt Schiele[2]   Mario Fritz[2]

[1]KU-Leuven/PSI, Toyota Motor Europe (TRACE)   [3]ETH Zurich
[2]Max Planck Institute for Informatics, Saarland Informatics Campus
{liqian.ma, sgeorgou, luc.vangool}@esat.kuleuven.be
vangool@vision.ee.ethz.ch
{qsun, schiele, mfritz}@mpi-inf.mpg.de

## Abstract

*Generating novel, yet realistic, images of persons is a challenging task due to the complex interplay between the different image factors, such as the foreground, background and pose information. In this work, we aim at generating such images based on a novel, two-stage reconstruction pipeline that learns a disentangled representation of the aforementioned image factors and generates novel person images at the same time. First, a multi-branched reconstruction network is proposed to disentangle and encode the three factors into embedding features, which are then combined to re-compose the input image itself. Second, three corresponding mapping functions are learned in an adversarial manner in order to map Gaussian noise to the learned embedding feature space, for each factor respectively. Using the proposed framework, we can manipulate the foreground, background and pose of the input image, and also sample new embedding features to generate such targeted manipulations, that provide more control over the generation process. Experiments on Market-1501 and Deepfashion datasets show that our model does not only generate realistic person images with new foregrounds, backgrounds and poses, but also manipulates the generated factors and interpolates the in-between states. Another set of experiments on Market-1501 shows that our model can also be beneficial for the person re-identification task.*

## 1. Introduction

The process of generating realistic-looking images of persons is of great importance for the computer vision community and finds application in different tasks, like image editing, person re-identification (re-ID), in-painting or



Figure 1: Left: image sampling results on Market-1501. Three factors, *i.e.* foreground, background and pose, can be sampled independently (1st-3rd rows) and jointly (4th row). Right: similar joint sampling results on DeepFashion (images are cut for the sake of display). This dataset contains almost no background, so we only disentangle the image into appearance and pose factors. Zoom in for details.

on-demand generated art for movie production. The recent advent of image generation models, such as variational autoencoders (VAE) [13], generative adversarial networks (GANs) [7] and autoregressive models (ARMs) (*e.g.* Pixel-RNN [33]), has provided powerful tools towards this goal. Several papers [25, 2, 1] have then exploited the ability of these networks to generate sharp images in order to synthesize realistic photos of faces and natural scenes. Most recently, Ma *et al*. [21] proposed an architecture to synthesize novel person images in arbitrary poses given as input

an image of that person and a new pose.

From an application perspective, however, the user usually wants to have more control over the generated images (*e.g.* change the background, a person's appearance and clothing, or its viewpoint), which is something that existing methods are essentially lacking still. In this work, we go beyond these constraints and investigate how to generate novel person images with a specific user intention in mind (*i.e.* foreground (FG), background (BG), pose manipulation). The key idea is to explicitly guide the generation process by an appropriate representation of that intention. Fig. 1 gives examples of the intended generated images.

To this end, we disentangle the input image into intermediate embedding features, *i.e.* person images can be reduced to a composition of foreground features, background features and pose features. Unlike existing approaches, however, we rely on a different technique to generate new samples. In particular, we aim at sampling from a standard distribution, *e.g.* Gaussian distribution, to first generate new embedding features and from them generate new images. To achieve this, *fake* embedding features $\tilde{e}$ are learned in an adversarial manner to match the distribution of the *real* embedding features $e$, where the encoded features from the input image are treated as *real* whilst the ones generated from the Gaussian noise as *fake* (see Fig. 2). Consequently, the newly sampled images come from learned *fake* embedding features $\tilde{e}$ rather than the original Gaussian noise as in the traditional GAN models. By doing so, the proposed technique enables us not only to sample a controllable input for the generator, but also to preserve the complexity of the composed images (*i.e.* realistic person images).

To sum up, our full pipeline proceeds in two stages as shown in Fig. 2. At stage-I, we use a person's image as input and disentangle the information into three main factors, namely foreground, background and pose. Each disentangled factor is modeled by embedding features through a reconstruction network. At stage-II, a mapping function is learned to map Gaussian distribution to feature embedding distribution.

Our contributions can be summarized as follows: (1) A new task of generating natural person images by disentangling the input into weakly correlated factors, namely foreground, background and pose. (2) A two-stage framework to learn manipulatable embedding features for all three factors. In stage-I, the encoder of the multi-branched reconstruction network serves conditional image generation tasks, whereas in stage-II the mapping functions learned through adversarial training (*i.e.* mapping noise $z$ to *fake* embedding features $emb$) serve sampling tasks (*i.e.* the input is sampled from a standard Gaussian distribution). (3) A technique to match the distribution of *real* and *fake* embedding features through adversarial training, which is not bound to the image generation task. (4) An approach to gen-

erate new sampled image pairs for the person re-ID task. As shown in Sec. 4, we construct a Virtual Market re-ID dataset by fixing the foreground features and changing the background features and pose keypoints to generate samples of one identity.
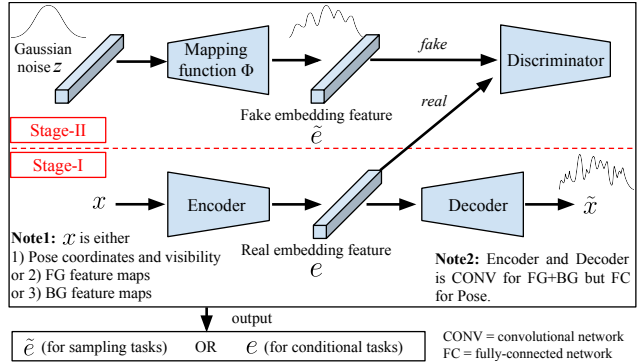
## 2. Related work



Figure 2: Our two-stage framework. In stage-I, we use a reconstruction network to obtain the *real* embedding features $e$ for each factor, *i.e.* foreground, background and pose. The architectural details of stage-I are shown in Figure 3. In stage-II, we propose a novel, two-step mapping technique for adversarial embedding feature learning that first map Gaussian noise $z$ to intermediate embedding features $\tilde{e}$ then to the data $\tilde{x}$. We use the pre-trained encoder and decoder of stage-I to guide the learning of mapping functions $\Phi$.

**Image generation from noise.** The ability of generative models, such as GANs [7], adversarial autoencoders (AAE) [22], VAEs [13] and ARMs (*e.g.* PixelRNN [33]), to synthesize realistic-looking, sharp images has led image generation research the last years. Traditional image generation works use GANs [7] or VAEs [13] to map a distribution generated by noise $z$ to the distribution of real data. Convolutional VAEs and AAEs [22] have shown how to transform an auto-encoder into a generator, but in this case, it is rather difficult to train the mapping function for complex data distributions, such as person images (as also mentioned in ARAE-GAN [11]). As such, traditional image generation methods are not optimal when it comes to the human body. For example, Zheng *et al.* [40] directly adopted the DCGAN architecture [25] to generate person images from noise, but as can be seen in Fig. 7(b) vanilla DCGAN leads to unrealistic results. Instead, we propose to use a two-step mapping technique in stage-II to guide the learning, *i.e.* $z \rightarrow e \rightarrow x$ (see Fig. 2). Similar to [11], we use a decoder to adversarially map the noise distribution to the feature embedding distribution learned by the reconstruction network.

