header_navigation

# Multi-scale Attention Guided Pose Transfer

Prasun Roy, Saumik Bhattacharya, Subhankar Ghosh, and Umapada Pal

abstract
**Abstract**— Pose transfer refers to the probabilistic image generation of a person with a previously unseen novel pose from another image of that person having a different pose. Due to potential academic and commercial applications, this problem is extensively studied in recent years. Among the various approaches to the problem, attention guided progressive generation is shown to produce state-of-the-art results in most cases. In this paper, we present an improved network architecture for pose transfer by introducing attention links at every resolution level of the encoder and decoder. By utilizing such dense multi-scale attention guided approach, we are able to achieve significant improvement over the existing methods both visually and analytically. We conclude our findings with extensive qualitative and quantitative comparisons against several existing methods on the DeepFashion dataset.

**Index Terms**—Pose transfer, Attention, GAN, DeepFashion.
abstract

◆

## 1 INTRODUCTION

SYNTHESIZING novel views of an object or a scene can have several applications in computer graphics, image reconstruction, photography, multi-view data generation etc. While view synthesis can be useful in many practical scenarios, this is a significantly challenging problem in computer vision due to a number of factors including occlusion, illumination, geometric and perspective distortions. This is particularly difficult for deformable objects because their shapes can also vary in a specific view which leads to a much larger number of possible views to predict from a single observation.

Pose transfer can be assumed as a fine-grained view synthesis problem where we like to estimate a view (image) of a deformable object (person) with a particular state (pose) from a single observation. Therefore, pose transfer requires an exceptionally robust generative algorithm to infer both the visible and occluded body parts in a new unobserved state (target pose) from a single observed state (source pose) while preserving the general appearance of the person including expression, skin tone, dress and background.

In Fig. 1, we demonstrate the general overview of pose transfer. Here, we have an image, $I_A$ of a person with an initial source pose, $P_A$. The task of the algorithm is to generate an image, $I_B$ of the same person corresponding to an unobserved target pose, $P_B$. We refer to $I_A$ and $I_B$ as *condition image* and *generated image* respectively.

Initial solution to this problem [1], [2] introduced a coarse to fine generation by dividing the problem into individual sub-tasks to handle foreground, background and pose separately. The complexity of such pipeline is later simplified to a unified approach by utilizing deformable GANs [3] and variational U-Net [4]. More recently, a more streamlined approach [5] is introduced by leveraging attention mechanism to progressively transfer pose. The key idea is to perform pose transfer on a local manifold at each intermediate step to avoid the difficulties arise due to much
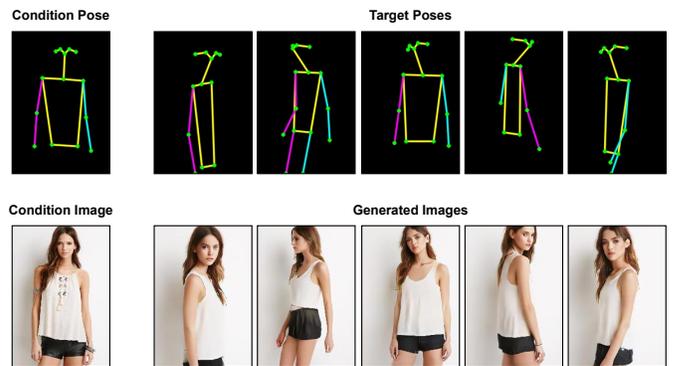


Fig. 1. General overview of pose transfer using the proposed method.

complex structures on the global manifold. In this technique, the condition image and the sparse representations of the poses are initially downsampled to a lower resolution by an encoder. Then, the attention guided progressive pose transfer is performed with the encoded images. Finally, the generated image is upsampled back to higher resolution by a decoder. This method is able to achieve significantly better generation quality than all previous approaches.

This effectiveness of an attention guided approach motivates us to further explore the potential improvements of generation quality. We hypothesized that progressively transferring pose only at a lower resolution followed by upsampling cause a significant loss of finer details present in the condition image. To mitigate this information loss, we propose attention links at every underlying resolution levels of encoder and decoder. In this approach, we introduce a network architecture with such links without requiring any cascaded pose transfer block [5].

In Fig. 2, we show the architecture of the proposed generator. We perform extensive comparative studies among various approaches of pose transfer on the DeepFashion dataset [6]. Our method exhibits significant improvements on the generation performance in both qualitative and analytical benchmarks.

