

Deformable GANs for Pose-based Human Image Generation

Aliaksandr Siarohin¹, Enver Sangineto¹, Stéphane Lathuilière², and Nicu Sebe¹

¹DISI, University of Trento, Italy, ²Inria Grenoble Rhone-Alpes, France

{aliaksandr.siarohin,enver.sangineto,niculae.sebe}@unitn.it, stephane.lathuiliere@inria.fr

Abstract

In this paper we address the problem of generating person images conditioned on a given pose. Specifically, given an image of a person and a target pose, we synthesize a new image of that person in the novel pose. In order to deal with pixel-to-pixel misalignments caused by the pose differences, we introduce deformable skip connections in the generator of our Generative Adversarial Network. Moreover, a nearest-neighbour loss is proposed instead of the common L_1 and L_2 losses in order to match the details of the generated image with the target image. We test our approach using photos of persons in different poses and we compare our method with previous work in this area showing state-of-the-art results in two benchmarks. Our method can be applied to the wider field of deformable object generation, provided that the pose of the articulated object can be extracted using a keypoint detector.

1. Introduction

In this paper we deal with the problem of generating images where the foreground object changes because of a viewpoint variation or a deformable motion, such as the articulated human body. Specifically, inspired by Ma et al. [12], our goal is to generate a human image conditioned on two different variables: (1) the appearance of a specific person in a given image and (2) the pose of the same person in another image. The task our networks need to solve is to preserve the appearance details (e.g., the texture) contained in the first variable while performing a deformation on the structure of the foreground object according to the second variable. We focus on the human body which is an articulated “object”, important for many applications (e.g., computer-graphics based manipulations or re-identification dataset synthesis). However, our approach can be used with other deformable objects such as human faces or animal bodies, provided that a significant number of keypoints can be automatically extracted from the object of interest in order to represent its pose.

Pose-based human-being image generation is motivated



Figure 1: (a) A typical “rigid” scene generation task, where the conditioning and the output image local structure is well aligned. (b) In a deformable-object generation task, the input and output are not spatially aligned.

by the interest in synthesizing videos [18] with non-trivial human movements or in generating rare poses for human pose estimation [1] or re-identification [22] training datasets. However, most of the recently proposed, deep-network based generative approaches, such as Generative Adversarial Networks (GANs) [3] or Variational Autoencoders (VAEs) [7] do not explicitly deal with the problem of articulated-object generation. Common conditional methods (e.g., conditional GANs or conditional VAEs) can synthesize images whose appearances depend on some conditioning variables (e.g., a label or another image). For instance, Isola et al. [4] recently proposed an “image-to-image translation” framework, in which an input image x is transformed into a second image y represented in another “channel” (see Fig. 1a). However, most of these methods have problems when dealing with large spatial deformations between the conditioning and the target image. For instance, the U-Net architecture used by Isola et al. [4] is based on *skip connections* which help preserving local information between x and y . Specifically, skip connections are used to copy and then concatenate the feature maps of the generator “encoder” (where information is downsam-

pled using convolutional layers) to the generator “decoder” (containing the upconvolutional layers). However, the assumption used in [4] is that x and y are roughly aligned with each other and they represent the same underlying structure. This assumption is violated when the foreground object in y undergoes to large spatial deformations with respect to x (see Fig. 1b). As shown in [12], skip connections cannot reliably cope with misalignments between the two poses.

Ma et al. [12] propose to alleviate this problem using a two-stage generation approach. In the first stage a U-Net generator is trained using a masked L_1 loss in order to produce an intermediate image conditioned on the target pose. In the second stage, a second U-Net based generator is trained using also an adversarial loss in order to generate an appearance difference map which brings the intermediate image closer to the appearance of the conditioning image. In contrast, the GAN-based method we propose in this paper is end-to-end trained by explicitly taking into account pose-related spatial deformations. More specifically, we propose *deformable skip connections* which “move” local information according to the structural deformations represented in the conditioning variables. These layers are used in our U-Net based generator. In order to move information according to a specific spatial deformation, we decompose the overall deformation by means of a set of local affine transformations involving subsets of joints, then we deform the convolutional feature maps of the encoder according to these transformations and we use common skip connections to transfer the transformed tensors to the decoder’s fusion layers. Moreover, we also propose to use a *nearest-neighbour loss* as a replacement of common pixel-to-pixel losses (such as, e.g., L_1 or L_2 losses) commonly used in conditional generative approaches. This loss proved to be helpful in generating local information (e.g., texture) similar to the target image which is not penalized because of small spatial misalignments.

We test our approach using the benchmarks and the evaluation protocols proposed in [12] obtaining higher qualitative and quantitative results in all the datasets. Although tested on the specific human-body problem, our approach makes few human-related assumptions and can be easily extended to other domains involving the generation of highly deformable objects. Our code and our trained models are publicly available¹.