**Conditional image generation.** Since the human body has a complex non-rigid structure with a lot of degrees of freedom [23], several works have used structure conditions
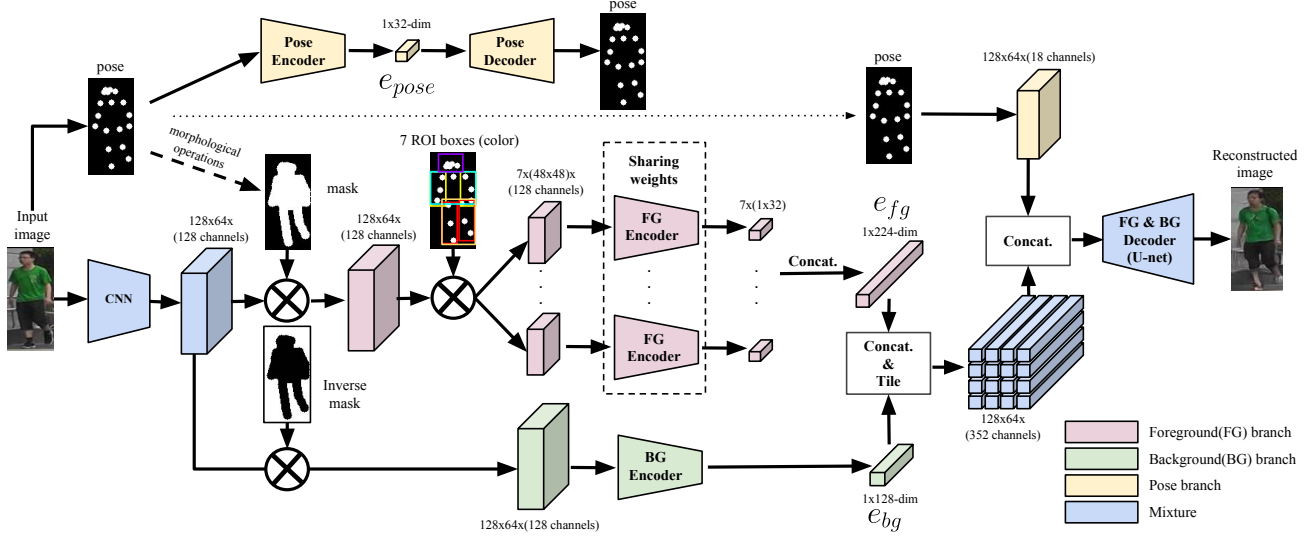
Figure 3: Stage-I: disentangled image reconstruction. This reconstruction framework is composed by three branches, namely foreground, background and pose. Note that we use a fully-connected auto-encoder network to reconstruct the pose (incl. keypoint coordinates and visibility), so that we can decode the embedded pose features to obtain the heatmaps at the sampling phase.

to generate person images. Reed *et al.* in [26] proposed the Generative Adversarial What-Where Network that uses pose keypoints and text descriptions as condition, whereas in [27] they used an extension of PixelCNN in addition to conditioning on part keypoints, segmentation masks and text to generate images on the MPI Human Pose dataset, among others. Lassner *et al.* [15] generated full-body images of persons in clothing by conditioning on fine-grained body and clothing segments, *e.g.* pose, shape or color. Zhao *et al.* [37] combined the strengths of GANs with variational inference to generate multi-view images of persons in clothing in a coarse-to-fine manner. Closer to our work, Ma *et al.* [21] proposed to condition on image and pose keypoints to transfer the human pose in a flexible way, but their method needs a training set of aligned person image pairs which costs expensive human annotations. Most recently, Zhu *et al.* [41] proposed the CycleGAN that uses cycle consistency to achieve unpaired image-to-image translation between domains. They achieve compelling results in appearance changes but show little success in geometric changes.

Since images themselves contain abundant context supervision information, some works have tried to tackle the same problem in an unsupervised way. Doersch *et al.* [4] explored the use of spatial context, *i.e.* relative position between two neighboring patches in an image, as a supervisory signal for unsupervised visual representation learning. Noroozi *et al.* [24] extended the task to a jigsaw puzzle problem solved by observing all the tiles at the same time, which can reduce the ambiguity among these local patch pairs. Lee *et al.* [16] utilized the context information in a

image generation task though inferring the spatial arrangement and generating the image at the same time. In contrast, we use the supervision in a different way. That is, to extract pose-invariant appearance features, we arrange the body part feature embeddings according to the region-of-interest (ROI) bounding boxes obtained with pose keypoints. Then, we explicitly utilize these pose keypoints as structure information to select the necessary appearance features for each body part and generate the entire person image.

In general, this paper studies a different problem compared to these supervised or unsupervised approaches and tries to solve the disentangled person image generation task in an unpaired, self-supervised manner, by leveraging foreground, background and pose sampling at the same time, in order to gain more control over the generation process.

**Disentangled image generation.** Few papers have already tried to work towards this direction by learning a disentangled representation of the input image. Chen *et al.* [2] proposed InfoGAN, an extension to GANs, to learn disentangled representations using mutual information in an unsupervised manner, like writing styles from digit shapes on the MNIST dataset, pose from lighting of 3D rendered images, and background digits from the central digit on the SVHN dataset. Cheung *et al.* [3] added a cross-covariance penalty in a semi-supervised autoencoder architecture in order to disentangle factors, like hand-writing style for digits and subject identity in faces. Tran *et al.* [32] proposed DR-GAN to learn both a generative and a discriminative representation from one or multiple face images to synthesize identity-preserving faces at target poses. In contrast, our

method gives an explicit representation of the main 3 axis of variation (foreground, background, pose). Moreover, training is facilitated without the need of expensive identity annotations - which is not readily available at scale.

# 3. Method

Our goal is to disentangle the appearance and structure factors contained in person images, so that we can manipulate the foreground, background and pose separately. To achieve this, we propose a two-stage pipeline shown in Fig. 2. In stage-I, we disentangle the foreground, background and pose factors using a reconstruction network in a divide-and-conquer manner. In particular, we reconstruct person images by first disentangling into intermediate embedding features of the three factors, then recover the input image by decoding these features. In stage-II, we treat these features as *real* to learn mapping functions $\Phi$ for mapping a Gaussian distribution to the embedding feature distribution adversarially.

## 3.1. Stage-I: Disentangled image reconstruction

At stage-I, we propose a multi-branched reconstruction architecture to disentangle the foreground, background and pose factors as shown in Fig. 3. Note that, to obtain the pose heatmaps and the coarse pose mask we adopt the same procedure as in [21], but we instead use them to guide the information flow in our multi-branched network.

**Foreground branch.** To separate the foreground and background information, we apply the coarse pose mask to the feature maps instead of the input image directly. By doing so, we can alleviate the inaccuracies of the coarse pose mask. Then, in order to further disentangle the foreground from the pose information, we encode pose invariant features with 7 Body Regions-Of-Interest instead of the whole image similar to [38]. Specifically, for each ROI we extract the feature maps resized to $48 \times 48$ and pass them into the weight sharing foreground encoder to increase the learning efficiency. Finally, the encoded 7 body ROI embedding features are concatenated into a 224-dim feature vector. In the following sections, we use BodyROI7 to denote our model which uses 7 body ROIs to extract foreground embedding features, and use WholeBody to denote our model that extracts foreground embedding features from the whole feature maps directly instead of extracting and resizing the ROI feature map.

**Background branch.** For the background branch, we apply the inverse pose mask to get the background feature maps and pass them into the background encoder to obtain a 128-dim embedding feature. Then, the foreground and background features are concatenated and tiled into $128 \times 64 \times 352$ appearance feature maps.

**Pose branch.** For the pose branch, we concatenate the 18-channel heatmaps with the appearance feature maps and

pass them into the a "U-Net"-based architecture [29], *i.e.*, convolutional autoencoder with skip connections, to generate the final person image following PG$^2$ (G1+D) [21]. Here, the combination of appearance and pose imposes a strong explicit disentangling constraint that forces the network to learn how to use pose structure information to select the useful appearance information for each pixel. For pose sampling, we use an extra fully-connected network to reconstruct the pose information, so that we can decode the embedded pose features to obtain the heatmaps. Since some body regions may be unseen due to occlusions, we introduce a visibility variable $\alpha_i \in \{0, 1\}, i = 1, ..., 18$ to represent the visibility state of each pose keypoint. Now, the pose information can be represented by a 54-dim vector (36-dim keypoint coordinates $\gamma$ and 18-dim keypoint visibility $\alpha$).