---

- *P. Roy, S. Ghosh and U. Pal are with the Computer Vision and Pattern Recognition Unit, Indian Statistical Institute, Kolkata, India.*
- *S. Bhattacharya is with Indian Institute of Technology, Kharagpur, India.*
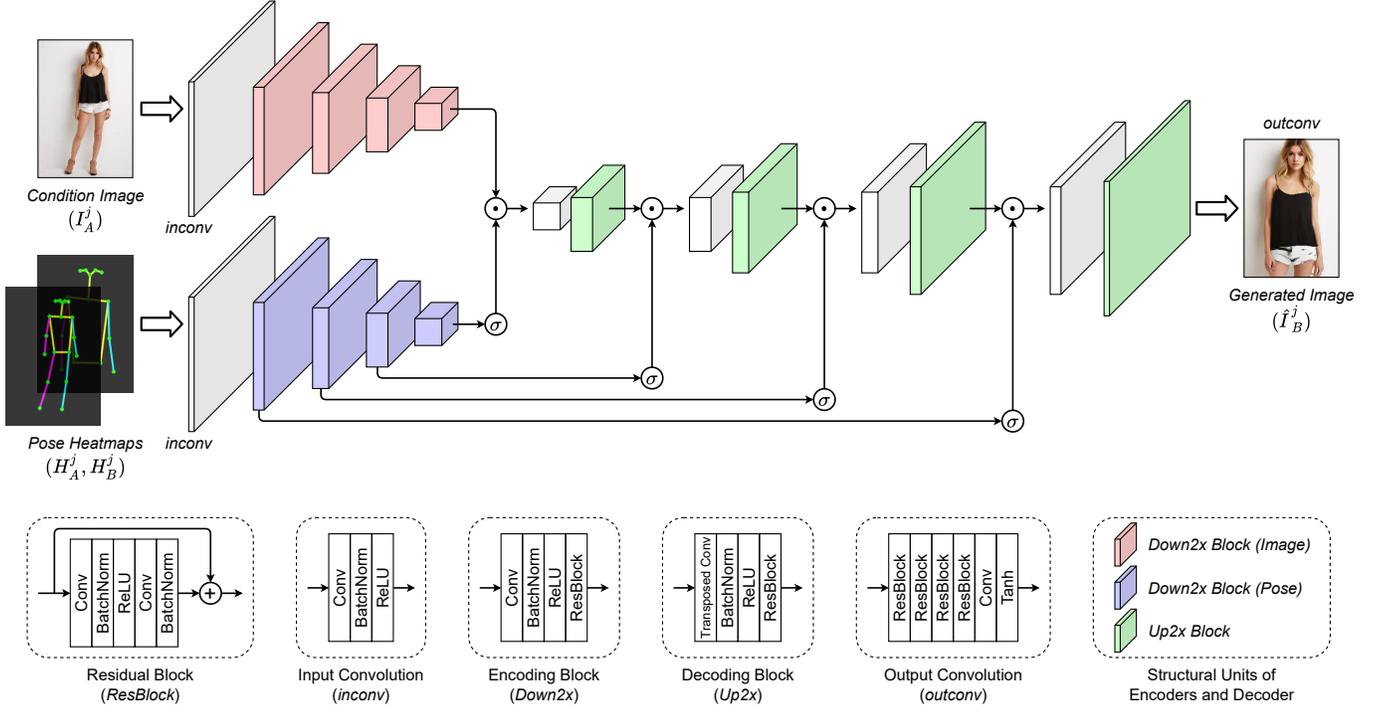- *Project repository: https://github.com/prasunroy/pose-transfer*

Fig. 2. Architecture of the proposed generator. The generator takes the condition image $I_A^j$ along with the channel-wise concatenated pose heatmaps $(H_A^j, H_B^j)$ as inputs and generates an estimate $\hat{I}_B^j$ of the target image $I_B^j$.

## 2 RELATED WORKS

Image synthesis is a fundamental problem in computer vision. Since the introduction of Generative Adversarial Networks (GANs) [7], exceptional improvements in realistic image synthesis have been achieved by several variations [8], [9], [10], [11], [12] of the core GAN architecture in recent years. Initial proposal of GAN [7] introduces unconstrained image generation by enforcing an adversarial learning scheme between a generator and a discriminator. Later, conditional GAN [8] extends the original idea to constrained image generation. Conditional GANs are successfully applied to several computer vision problems including image-to-image translation [13], [14], [15], inpainting [16], super-resolution [17], [18]. Lassner *et al.* [11] propose a generation scheme conditioned over a variational auto-encoder (VAE) to render a 3D model of a person as a fully clothed image. Zhao *et al.* [19] propose a method for generating multi-view images of a person from a single observation in a coarse to fine approach. Balakrishnan *et al.* [20] introduce a pose transfer scheme to segment and generate foreground and background separately. A similar computer vision problem of *Virtual Try On* is also introduced where a human model can be conditionally rendered with a given set of apparel. Han *et al.* [21] propose a thin plate spline transformation guided generation to transfer a specific cloth onto a person in an image. Wang *et al.* [22] propose a geometric matching module in a characteristic preserving generative network.