2. Related work

Most common deep-network-based approaches for visual content generation can be categorized as either Variational Autoencoders (VAEs) [7] or Generative Adversarial Networks (GANs) [3]. VAEs are based on probabilistic graphical models and are trained by maximizing a lower

bound of the corresponding data likelihood. GANs are based on two networks, a generator and a discriminator, which are trained simultaneously such that the generator tries to “fool” the discriminator and the discriminator learns how to distinguish between real and fake images.

Isola et al. [4] propose a conditional GAN framework for image-to-image translation problems, where a given scene representation is “translated” into another representation. The main assumption behind this framework is that there exists a spatial correspondence between the low-level information of the conditioning and the output image. VAEs and GANs are combined in [20] to generate realistic-looking multi-view clothes images from a single-view input image. The target view is filled to the model via a viewpoint label as *front* or *left side* and a two-stage approach is adopted: pose integration and image refinement. Adopting a similar pipeline, Lassner et al. [8] generate images of people with different clothes in a given pose. This approach is based on a costly annotation (fine-grained segmentation with 18 clothing labels) and a complex 3D pose representation.

Ma et al. [12] propose a more general approach which allows to synthesize person images in any arbitrary pose. Similarly to our proposal, the input of their model is a conditioning image of the person and a target new pose defined by 18 joint locations. The target pose is described by means of binary maps where small circles represent the joint locations. Similarly to [8, 20], the generation process is split in two different stages: pose generation and texture refinement. In contrast, in this paper we show that a single-stage approach, trained end-to-end, can be used for the same task obtaining higher qualitative results.

Jaderberg et al. [5] propose a spatial transformer layer, which learns how to transform a feature map in a “canonical” view, conditioned on the feature map itself. However only a global, parametric transformation can be learned (e.g., a global affine transformation), while in this paper we deal with non-parametric deformations of articulated objects which cannot be described by means of a unique global affine transformation.

Generally speaking, U-Net based architectures are frequently adopted for pose-based person-image generation tasks [8, 12, 18, 20]. However, common U-Net skip connections are not well-designed for large spatial deformations because local information in the input and in the output images is not aligned (Fig. 1). In contrast, we propose deformable skip connections to deal with this misalignment problem and “shuttle” local information from the encoder to the decoder driven by the specific pose difference. In this way, differently from previous work, we are able to simultaneously generate the overall pose and the texture-level refinement.

Finally, our nearest-neighbour loss is similar to the perceptual loss proposed in [6] and to the style-transfer spatial-

¹<https://github.com/AliaksandrSiarohin/pose-gan>

analogy approach recently proposed in [9]. However, the perceptual loss, based on an element-by-element difference computed in the feature map of an external classifier [6], does not take into account spatial misalignments. On the other hand, the patch-based similarity, adopted in [9] to compute a dense feature correspondence, is very computationally expensive and it is not used as a loss.

3. The network architectures

In this section we describe the architectures of our generator (G) and discriminator (D) and the proposed deformable skip connections. We first introduce some notation. At testing time our task, similarly to [12], consists in generating an image \hat{x} showing a person whose appearance (e.g., clothes, etc.) is similar to an input, conditioning image x_a but with a body pose similar to $P(x_b)$, where x_b is a different image of the same person and $P(x) = (\mathbf{p}_1, \dots, \mathbf{p}_k)$ is a sequence of k 2D points describing the locations of the human-body joints in x . In order to allow a fair comparison with [12], we use the same number of joints ($k = 18$) and we extract $P()$ using the same Human Pose Estimator (HPE) [1] used in [12]. Note that this HPE is used both at testing and at training time, meaning that we do not use manually-annotated poses and the so extracted joint locations may have some localization errors or missing detections/false positives.

At training time we use a dataset $\mathcal{X} = \{(x_a^{(i)}, x_b^{(i)})\}_{i=1, \dots, N}$ containing pairs of conditioning-target images of the same person in different poses. For each pair (x_a, x_b) , a conditioning and a target pose $P(x_a)$ and $P(x_b)$ is extracted from the corresponding image and represented using two tensors $H_a = H(P(x_a))$ and $H_b = H(P(x_b))$, each composed of k heat maps, where H_j ($1 \leq j \leq k$) is a 2D matrix of the same dimension as the original image. If \mathbf{p}_j is the j -th joint location, then:

$$H_j(\mathbf{p}) = \exp\left(-\frac{\|\mathbf{p} - \mathbf{p}_j\|^2}{\sigma^2}\right), \quad (1)$$

with $\sigma = 6$ pixels (chosen with cross-validation). Using blurring instead of a binary map is useful to provide widespread information about the location \mathbf{p}_j .