## 3.2. Stage-II: Embedding feature mapping

Natural signals, such as images, can be represented by a low-dimensional, continuous feature embedding space. In particular, in [34, 30, 35, 5] it has been shown that they lie on or near a low-dimensional manifold of the original high-dimensional space. Therefore, the distribution of this feature embedding space should be more continuous and easier to learn compared to the real data distribution. Some works [36, 8, 28] have then attempted to use the intermediate feature representations of a pre-trained DNN to guide another DNN. Inspired by these ideas, we propose a two-step mapping technique as illustrated in Fig. 2. Instead of directly learning to decode Gaussian noise to the image space, we first learn a mapping function $\Phi$ that maps a Gaussian space **Z** into a continuous feature embedding space **E**, and then use the pre-trained decoder to map the feature embedding space **E** into the real image space **X**. The encoder learned in stage-I encodes the FG, BG and Pose factors $x$ into low-dimensional *real* embedding features $e$. Then, we treat the features mapped from Gaussian noise $z$ as *fake* embedding features $\tilde{e}$ and learn the mapping function $\Phi$ adversarially. In this way, we can sample *fake* embedding features from noise and then map them back to images using the decoder learned in stage-I. The proposed two-step mapping technique is easy to train in a piecewise style and most importantly can be useful for other image generation applications.

## 3.3. Person image sampling

As explained above, each image factor can not only be encoded from the input information, but also be sampled from Gaussian noise. Regarding the latter, to sample a new foreground, background or pose, we combine the decoders learned in stage-I and mapping functions learned in stage-II to construct a $z \rightarrow \tilde{e} \rightarrow \tilde{x}$ sampling pipeline as shown in Fig. 4. Note that, for foreground and background sampling the decoder is a convolutional "U-net"-based architecture, while for pose sampling the decoder is a fully-connected
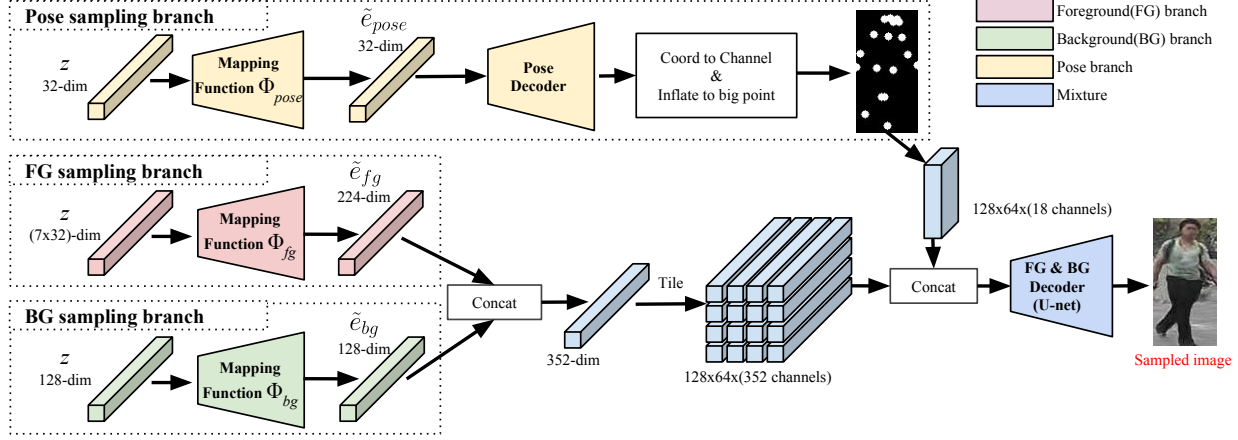
Figure 4: Sampling phase: Sample foreground, background and pose from Gaussian noise to compose new person images.

architecture. In our experiments, we show that our framework performs well when used in both a conditional and an unconditional way.

### 3.4. Network architecture

Here, we describe the proposed architecture. For both stages, we use residual blocks to make the training easier. All convolution layers consist of $3 \times 3$ filters and the number of filters increases linearly with each block. All fully-connected layers consist of 512-dim, except for the bottleneck layers. We apply rectified linear units (ReLU) to each layer, except for the bottleneck and the output layers.

For the foreground and background branches in stage-I, the input image is fed into a convolutional residual block and the pose mask is used to extract the foreground and background feature maps. Then, the masked foreground and background feature maps are passed into an encoder consisting of $N$ convolutional residual blocks, respectively, where $N$ depends on the size of the input. Similar to [21], each residual block consists of two convolution layers with stride=1, followed by one sub-sampling convolution layer with stride=2, except for the last block. For the decoder, an "U-Net"-based architecture [29] is used with $N$ convolutional residual blocks before and after the bottlenecks, respectively, following PG$^2$ (G1+D) [21].

For pose reconstruction, we use an auto-encoder architecture where both encoder and decoder consist of 4 fully-connected residual blocks with 32-dim bottleneck layers. As in [9], we use a densely-connected-like architecture, *i.e.* each residual block consists of two fully-connected layers.

For each mapping function in stage-II, we use a fully-connected network consisting of 4 fully-connected residual blocks to map $K$-dim Gaussian noise $z$ to $K$-dim embedding features $e$. For the discriminator, we adopt a fully-connected network with 4 fully-connected layers.

### 3.5. Optimization strategy

The training procedures of stage-I and stage-II are separated, since the mapping functions $\Phi_{fg}$, $\Phi_{bg}$ and $\Phi_{pose}$ in stage-II can be trained in a piecewise style. In stage-I, we use both L1 and adversarial loss to optimize the image (*i.e.* foreground and background) reconstruction network. This choice is known to result in sharper and more realistic images. In particular, we use $G_1$ and $D_1$ to denote the image reconstruction network and the corresponding discriminator in stage-I. The overall losses for $G_1$ and $D_1$ are as follows,

$$
\begin{aligned}
\mathcal{L}_R^{D_1} =\ & \mathbb{E}_{x \sim p_{data}(x)}\big[\log D_1(x)\big] + \\
& \mathbb{E}_{x \sim p_{data}(x)}\big[\log\big(1 - D_1(G_1(x,h))\big)\big], \quad (1)
\end{aligned}
$$

$$
\begin{aligned}
\mathcal{L}_R^{G_1} =\ & \mathbb{E}_{x \sim p_{data}(x)}\big[\log\big(D_1(G_1(x,h))\big)\big] + \\
& \lambda \|(G_1(x,h) - x)\|_1, \quad (2)
\end{aligned}
$$

where $x$ denotes the person image, $h$ denotes the pose heatmaps, and $\lambda$ is the weight of L1 loss controlling how close the reconstruction looks like to the input image at low frequencies. For pose reconstruction, we use the L2 loss to reconstruct the input pose information including keypoint coordinates $\gamma$ and visibility $\alpha$ mentioned in Sec. 3.1,

$$
\mathcal{L}_R^{Pose} = \mathbb{E}_{(\gamma,\alpha) \sim p_{data}(\gamma,\alpha)} \|(G_1(\gamma,\alpha) - (\gamma,\alpha)\|_2^2, \quad (3)
$$

After training the reconstruction network in stage-I, we fix it and use the Wasserstein GAN [1] loss to optimize the fully-connected network of mapping functions in stage-II. We use $\Phi$ and $D_2$ to denote the mapping functions (incl. $\Phi_{fg}$, $\Phi_{bg}$ and $\Phi_{pose}$) and the corresponding discriminators in stage-II. The overall losses for $\Phi$ and $D_2$ are as follows,

$$
\mathcal{L}_M^{D_2} = \mathbb{E}_{e \sim p_{emb}(e)}\big[D_2(e)\big] - \mathbb{E}_{z \sim p_z(z)}\big[D_2(\Phi(z))\big], \quad (4)
$$

$$
\mathcal{L}_M^{\Phi} = \mathbb{E}_{z \sim p_z(z)}\big[D_2(\Phi(z))\big], \quad (5)
$$

where $e$ denotes the embedding features extracted from the reconstruction network in stage-I, $z$ denotes the Gaussian

noise. Note that, we also tried the vanilla GAN loss but suffered a model collapse. For adversarial training, we optimize the discriminator and generator alternatively.

# 4. Experiments

The proposed pipeline enables many interesting applications, including image manipulation, pose-guided person image generation, image interpolation, image sampling and person re-identification.[1]

## 4.1. Datasets and metrics

Our main experiments use the challenging re-ID dataset Market-1501 [39], containing 32,668 images of 1,501 persons captured from six disjoint surveillance cameras. All images are resized to $128 \times 64$ pixels. We use the same train/test split (12,936/19,732) as in [39], but use all the images in the train set for training without any identity label. For pose-guided person image generation, we randomly select 12,800 pairs in the test set for testing, following [21]. For re-ID, we follow the same testing protocol as in [39].