The generator G is fed with: (1) a noise vector z , drawn from a noise distribution \mathcal{Z} and implicitly provided using dropout [4] and (2) the triplet (x_a, H_a, H_b) . Note that, at testing time, the target pose is known, thus $H(P(x_b))$ can be computed. Note also that the joint locations in x_a and H_a are spatially aligned (by construction), while in H_b they are different. Hence, differently from [12, 4], H_b is not concatenated with the other input tensors. Indeed the convolutional-layer units in the encoder part of G have a small receptive field which cannot capture large spatial displacements. For instance, a large movement of a body limb in x_b with respect to x_a , is represented in different locations

in x_a and H_b which may be too far apart from each other to be captured by the receptive field of the convolutional units. This is emphasized in the first layers of the encoder, which represent low-level information. Therefore, the convolutional filters cannot simultaneously process texture-level information (from x_a) and the corresponding pose information (from H_b).

For this reason we independently process x_a and H_a from H_b in the encoder. Specifically, x_a and H_a are concatenated and processed using a convolutional stream of the encoder while H_b is processed by means of a second convolutional stream, without sharing the weights (Fig. 2). The feature maps of the first stream are then fused with the layer-specific feature maps of the second stream in the decoder after a pose-driven spatial deformation performed by our deformable skip connections (see Sec. 3.1).

Our discriminator network is based on the conditional, fully-convolutional discriminator proposed by Isola et al. [4]. In our case, D takes as input 4 tensors: (x_a, H_a, y, H_b) , where either $y = x_b$ or $y = \hat{x} = G(z, x_a, H_a, H_b)$ (see Fig. 2). These four tensors are concatenated and then given as input to D . The discriminator's output is a scalar value indicating its confidence on the fact that y is a real image.

3.1. Deformable skip connections

As mentioned above and similarly to [4], the goal of the deformable skip connections is to "shuttle" local information from the encoder to the decoder part of G . The local information to be transferred is, generally speaking, contained in a tensor F , which represents the feature map activations of a given convolutional layer of the encoder. However, differently from [4], we need to "pick" the information to shuttle taking into account the object-shape deformation which is described by the difference between $P(x_a)$ and $P(x_b)$. To do so, we decompose the global deformation in a set of local affine transformations, defined using subsets of joints in $P(x_a)$ and $P(x_b)$. Using these affine transformations and local masks constructed using the specific joints, we deform the content of F and then we use common skip connections to copy the transformed tensor and concatenate it with the corresponding tensor in the destination layer (see Fig. 2). Below we describe in more detail the whole pipeline.

Decomposing an articulated body in a set of rigid sub-parts. The human body is an articulated "object" which can be roughly decomposed into a set of rigid sub-parts. We chose 10 sub-parts: the head, the torso, the left/right upper/lower arm and the left/right upper/lower leg. Each of them corresponds to a subset of the 18 joints defined by the HPE [1] we use for extracting $P()$. Using these joint locations we can define rectangular regions which enclose the specific body part. In case of the head, the region is simply chosen to be the axis-aligned enclosing rectangle of all the corresponding joints. For the torso, which is

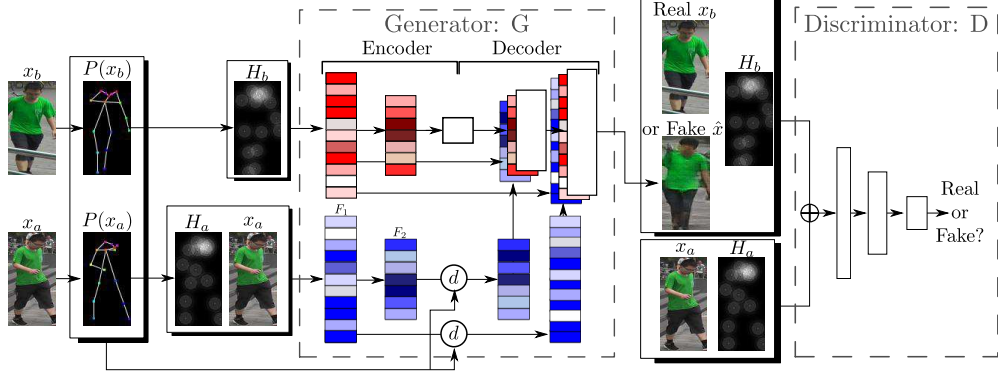


Figure 2: A schematic representation of our network architectures. For the sake of clarity, in this figure we depict $P(\cdot)$ as a skeleton and each tensor H as the average of its component matrices H_j ($1 \leq j \leq k$). The white rectangles in the decoder represent the feature maps directly obtained using up-convolutional filters applied to the previous-layer maps. The reddish rectangles represent the feature maps “shuttled” by the skip connections from the H_b stream. Finally, blueish rectangles represent the deformed tensors $d(F)$ “shuttled” by the deformable skip connections from the (x_a, H_a) stream.

the largest area, we use a region which includes the whole image, in such a way to shuttle texture information for the background pixels. Concerning the body limbs, each limb corresponds to only 2 joints. In this case we define a region to be a rotated rectangle whose major axis (r_1) corresponds to the line between these two joints, while the minor axis (r_2) is orthogonal to r_1 and with a length equal to one third of the mean of the torso’s diagonals (this value is used for all the limbs). In Fig. 3 we show an example. Let $R_h^a = \{\mathbf{p}_1, \dots, \mathbf{p}_4\}$ be the set of the 4 rectangle corners in x_a defining the h -th body region ($1 \leq h \leq 10$). Note that these 4 corner points are not joint locations. Using R_h^a we can compute a binary mask $M_h(\mathbf{p})$ which is zero everywhere except those points \mathbf{p} lying inside R_h^a . Moreover, let $R_h^b = \{\mathbf{q}_1, \dots, \mathbf{q}_4\}$ be the corresponding rectangular region in x_b . Matching the points in R_h^a with the corresponding points in R_h^b we can compute the parameters of a body-part specific affine transformation (see below). In either x_a or x_b , some of the body regions can be occluded, truncated by the image borders or simply miss-detected by the HPE. In this case we leave the corresponding region R_h empty and the h -th affine transform is not computed (see below).

Note that our body-region definition is the only human-specific part of the proposed approach. However, similar regions can be easily defined using the joints of other articulated objects such as those representing an animal body or a human face.

Computing a set of affine transformations. During the forward pass (i.e., both at training and at testing time) we decompose the global deformation of the conditioning pose with respect to the target pose by means of a set of local affine transformations, one per body region. Specifically, given R_h^a in x_a and R_h^b in x_b (see above), we compute the 6 parameters \mathbf{k}_h of an affine transformation $f_h(\cdot; \mathbf{k}_h)$ using Least Squares Error:

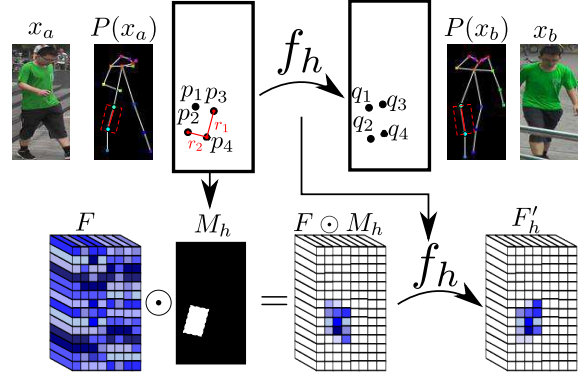


Figure 3: For each specific body part, an affine transformation f_h is computed. This transformation is used to “move” the feature-map content corresponding to that body part.

$$\min_{\mathbf{k}_h} \sum_{\mathbf{p}_j \in R_h^a, \mathbf{q}_j \in R_h^b} \|\mathbf{q}_j - f_h(\mathbf{p}_j; \mathbf{k}_h)\|_2^2 \quad (2)$$

The parameter vector \mathbf{k}_h is computed using the original image resolution of x_a and x_b and then adapted to the specific resolution of each involved feature map F . Similarly, we compute scaled versions of each M_h . In case either R_h^a or R_h^b is empty (i.e., when any of the specific body-region joints has not been detected using the HPE, see above), then we simply set M_h to be a matrix with all elements equal to 0 (f_h is not computed).

Note that $(f_h(\cdot), M_h)$ and their lower-resolution variants need to be computed only once per each pair of real images $(x_a, x_b) \in \mathcal{X}$ and, in case of the training phase, this can be done before starting training the networks (but in our current implementation this is done on the fly).

Combining affine transformations to approximate the object deformation. Once $(f_h(\cdot), M_h)$, $h = 1, \dots, 10$

are computed for the specific spatial resolution of a given tensor F , the latter can be transformed in order to approximate the global pose-dependent deformation. Specifically, we first compute for each h :

$$F'_h = f_h(F \odot M_h), \quad (3)$$

where \odot is a point-wise multiplication and $f_h(F(\mathbf{p}))$ is used to “move” all the channel values of F corresponding to point \mathbf{p} . Finally, we merge the resulting tensors using:

$$d(F(\mathbf{p}, c)) = \max_{h=1, \dots, 10} F'_h(\mathbf{p}, c), \quad (4)$$

where c is a specific channel. The rationale behind Eq. 4 is that, when two body regions partially overlap each other, the final deformed tensor $d(F)$ is obtained by picking the maximum-activation values. Preliminary experiments performed using average pooling led to slightly worse results.

4. Training

D and G are trained using a combination of a standard conditional adversarial loss \mathcal{L}_{cGAN} with our proposed nearest-neighbour loss \mathcal{L}_{NN} . Specifically, in our case \mathcal{L}_{cGAN} is given by:

$$\mathcal{L}_{cGAN}(G, D) = \mathbb{E}_{(x_a, x_b) \in \mathcal{X}} [\log D(x_a, H_a, x_b, H_b)] + \mathbb{E}_{(x_a, x_b) \in \mathcal{X}, z \in \mathcal{Z}} [\log(1 - D(x_a, H_a, \hat{x}, H_b))], \quad (5)$$

where $\hat{x} = G(z, x_a, H_a, H_b)$.