We also experiment with a high-resolution dataset, namely DeepFashion (In-shop Clothes Retrieval Benchmark) [20], that consists of 52,712 in-shop clothes images and 200,000 cross-pose/scale pairs. Following [21], we use the up-body person images and filter out failure cases in pose estimation for both training and testing. Thus, we have 15,079 training images and 7,996 testing images. We also randomly select 12,800 pairs from the test set for pose-guided person image generation testing.

**Implementation details.** For Market-1501, our method is applied to disentangle the image into three factors: foreground, background and pose. We set the number of convolutional residual blocks $N = 5$ for foreground and background encoders and decoders. For DeepFashion, since the data contains almost no background, our method is applied to disentangle the image into only two factors: appearance and pose. We set the number of convolution blocks $N = 7$ for the foreground encoder and decoder. On both datasets, we do a left-right flip data augmentation.

## 4.2. Image manipulation

As explained, a person's image can be disentangled into three factors: FG, BG and Pose. Each factor can then be generated either from a Gaussian signal (sampling) or conditioned on input data, namely image and pose (conditioning). The conditional case contains at least one other factor sampled from Gaussian signals. In Fig. 1, the left-top3 rows show examples with one-factor sampling and two-factor conditioning for FG, BG and Pose on Market-1501, respectively. Our framework successfully manipulates each

---

[1]More generated results, parameters of our network architecture and training details are given in the supplementary material.

| Model | DeepFashion | | Market-1501 | | | |
|---|---|---|---|---|---|---|
| | SSIM | IS | SSIM | IS | mask-SSIM | mask-IS |
| PG$^2$[21] | 0.762 | 3.090 | 0.253 | 3.460 | 0.792 | 3.435 |
| Ours | 0.614 | 3.228 | 0.099 | 3.483 | 0.614 | 3.491 |

Table 1: Quantitative evaluation. Higher scores are better.

intended factor while keeping the others unchanged. In the first row, we sample foreground with $z_{fg} \to \tilde{e}_{fg}$ and condition background and pose with $x \to e$, so that different cloth colors, styles and hair styles can be generated while the pose and background stay mostly the same. Similarly, we can manipulate the background and pose independently as shown in the left-second/third row. The left-last row shows a sampling example without any conditioning. In this way, we can sample novel person images from noise and still generate realistic images compared to vanilla VAE and DCGAN as shown in Sec. 4.5. Finally, on the right rows we show that our method can also sample $256 \times 256$ images with realistic cloth and hair details on DeepFashion.

## 4.3. Pose-guided person image generation

We compare our method with PG$^2$ on pose-conditional person image generation. Unlike PG$^2$, our method does not need paired training images. As shown in Fig. 5, our method can generate more realistic details and less artifacts. Especially, the arms and legs are better shaped on both datasets, and the hair details are more clear on Deep-Fashion. This is in agreement with the Inception Score (IS) and mask Inception Score (mask-IS) in Table 1. The SSIM score of our method is lower than PG$^2$ mainly for two reasons. 1) In stage-I, there are no skip-connections between encoder and decoder, and as such our method has to generate images from compressed embedding features instead of pixel level transforms like in PG$^2$, which is a harder task. 2) Our method generates sharper images which might decrease the SSIM score, as also observed in [21, 10, 31].

## 4.4. Image interpolation

Interpolation is possible for sampled and real images.
**Sampling interpolation.** For sampling interpolation, we directly interpolate in Gaussian space and generate images in a $z \to \tilde{e} \to \tilde{x}$ manner. In particular, we first interpolate linearly between two Gaussian codes $z_1$ and $z_2$ to obtain intermediate codes $z_i$, which in turn are mapped into embedding features $\tilde{e}_i$ using the learned mapping functions. The person's image is then generated from the embedding features $\tilde{e}_i$. As illustrated in Fig. 14(a)(b)(c), our method can smoothly interpolate each factor in Gaussian space separately, which validates that: 1) our method can learn the encoders for foreground, background and pose in a disentangled way; 2) the encoders can map real high-dimensional data distributions into continuous low-dimensional feature
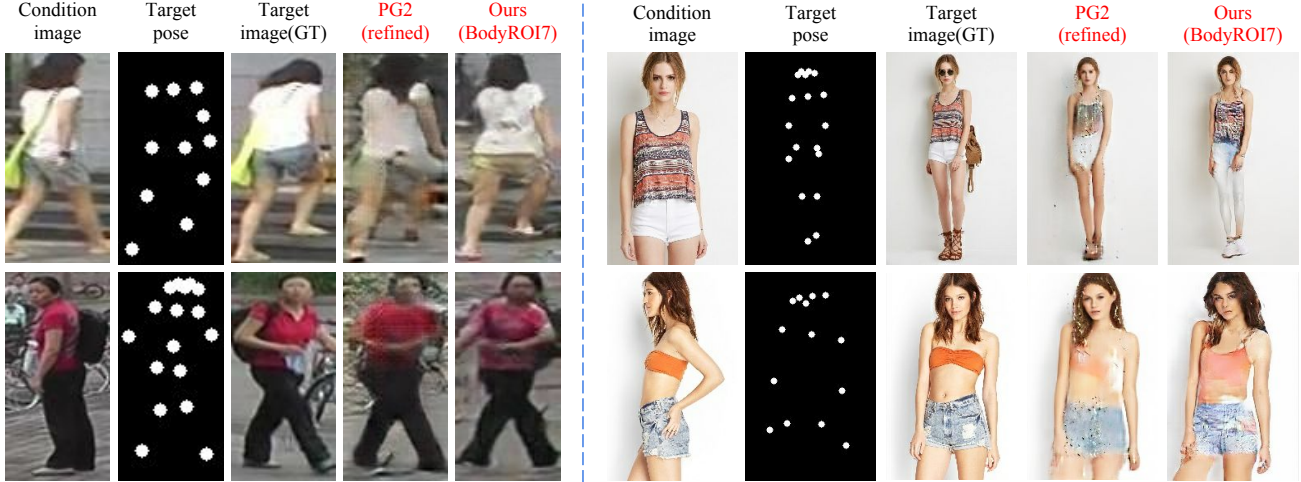
| Condition image | Target pose | Target image(GT) | PG2 (refined) | Ours (BodyROI7) | | Condition image | Target pose | Target image(GT) | PG2 (refined) | Ours (BodyROI7) |

Figure 5: Comparison to PG$^2$. Left: results on Market-1501. Right: results on DeepFashion. Zoom in for details.



(a) Foreground interpolation

(b) Background interpolation

(c) Pose interpolation

x1    z1    z1+d    z1+2d    ...    ⟶    ...    z1+(n-2)d    z1+(n-1)d    z2    x2

(d) Inverse interpolation between two images.

Figure 6: Factor interpolation. (a)(b)(c) We randomly select two Gaussian codes $z_1$ and $z_2$ and interpolate codes between $z_1$ and $z_2$ linearly; we then generate the interpolated images accordingly. (d) We invert an image pair first to embedding features $e_1$ and $e_2$, then to Gaussian codes $z_1$ and $z_2$. We then follow the same procedure as in (a)(b)(c).

embedding distributions; 3) the mapping functions trained adversarially can map Gaussian distributions to feature embedding distributions; 4) the decoder can map feature embedding distributions back to real data distributions.

**Inverse interpolation** To interpolate between real data (incl. image and pose keypoints), we proceed in 3 steps. 1) $x \rightarrow e$: Use the learned encoders to encode real data $x$ into embedding features $e$. 2) $e \rightarrow z$: Use gradient-based minimization [19] to find the corresponding Gaussian codes $z$. 3) $z \rightarrow \tilde{e} \rightarrow \tilde{x}$: Interpolate linearly between two Gaussian codes, then map intermediate codes into embedding features - using the learned mapping functions - to generate the person image. As shown in Fig. 14, our method interpolates reasonable frames between the input pair showing a person with different poses. The interpolated sequence shows realistic intermediate states and can be used to predict potential behaviors.