Previous works on conditional GANs combine the adversarial loss with either an L_2 [13] or an L_1 -based loss [4, 12] which is used only for G . For instance, the L_1 distance computes a pixel-to-pixel difference between the generated and the real image, which, in our case, is:

$$L_1(\hat{x}, x_b) = \|\hat{x} - x_b\|_1. \quad (6)$$

However, a well-known problem behind the use of L_1 and L_2 is the production of blurred images. We hypothesize that this is also due to the inability of these losses to tolerate small spatial misalignments between \hat{x} and x_b . For instance, suppose that \hat{x} , produced by G , is visually plausible and semantically similar to x_b , but the texture details on the clothes of the person in the two compared images are not pixel-to-pixel aligned. Both the L_1 and the L_2 loss will penalize this inexact pixel-level alignment, although not semantically important from the human point of view. Note that these misalignments do *not* depend on the global deformation between x_a and x_b , because \hat{x} is supposed to have the same pose as x_b . In order to alleviate this problem, we propose to use a *nearest-neighbour* loss \mathcal{L}_{NN} based on the following definition of image difference:

$$L_{NN}(\hat{x}, x_b) = \sum_{\mathbf{p} \in \hat{x}} \min_{\mathbf{q} \in \mathcal{N}(\mathbf{p})} \|g(\hat{x}(\mathbf{p})) - g(x_b(\mathbf{q}))\|_1, \quad (7)$$

where $\mathcal{N}(\mathbf{p})$ is a $n \times n$ local neighbourhood of point \mathbf{p} (we use 5×5 and 3×3 neighbourhoods for the DeepFashion and the Market-1501 dataset, respectively, see Sec. 6). $g(x(\mathbf{p}))$ is a vectorial representation of a patch around point \mathbf{p} in image x , obtained using convolutional filters (see below for more details). Note that $L_{NN}()$ is not a metrics because it is not symmetric. In order to efficiently compute Eq. 7, we compare patches in \hat{x} and x_b using their representation ($g()$) in a convolutional map of an externally trained network. In more detail, we use VGG-19 [15], trained on ImageNet and, specifically, its second convolutional layer (called *conv1.2*). The first two convolutional maps in VGG-19 (*conv1.1* and *conv1.2*) are both obtained using a convolutional stride equal to 1. For this reason, the feature map (C_x) of an image x in *conv1.2* has the same resolution of the original image x . Exploiting this fact, we compute the nearest-neighbour field directly on *conv1.2*, without losing spatial precision. Hence, we define: $g(x(\mathbf{p})) = C_x(\mathbf{p})$, which corresponds to the vector of all the channel values of C_x with respect to the spatial position \mathbf{p} . $C_x(\mathbf{p})$ has a receptive field of 5×5 in x , thus effectively representing a patch of dimension 5×5 using a cascade of two convolutional filters. Using C_x , Eq. 7 becomes:

$$L_{NN}(\hat{x}, x_b) = \sum_{\mathbf{p} \in \hat{x}} \min_{\mathbf{q} \in \mathcal{N}(\mathbf{p})} \|C_{\hat{x}}(\mathbf{p}) - C_{x_b}(\mathbf{q})\|_1, \quad (8)$$

In the Supplementary Material we show how Eq. 8 can be efficiently implemented using GPU-based parallel computing. The final L_{NN} -based loss is:

$$\mathcal{L}_{NN}(G) = \mathbb{E}_{(x_a, x_b) \in \mathcal{X}, z \in \mathcal{Z}} L_{NN}(\hat{x}, x_b). \quad (9)$$

Combining Eq. 5 and Eq. 9 we obtain our objective:

$$G^* = \arg \min_G \max_D \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{NN}(G), \quad (10)$$

with $\lambda = 0.01$ used in all our experiments. The value of λ is small because it acts as a normalization factor in Eq. 8 with respect to the number of channels in C_x and the number of pixels in \hat{x} (more details in the Supplementary Material).

5. Implementation details

We train G and D for 90k iterations, with the Adam optimizer (learning rate: $2 * 10^{-4}$, $\beta_1 = 0.5$, $\beta_2 = 0.999$). Following [4] we use instance normalization [17]. In the following we denote with: (1) C_m^s a convolution-ReLU layer

with m filters and stride s , (2) CN_m^s the same as C_m^s with instance normalization before ReLU and (3) CD_m^s the same as CN_m^s with the addition of dropout at rate 50%. Differently from [4], we use dropout only at training time. The encoder part of the generator is given by two streams (Fig. 2), each of which is composed of the following sequence of layers:

$$CN_{64}^1 - CN_{128}^2 - CN_{256}^2 - CN_{512}^2 - CN_{512}^2 - CN_{512}^2.$$

The decoder part of the generator is given by:

$$CD_{512}^2 - CD_{512}^2 - CD_{512}^2 - CN_{256}^2 - CN_{128}^2 - C_3^1.$$

In the last layer, ReLU is replaced with \tanh .