## 4.5. Sampling results comparison

In this experiment, we compare sampling results from our method and baseline models, *i.e.* VAE [13] and DC-GAN [25]. As illustrated in Fig. 7, VAE generates blurry images and DCGAN sharp but unrealistic person images. In contrast, our model generates more realistic images (see Fig. 7(c)(d)(e)). By comparing (d) and (c), we observe that our model using body ROI generates more sharp and realistic images whose colors on each body part are more natural. A similar tendency can be observed for re-ID. By comparing (e) and (d), we see that when sampling foreground and background but using the real pose keypoints randomly selected from the training data, we generate better results. Therefore, we use this setting in (e) to sample virtual data for the following re-ID experiment.

7

(a) VAE [13]       (b) DCGAN [25]       (c) Ours - Whole Body







(d) Ours - BodyROI7   (e) Ours - BodyROI7 with real pose from training set   (f) Real data

Figure 7: Sampling results comparison. From left to right and from top to bottom: (a) VAE [13] (b) DCGAN [25] (c) Ours - Whole Body (d) Ours - BodyROI7 (e) Ours - BodyROI7 with real pose from training set (f) Real data.



Figure 8: Virtual identities for re-ID model training. Each column contains a pair of images of one identity (one FG). BG and Pose are randomly selected from training data.

| Model | Training data | Rank-1 | mAP |
|---|---|---|---|
| Bow [39] | Market | 0.344 | 0.141 |
| Bow* [39] | Market | 0.358 | 0.148 |
| LOMO* [18] | / | 0.272 | 0.08 |
| WholeBody feature | Market | 0.307 | 0.100 |
| BodyROI7 feature | Market | 0.338 | 0.107 |
| BodyROI7 feature PCA | Market | 0.355 | 0.114 |
| Res50* [6] | *CUHK03* (labeled) | 0.300 | 0.115 |
| Res50* [6] | *Duke* (labeled) | 0.361 | 0.142 |
| Res50 | VM | 0.338 | 0.134 |
| Res50+PUL | VM+Market | 0.369 | 0.156 |
| Res50+PUL+KISSME | VM+Market | 0.375 | 0.154 |

Table 2: Re-ID results on Market-1501. Top: using embedding features. Bottom: using VM and Market-1501 dataset without labels. Higher scores are better. *Results are reported in [6].

## 4.6. Person re-identification

Person re-ID refers to associating the person who appears under different cameras or at different time. Given the query person image, re-ID is expected to provide matching images of the same identity. We propose to use the re-ID performance as a quantitative metric for our generation approach. We adopt the re-ID model in [6] and use rank-1 matching rate and mean Average Precision (mAP) following [39]. We show that our approach can be evaluated in two ways: (1) use FG features extracted in stage-I for re-ID;

(2) generate virtual image pairs to train re-ID model. The virtual market data is denoted as "VM" generated with our BodyROI7 model. Note that, *CUHK03* [17] and *Duke* [40] datasets are used with identity labels, while *Market-1501* and *VM* datasets are used with no labels.

**Using embedding features.** We use the FG encoder to extract the features for re-ID and use the re-ID performance to evaluate the reconstruction network in stage-I. Intuitively, the re-ID performance will be higher if the encoded features are more representative. Euclidean distance is used to calculate the extracted features after $l_2$-norm normalization [6]. As shown in the top rows of Table 2, our BodyROI7 model achieves 0.338 and 0.355 (with PCA) rank-1 performance, higher than our WholeBody model, which is in accordance with the sampling results in Sec. 4.5. Besides, our method can achieve comparable performance with the unsupervised baseline methods, which indicates that our encoder can extract not only generative but also discriminative features.

**Using generated virtual image pairs.** We use the generated image pairs to train the re-ID model and use the re-ID performance to evaluate our generation framework in an indirect manner. We first generate the VM re-ID dataset consisting of 500 identities with 24 images for each ID as illustrated in Fig. 8. For each identity, we randomly sample one foreground feature and 24 background features and randomly select 24 pose keypoint heatmaps from the Market-1501 training data. Then, we use the same re-ID model and training procedure as in [6], but with different training data. As shown in the bottom rows of Table 2, using our VM data the model can achieve the rank-1 performance 0.338 which is comparable to the model trained using another Duke re-ID dataset. When using the post-processing progressive unsupervised learning (PUL) proposed in [6], the rank-1 performance is improved to 0.369. Additionally, using our VM data, we can train a metric model, *e.g.* KISSME [14], and further improve the rank-1 performance to 0.375. Compared to the model trained using *CUHK03* (rank-1 0.300) or *Duke* (rank-1 0.361) re-ID dataset with expensive human annotations, our method achieves better performance using

only Market dataset without identity labels. These results show that our disentangled generated images are similar to the real data and can be further beneficial to re-ID tasks.

## 5. Conclusion

We propose a novel two-stage pipeline for addressing the disentangled person image generation task. Stage-I disentangles and encodes three modes of variation in the input image, namely foreground, background and pose, into embedding features then decodes them back to an image using a multi-branched reconstruction network. Stage-II learns mapping functions in an adversarial manner for mapping noise distributions to feature embedding distributions guided by the decoders learned in stage-I. Our experimental results on two datasets demonstrate the ability of our method to manipulate the foreground, background and pose of an input image, and sample new embedding features to generate intended manipulations of the different factors, thus providing more control over the generation process. In the future, we plan to apply our method to faces and rigid object images with different types of structure.

## Acknowledgments

## References

[1] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein gan. In *ICLR*, 2017. 1, 5

[2] X. Chen, Y. Duan, R. Houthooft, J. Schulman, I. Sutskever, and P. Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *NIPS*, 2016. 1, 3

[3] B. Cheung, J. A. Livezey, A. K. Bansal, and B. A. Olshausen. Discovering hidden factors of variation in deep networks. In *ICLR workshop*, 2015. 3

[4] C. Doersch, A. Gupta, and A. A. Efros. Unsupervised visual representation learning by context prediction. In *ICCV*, 2015. 3

[5] P. Dollár, V. Rabaud, and S. J. Belongie. Learning to traverse image manifolds. In *NIPS*, pages 361–368, 2007. 4

[6] H. Fan, L. Zheng, and Y. Yang. Unsupervised person re-identification: Clustering and fine-tuning. *arXiv preprint arXiv:1705.10444*, 2017. 8

[7] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, pages 2672–2680, 2014. 1, 2

[8] S. Gupta, J. Hoffman, and J. Malik. Cross modal distillation for supervision transfer. In *CVPR*, pages 2827–2836, 2016. 4

[9] G. Huang, Z. Liu, K. Q. Weinberger, and L. van der Maaten. Densely connected convolutional networks. In *CVPR*, 2017. 5

[10] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 6

[11] Y. Kim, K. Zhang, A. M. Rush, Y. LeCun, et al. Adversarially regularized autoencoders for generating discrete structures. *arXiv preprint arXiv:1706.04223*, 2017. 2

[12] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv*, 1412.6980, 2014. 11

[13] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 1, 2, 7, 8

[14] M. Köstinger, M. Hirzer, P. Wohlhart, P. M. Roth, and H. Bischof. Large Scale Metric Learning from Equivalence Constraints. In *CVPR*, 2012. 8

[15] C. Lassner, G. Pons-Moll, and P. V. Gehler. A generative model of people in clothing. In *ICCV*, 2017. 3

[16] D. Lee, S. Yun, S. Choi, H. Yoo, M.-H. Yang, and S. Oh. Unsupervised holistic image generation from key local patches. *arXiv preprint arXiv:1703.10730*, 2017. 3

[17] W. Li, R. Zhao, T. Xiao, and X. Wang. Deepreid: Deep filter pairing neural network for person re-identification. In *CVPR*, pages 152–159, 2014. 8

[18] S. Liao, Y. Hu, X. Zhu, and S. Z. Li. Person re-identification by local maximal occurrence representation and metric learning. In *CVPR*, 2015. 8

[19] Z. C. Lipton and S. Tripathi. Precise recovery of latent vectors from generative adversarial networks. In *ICLR workshop*, 2017. 7