The discriminator architecture is:

$$C_{64}^2 - C_{128}^2 - C_{256}^2 - C_{512}^2 - C_1^2,$$

where the ReLU of the last layer is replaced with *sigmoid*.

The generator for the DeepFashion dataset has one additional convolution block (CN_{512}^2) both in the encoder and in the decoder, because images in this dataset have a higher resolution.

6. Experiments

Datasets The person re-identification Market-1501 dataset [21] contains 32,668 images of 1,501 persons captured from 6 different surveillance cameras. This dataset is challenging because of the low-resolution images (128×64) and the high diversity in pose, illumination, background and viewpoint. To train our model, we need pairs of images of the same person in two different poses. As this dataset is relatively noisy, we first automatically remove those images in which no human body is detected using the HPE, leading to 263,631 training pairs. For testing, following [12], we randomly select 12,000 pairs. No person is in common between the training and the test split.

The DeepFashion dataset (*In-shop Clothes Retrieval Benchmark*) [11] is composed of 52,712 clothes images, matched each other in order to form 200,000 pairs of identical clothes with two different poses and/or scales of the persons wearing these clothes. The images have a resolution of 256×256 pixels. Following the training/test split adopted in [12], we create pairs of images, each pair depicting the same person with identical clothes but in different poses. After removing those images in which the HPE does not detect any human body, we finally collect 89,262 pairs for training and 12,000 pairs for testing.

Metrics Evaluation in the context of generation tasks is a problem in itself. In our experiments we adopt a redundancy of metrics and a user study based on human judgments. Following [12], we use Structural Similarity (*SSIM*) [19], Inception Score (*IS*) [14] and their corresponding masked versions *mask-SSIM* and *mask-IS* [12]. The latter are obtained by masking-out the image background and the rationale behind this is that, since no background information of the

target image is input to G , the network cannot guess what the target background looks like. Note that the evaluation masks we use to compute both the mask-IS and the mask-SSIM values do not correspond to the masks ($\{M_h\}$) we use for training. The evaluation masks have been built following the procedure proposed in [12] and adopted in that work for both training and evaluation. Consequently, the mask-based metrics may be biased in favor of their method. Moreover, we observe that the IS metrics [14], based on the entropy computed over the classification neurons of an external classifier [16], is not very suitable for domains with only one object class. For this reason we propose to use an additional metrics that we call Detection Score (*DS*). Similarly to the classification-based metrics (*FCN-score*) used in [4], DS is based on the detection outcome of the state-of-the-art object detector SSD [10], trained on Pascal VOC 07 [2] (and not fine-tuned on our datasets). At testing time, we use the person-class detection scores of SSD computed on each generated image \hat{x} . $DS(\hat{x})$ corresponds to the maximum-score box of SSD on \hat{x} and the final DS value is computed by averaging the scores of all the generated images. In other words, DS measures the confidence of a person detector in the presence of a person in the image. Given the high accuracy of SSD in the challenging Pascal VOC 07 dataset [10], we believe that it can be used as a good measure of how much realistic (person-like) is a generated image.

Finally, in our tables we also include the value of each metrics computed using the *real* images of the test set. Since these values are computed on real data, they can be considered as a sort of an upper-bound to the results a generator can obtain. However, these values are not actual upper bounds in the strict sense: for instance the DS metrics on the real datasets is not 1 because of SSD failures.

6.1. Comparison with previous work

In Tab. 1 we compare our method with [12]. Note that there are no other works to compare with on this task yet. The mask-based metrics are not reported in [12] for the DeepFashion dataset. Concerning the DS metrics, we used the publicly available code and network weights released by the authors of [12] in order to generate new images according to the common testing protocol and ran the SSD detector to get the DS values.

On the Market-1501 dataset our method reports the highest performance with all but the IS metrics. Specifically, our DS values are much higher than those obtained by [12]. Conversely, on the DeepFashion dataset, our approach significantly improves the IS value but returns a slightly lower SSIM value.

6.2. User study

In order to further compare our method with the state-of-the-art approach [12] we implement a user study following

Table 1: Comparison with the state of the art. (*) These values have been computed using the code and the network weights released by Ma et al. [12] in order to generate new images.