[20] Z. Liu, P. Luo, S. Qiu, X. Wang, and X. Tang. Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In *CVPR*, pages 1096–1104, 2016. 6

[21] L. Ma, J. Xu, Q. Sun, B. Schiele, T. Tuytelaars, and L. Van Gool. Pose guided person image generation. In *NIPS*, 2017. 1, 3, 4, 5, 6

[22] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015. 2

[23] T. B. Moeslund, A. Hilton, and V. Krüger. A survey of advances in vision-based human motion capture and analysis. *Computer vision and image understanding*, 104(2):90–126, 2006. 2

[24] M. Noroozi and P. Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *ECCV*, 2016. 3

[25] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *ICLR*, 2016. 1, 2, 7, 8

[26] S. Reed, Z. Akata, S. Mohan, S. Tenka, B. Schiele, and H. Lee. Learning what and where to draw. In *NIPS*, 2016. 3

[27] S. Reed, A. van den Oord, N. Kalchbrenner, V. Bapst, M. Botvinick, and N. de Freitas. Generating interpretable images with controllable structure. Technical report, 2016. 3

[28] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio. Fitnets: Hints for thin deep nets. In *ICLR*, 2015. 4

[29] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 234–241. Springer, 2015. 4, 5, 12

[30] L. K. Saul and S. T. Roweis. Think globally, fit locally: unsupervised learning of low dimensional manifolds. *Journal of machine learning research*, 4(Jun):119–155, 2003. 4

[31] W. Shi, J. Caballero, F. Huszar, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. 6

[32] L. Tran, X. Yin, and X. Liu. Disentangled representation learning gan for pose-invariant face recognition. In *CVPR*, 2017. 3

[33] A. Van Den Oord, N. Kalchbrenner, and K. Kavukcuoglu. Pixel recurrent neural networks. In *ICML*, 2016. 1, 2

[34] K. Q. Weinberger and L. K. Saul. Unsupervised learning of image manifolds by semidefinite programming. *International journal of computer vision*, 70(1):77–90, 2006. 4

[35] K. Yu, T. Zhang, and Y. Gong. Nonlinear learning using local coordinate coding. In *NIPS*, pages 2223–2231, 2009. 4

[36] H. Zhang, T. Xu, H. Li, S. Zhang, X. Huang, X. Wang, and D. Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *ICCV*, 2017. 4

[37] B. Zhao, X. Wu, Z. Cheng, H. Liu, and J. Feng. Multi-view image generation from a single-view. *arXiv*, 1704.04886, 2017. 3

[38] H. Zhao, M. Tian, S. Sun, J. Shao, J. Yan, S. Yi, X. Wang, and X. Tang. Spindle net: Person re-identification with human body region guided feature decomposition and fusion. In *CVPR*, pages 1077–1085, 2017. 4

[39] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian. Scalable person re-identification: A benchmark. In *ICCV*, pages 1116–1124, 2015. 6, 8

[40] Z. Zheng, L. Zheng, and Y. Yang. Unlabeled samples generated by gan improve the person re-identification baseline in vitro. In *ICCV*, 2017. 2, 8

[41] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 3

# Supplementary materials

This supplementary material includes additional details regarding the network architecture (§A) and training (§B), as well as extended results for image manipulation (§C), pose-guided person image generation (§D), inverse interpolation (§E) and image sampling (§F), respectively.

## A. Network architecture

In this section, we provide details regarding the network architectures in our two-stage framework used on the Market-1501 dataset. Fig. 10 shows 4 network architectures used at stage-I: 1) FG encoder consists of 5 convolutional residual blocks; 2) BG encoder consists of 5 convolutional residual blocks; 3) FG & BG decoder follows a "U-net"-based architecture; 4) Pose auto-encoder follows a fully-connected auto-encoder architecture. Fig. 9 shows the network architecture of the mapping functions $\Phi$ used at stage-II. It contains 4 fully-connected residual modules.
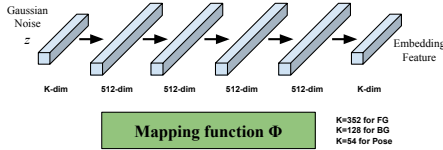


Figure 9: Network architecture of the mapping functions for FG, BG and Pose in stage-II.

## B. Training details

On Market-1501, our method is applied to disentangle the image into three factors: foreground, background and pose. We train the foreground and background models with a mini-batch of size 16 for $\sim70k$ iterations at stage-I and with a mini-batch of size 32 for $\sim30k$ iterations at stage-II. The pose models are trained with a mini-batch of size 64 for $\sim30k$ iterations at stage-I and with a mini-batch of size 32 for $\sim60k$ iterations at stage-II.

DeepFashion data contain clean background, therefore, our method is applied to disentangle the image into only two factors: appearance (*i.e.* foreground) and pose. We train the appearance model with a minibatch of size 6 for $\sim100k$ iterations at stage-I and with a minibatch of size 16 for $\sim60k$ iterations at stage-II. The pose models are trained with a minibatch of size 32 for $\sim30k$ iterations at stage-I and with a minibatch of size 32 for $\sim60k$ iterations at stage-II.

On both datasets, we use the Adam optimizer [12] with weights $\beta_1 = 0.5$ and $\beta_2 = 0.999$. The initial learning rate is set to $2e$-5. For adversarial training, we optimize the discriminator and generator alternatively.

## C. Image manipulation results

In Fig. 11 and Fig. 12, we provide results on appearance sampling and pose sampling for the DeepFashion dataset as an extension of Fig. 1 in the main paper. For each factor, we sample the embedding feature from Gaussian noise and fix the other factors by using the embedding feature extracted from the real data as explained in Sec. 4.2 in the main paper.

## D. Pose-guided person image generation results

For pose-guided person image generation, we provide more generated results. As an extension of Fig. 5 in the main paper, Fig. 13 shows the generated images of one appearance with various real poses selected randomly from DeepFashion.

## E. Inverse interpolation results

In this section, we provide more inverse interpolation results in Fig.14 as an extension of Fig. 6 in the main paper. For two images $x_1$ and $x_2$, we find the corresponding Gaussian codes $z_1$ and $z_2$ as explained in the Sec. 4.4 of the main paper. As shown in Fig. 14(a)(b), our method successfully generates the intermediate states between two images of the same person. Note that, the inverse interpolation between two images of different persons is more challenging (see Fig. 14(c)) since we need to interpolate both the appearance and pose.

## F. Image sampling results

We also give more sampling results as extensions of Fig.7 in the main paper. Fig. 15 shows the sampling results (a-e) and real images (f) on Market-1501 dataset. VAE generates blurry images and DCGAN sharp but unrealistic person images. In contrast, our model generates more realistic images (c)(d)(e). By comparing (d) and (c), we observe that our model using body ROI generates more sharp and realistic images whose colors on each body part are more natural. By comparing (e) and (d), we see that when sampling foreground and background but using the real pose keypoints randomly selected from the training data, we generate better results.
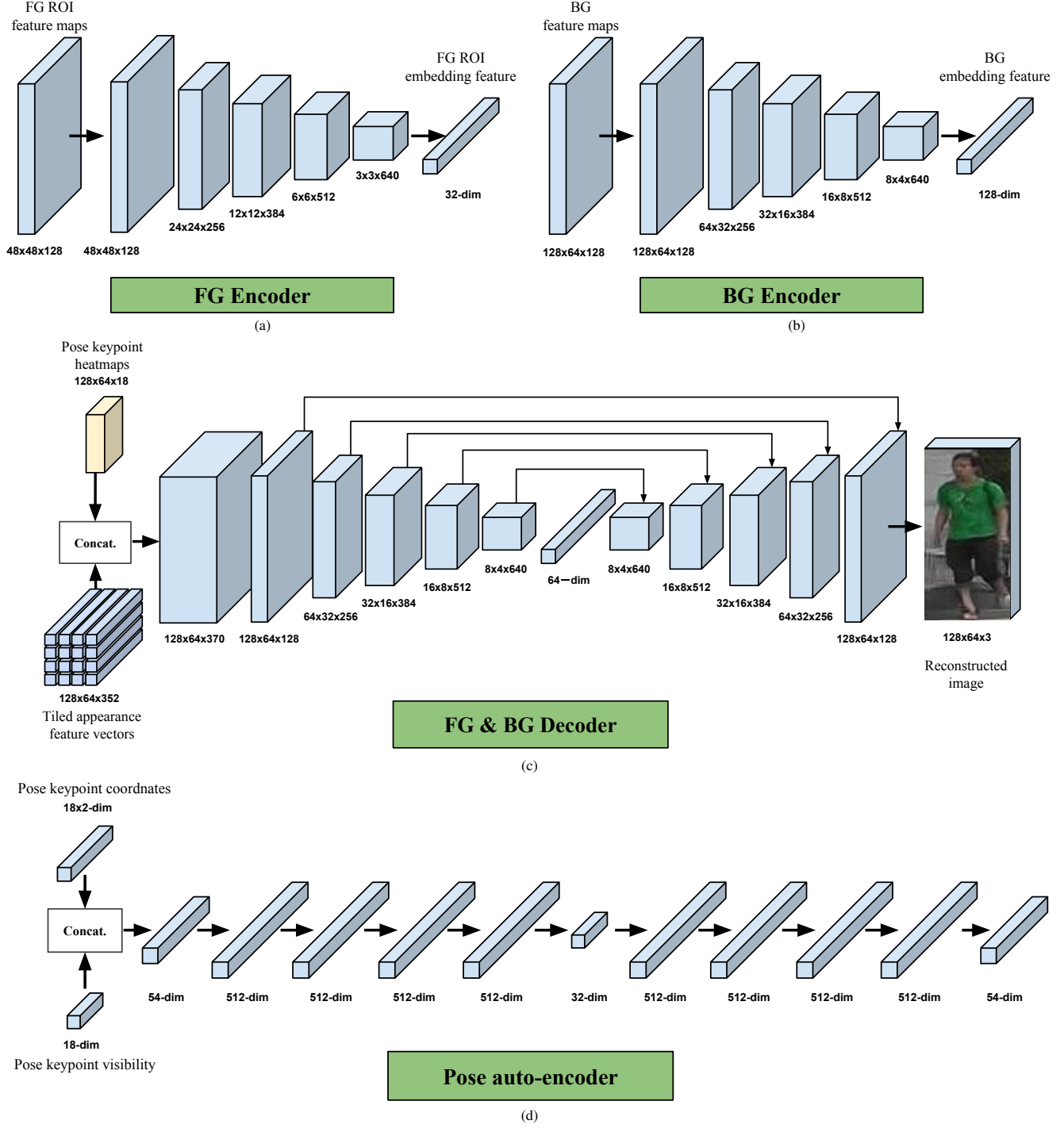
Figure 10: Network architectures of stage-I. (a) FG encoder, fed with the extracted 7 FG body ROI feature maps and outputting 7 FG embedding features of 32-dim after 5 convolutional residual blocks. (b) BG encoder, fed with the BG feature maps and outputting a BG embedding feature of 128-dim after 5 convolutional residual blocks. (c) FG and BG decoder, fed with the concatenated appearance and pose feature maps and outputting the generated image after the "U-net"-based [29] architecture. (d) Pose auto-encoder, fed with the concatenated keypoint coordinates and visibility vector and outputting the reconstructed vector after the auto-encoder.

Condition pose                                Generated images
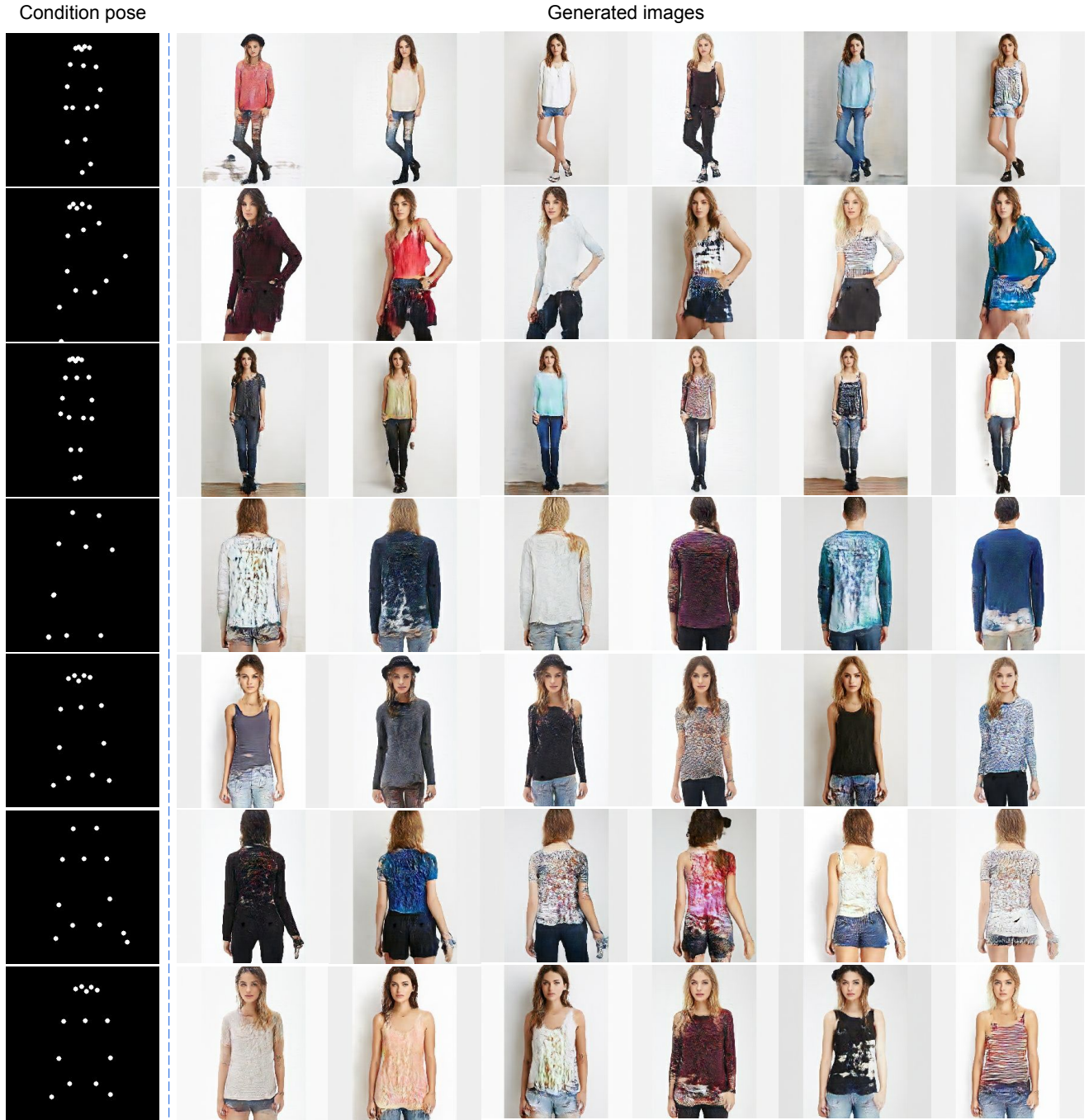