Model	Market-1501					DeepFashion		
	<i>SSIM</i>	<i>IS</i>	<i>mask-SSIM</i>	<i>mask-IS</i>	<i>DS</i>	<i>SSIM</i>	<i>IS</i>	<i>DS</i>
Ma et al. [12]	0.253	3.460	0.792	3.435	0.39*	0.762	3.090	0.95*
<i>Ours</i>	0.290	3.185	0.805	3.502	0.72	0.756	3.439	0.96
<i>Real-Data</i>	1.00	3.86	1.00	3.36	0.74	1.000	3.898	0.98

the protocol of Ma et al. [12]. For each dataset, we show 55 real and 55 generated images in a random order to 30 users for one second. Differently from Ma et al. [12], who used Amazon Mechanical Turk (AMT), we used “expert” (voluntary) users: PhD students and Post-docs working in Computer Vision and belonging to two different departments. We believe that expert users, who are familiar with GAN-like images, can more easily distinguish real from fake images, thus confusing our users is potentially a more difficult task for our GAN. The results² in Tab. 2 confirm the significant quality boost of our images with respect to the images produced in [12]. For instance, on the Market-1501 dataset, the *G2R* human “confusion” is one order of magnitude higher than in [12].

Finally, in the Supplementary Material we show some example images, directly comparing with [12]. We also show the results obtained by training different person re-identification systems after augmenting the training set with images generated by our method. These experiments indirectly confirm that the degree of realism and diversity of our images is very significant.

Table 2: User study (%). (*) These results are reported in [12] and refer to a similar study with AMT users.

Model	Market-1501		DeepFashion	
	<i>R2G</i>	<i>G2R</i>	<i>R2G</i>	<i>G2R</i>
Ma et al. [12] (*)	11.2	5.5	9.2	14.9
<i>Ours</i>	22.67	50.24	12.42	24.61

6.3. Ablation study and qualitative analysis

In this section we present an ablation study to clarify the impact of each part of our proposal on the final performance. We first describe the compared methods, obtained by “amputating” important parts of the full-pipeline presented in Sec. 3-4. The discriminator architecture is the same for all the methods.

- *Baseline*: We use the standard U-Net architecture [4] *without* deformable skip connections. The inputs of *G* and *D* and the way pose information is represented

²*R2G* means #Real images rated as generated / #Real images; *G2R* means #Generated images rated as Real / #Generated images.



Figure 4: Qualitative results on the Market-1501 dataset. Columns 1, 2 and 3 represent the input of our model. We plot $P(\cdot)$ as a skeleton for the sake of clarity, but actually no joint-connectivity relation is exploited in our approach. Column 4 corresponds to the ground truth. The last four columns show the output of our approach with respect to different baselines.

(see the definition of tensor H in Sec. 3) is the same as in the full-pipeline. However, in G , x_a , H_a and H_b are concatenated at the input layer. Hence, the encoder of G is composed of only one stream, whose architecture is the same as the two streams described in Sec.5.

- *DSC*: G is implemented as described in Sec. 3, introducing our Deformable Skip Connections (DSC). Both

Table 3: Quantitative ablation study on the Market-1501 and the DeepFashion dataset.

Model	Market-1501					DeepFashion	
	<i>SSIM</i>	<i>IS</i>	<i>mask-SSIM</i>	<i>mask-IS</i>	<i>DS</i>	<i>SSIM</i>	<i>IS</i>
<i>Baseline</i>	0.256	3.188	0.784	3.580	0.595	0.754	3.351
<i>DSC</i>	0.272	3.442	0.796	3.666	0.629	0.754	3.352
<i>PercLoss</i>	0.276	3.342	0.788	3.519	0.603	0.744	3.271
<i>Full</i>	0.290	3.185	0.805	3.502	0.720	0.756	3.439
<i>Real-Data</i>	1.00	3.86	1.00	3.36	0.74	1.000	3.898

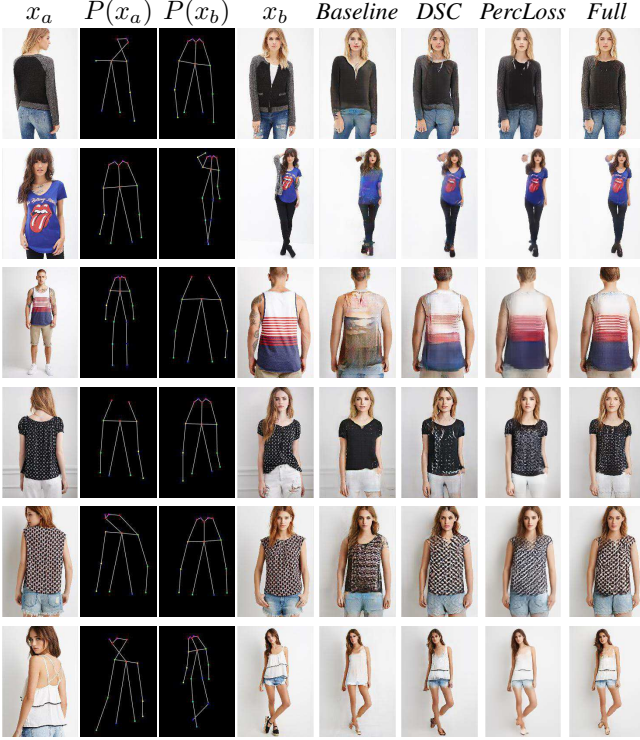


Figure 5: Qualitative results on the DeepFashion dataset with respect to different baselines. Some images have been cropped for visualization purposes.

in DSC and in Baseline, training is performed using an L_1 loss together with the adversarial loss.