Figure 11: Appearance sampling (fixed Pose) results on the DeepFashion dataset. In each row, 6 different appearance factors are sampled from Gaussian noise and the pose factor is fixed to a real one.
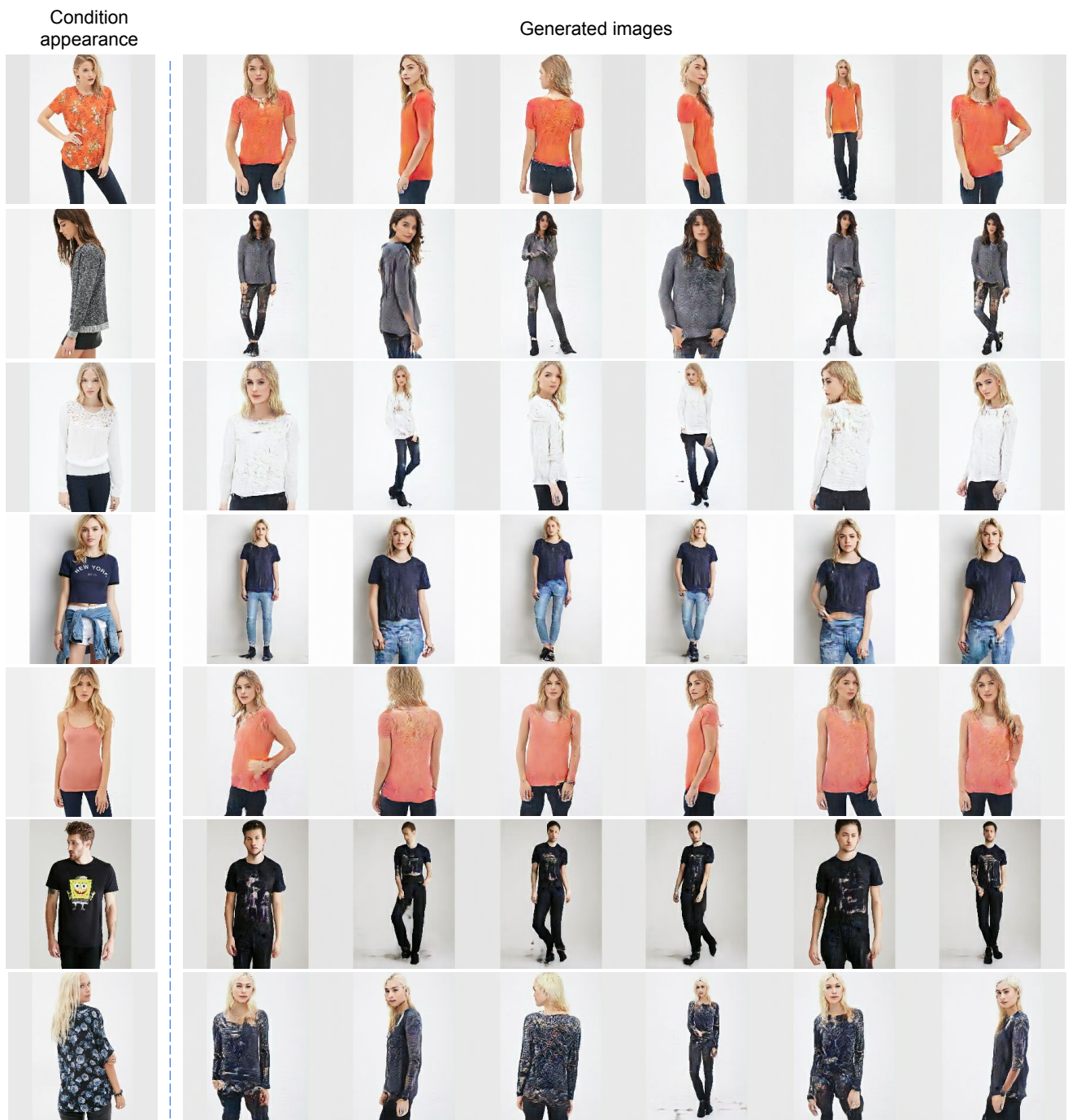
13

Figure 12: Pose sampling (fixed Appearance) results on the DeepFashion dataset. In each row, 6 different pose factors are sampled from Gaussian noises and the appearance factor is fixed to a real one.

Figure 13: Generated results for one appearance with various poses on the DeepFashion dataset.

x1  z1  z1+d  ...  →  ...  z1+(n-1)d  z2  x2

(a)

x1  →  x2  →  x3

(b)

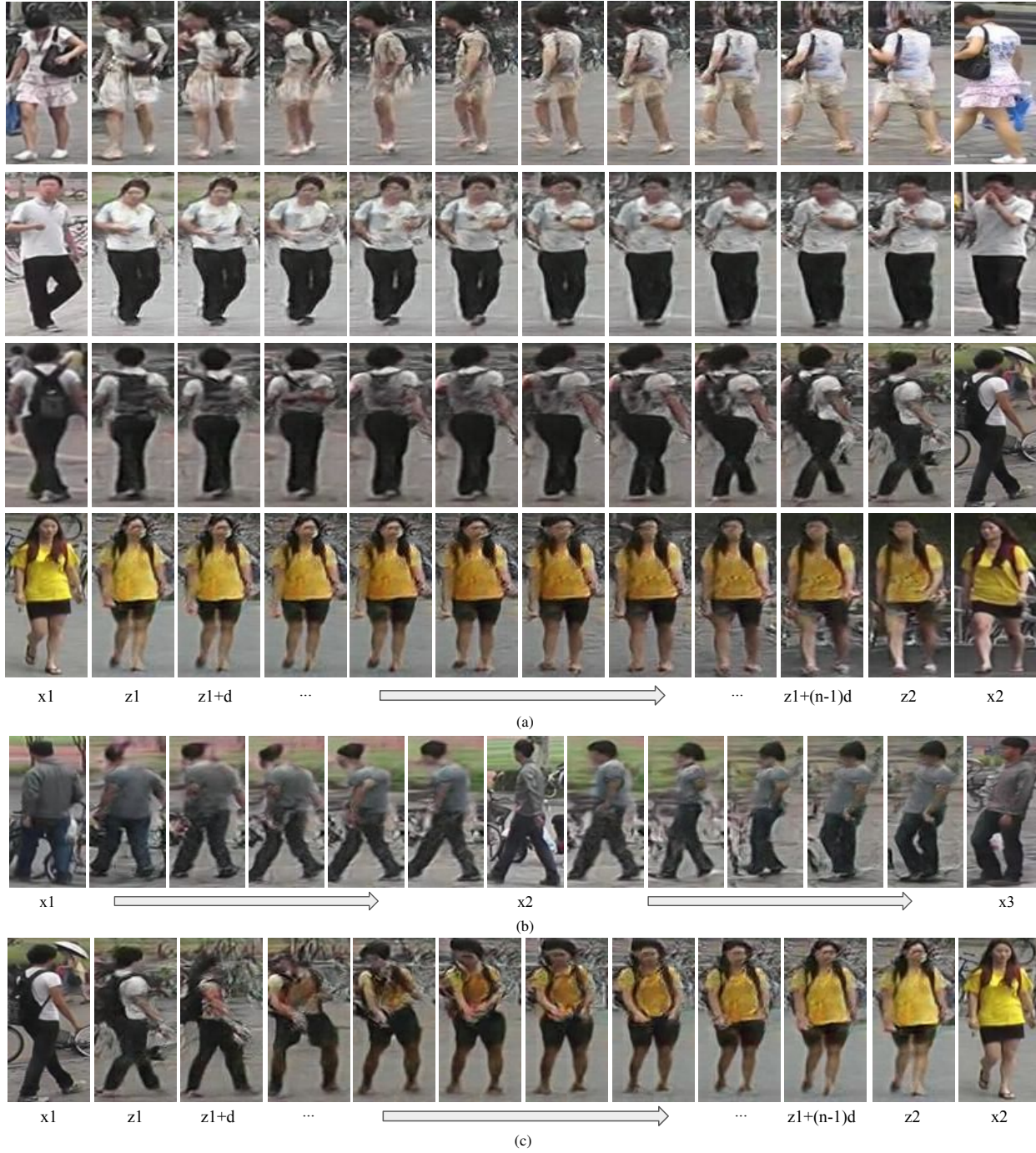x1  z1  z1+d  ...  →  ...  z1+(n-1)d  z2  x2

(c)

Figure 14: Inverse interpolation results on Market-1501. (a) Interpolation between two images of the same person. (b) Interpolation between three images of the same person. (c) Interpolation between two images of different persons.

(a) Vanilla VAE    (b) Vanilla DCGAN    (c) Ours - Whole Body

(d) Ours - BodyROI7    (e) Ours - BodyROI7 with real pose from training set    (f) Real data

Figure 15: Sampling results. (a) Vanilla VAE; (b) Vanilla DCGAN; (c) Ours - Whole Body; (d) Ours - BodyROI7; (e) Ours - BodyROI7 pose with real pose from training set; (f) Real data.