- *PercLoss*: This is DSC in which the L_1 loss is replaced with the Perceptual loss proposed in [6]. This loss is computed using the layer *conv2_1* of [15], chosen to have a receptive field the closest possible to $\mathcal{N}(\mathbf{p})$ in Eq. 8, and computing the element-to-element difference in this layer *without* nearest neighbor search.
- *Full*: This is the full-pipeline whose results are reported in Tab. 1, and in which we use the proposed \mathcal{L}_{NN} loss (see Sec. 4).

In Tab. 3 we report a quantitative evaluation on the

Market-1501 and on the DeepFashion dataset with respect to the four different versions of our approach. In most of the cases, there is a progressive improvement from Baseline to DSC to Full. Moreover, Full usually obtains better results than PercLoss. These improvements are particularly evident looking at the DS metrics, which we believe it is a strong evidence that the generated images are realistic. DS values on the DeepFashion dataset are omitted because they are all close to the value ~ 0.96 .

In Fig. 4 and Fig. 5 we show some qualitative results. These figures show the progressive improvement through the four baselines which is quantitatively presented above. In fact, while pose information is usually well generated by all the methods, the texture generated by Baseline often does not correspond to the texture in x_a or is blurred. In same cases, the improvement of Full with respect to Baseline is quite drastic, such as the drawing on the shirt of the girl in the second row of Fig. 5 or the stripes on the clothes of the persons in the third and in the fourth row of Fig. 4. Further examples are shown in the Supplementary Material.

7. Conclusions

In this paper we presented a GAN-based approach for image generation of persons conditioned on the appearance and the pose. We introduced two novelties: deformable skip connections and nearest-neighbour loss. The first is used to solve common problems in U-Net based generators when dealing with deformable objects. The second novelty is used to alleviate a different type of misalignment between the generated image and the ground-truth image.

Our experiments, based on both automatic evaluation metrics and human judgments, show that the proposed method is able to outperform previous work on this task. Despite the proposed method was tested on the specific task of human-generation, only few assumptions are used which refer to the human body and we believe that our proposal can be easily extended to address other deformable-object generation tasks.

Acknowledgements

We want to thank the NVIDIA Corporation for the donation of the GPUs used in this project.

References

- [1] Z. Cao, T. Simon, S. Wei, and Y. Sheikh. Realtime multi-person 2D pose estimation using part affinity fields. In *CVPR*, 2017.
- [2] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results.
- [3] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014.
- [4] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. *CVPR*, 2017.
- [5] M. Jaderberg, K. Simonyan, A. Zisserman, et al. Spatial transformer networks. In *NIPS*, 2015.
- [6] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016.
- [7] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *ICLR*, 2014.
- [8] C. Lassner, G. Pons-Moll, and P. V. Gehler. A generative model of people in clothing. In *ICCV*, 2017.
- [9] J. Liao, Y. Yao, L. Yuan, G. Hua, and S. B. Kang. Visual attribute transfer through deep image analogy. *ACM Trans. Graph.*, 36(4), 2017.
- [10] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. SSD: Single shot multibox detector. In *ECCV*, 2016.
- [11] Z. Liu, P. Luo, S. Qiu, X. Wang, and X. Tang. Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In *CVPR*, 2016.
- [12] L. Ma, X. Jia, Q. Sun, B. Schiele, T. Tuytelaars, and L. Van Gool. Pose guided person image generation. In *NIPS*, 2017.
- [13] D. Pathak, P. Krähenbühl, J. Donahue, T. Darrell, and A. A. Efros. Context encoders: Feature learning by inpainting. *CVPR*, 2016.
- [14] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training gans. In *NIPS*, 2016.
- [15] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv:1409.1556*, 2014.
- [16] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. *arXiv:1512.00567*, 2015.
- [17] D. Ulyanov, A. Vedaldi, and V. S. Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv:1607.08022*, 2016.
- [18] J. Walker, K. Marino, A. Gupta, and M. Hebert. The pose knows: Video forecasting by generating pose futures. In *ICCV*, 2017.
- [19] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Trans. on Image Processing*, 2004.
- [20] B. Zhao, X. Wu, Z. Cheng, H. Liu, and J. Feng. Multi-view image generation from a single-view. *arXiv:1704.04886*, 2017.
- [21] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian. Scalable person re-identification: A benchmark. In *ICCV*, 2015.
- [22] Z. Zheng, L. Zheng, and Y. Yang. Unlabeled samples generated by GAN improve the person re-identification baseline in vitro. In *ICCV*, 2017.