Initial approach to pose transfer [1] introduces a two stage method in a coarse to fine manner. This technique is further improved in a disentangled person image generation scheme [2] by introducing separate branches for foreground, background and pose in the network architecture. In [23],

authors propose a pose transfer method by estimating a 3D mesh from a single image. In [24], authors propose an unsupervised multi-level generation method with a pose conditioned bidirectional generator. Siarohin *et al.* [3] propose deformable GANs with nearest neighbour loss for pose transfer. Esser *et al.* [4] propose a variational U-Net architecture to generate different poses of an object. More recently, an attention guided progressive generation scheme [5] is introduced by leveraging attention mechanism to progressively transfer pose. In this approach, pose transfer is performed on a local manifold at each intermediate step to overcome several difficulties arise due to much complex structures on the global manifold.

## 3 PROPOSED METHOD

Let us assume that we have a condition image $I_A^j$ of a person with pose $P_A^j$, where $j$ denotes the index of the person in the dataset. Our goal is to generate an image $I_B^j$ of the same person with a new pose $P_B^j$. Similar to the previous approaches [1], [2], [3], [5], we represent a pose as a set of 2D coordinates of body keypoints. We employ a pre-trained Human Pose Estimator (HPE) [25] to estimate 18 keypoints on a human body to represent each pose. The HPE estimates each keypoint as a triplet $(x_i, y_i, v_i)$ where, $(x_i, y_i)$ represents the 2D coordinate and $v_i$ denotes a binary state for visibility of the $i$-th keypoint. More specifically, the value of $v_i$ is 1 for visible keypoints and 0 for occluded keypoints. To represent a pose as a sparse heatmap, we construct a $(h \times w \times 18)$ tensor where, $h$ and $w$ denotes the height and width of the corresponding person image respectively. Each of the 18 channels of the tensor corresponds to one specific

keypoint. For a visible keypoint $(x_k, y_k, 1)$, we put a value of 1 at the location corresponds to the spatial coordinate $(x_k, y_k)$ on the $k$-th channel of the tensor. We put a value of 0 everywhere else in the tensor. We denote the sparse pose heatmaps as $H_A^j$ and $H_B^j$ corresponding to the poses $P_A^j$ and $P_B^j$ respectively.

The proposed generator takes the RGB condition image $I_A^j$ of dimension $(h \times w \times 3)$ and the channel-wise concatenated pose heatmaps $H_A^j$ and $H_B^j$ of dimension $(h \times w \times 36)$ as inputs and generates an RGB image $\hat{I}_B^j$. We employ a PatchGAN discriminator [13] for determining visual correctness of the generated images. The discriminator takes two channel-wise concatenated RGB images, either $(I_A^j, I_B^j)$ or $(I_A^j, \hat{I}_B^j)$, of dimension $(h \times w \times 6)$ as input and estimates a binary class probability map for the input patches.

## 3.1 Generator

In Fig. 2, we illustrate the architecture of the proposed generator. The generator is composed of two downsampling paths (encoders) followed by an upsampling path (decoder) with attention links between feature maps at every underlying resolution levels. We begin by describing the essential components of the generator followed by encoders, decoder and the attention mechanism.

### 3.1.1 Generator Components

**Conv1x1.** A point-wise 2D convolution operation that preserves the size of the input. We perform a single 2D convolution with $1 \times 1$ kernel, stride = 1, padding = 0 and without adding any bias.

**Conv3x3.** A 2D convolution operation that preserves the size of the input. We perform a single 2D convolution with $3 \times 3$ kernel, stride = 1, padding = 1 and without adding any bias.

**Residual Block.** A basic residual block [26] that preserves the size of the input along with the number of channels. A residual block is composed of 5 sequential layers – Conv3x3, Batch Normalization [27], ReLU [28], Conv3x3, Batch Normalization. The input is passed through these layers and the result is added with the original input to produce the final output.

**Down2x Block.** A downsampling block in the encoder for compressing the input size by a factor of 2. We perform a 2D convolution with $4 \times 4$ kernel, stride = 2, padding = 1 and without adding any bias to downsample the input. The resulting feature maps are passed through sequential layers of Batch Normalization, ReLU and a Residual Block to produce the final output.

**Up2x Block.** An upsampling block in the decoder for expanding the input size by a factor of 2. We perform a 2D transposed convolution with $4 \times 4$ kernel, stride = 2, padding = 1 and without adding any bias to upsample the input. The resulting feature maps are passed through sequential layers of Batch Normalization, ReLU and a Residual Block to produce the final output.

### 3.1.2 Encoders

In our approach, we have two parallel downstream branches in the encoder. The *image branch* corresponds to the condition image $I_A^j$ and the *pose branch* corresponds to the concatenated pose heatmaps $(H_A^j, H_B^j)$. Initially, we perform a 2D convolution operation Conv3x3 to project each input into a feature space of $(h \times w \times N_f)$ where, $N_f$ denotes the initial number of feature maps. The convolution is followed by Batch Normalization and ReLU activation to produce the initial input feature maps at each branch. The input feature maps are then passed through $N$ subsequent Down2x Blocks at each branch. At each Down2x Block, we downsampled the input size by a factor of 2 while expanding the number of feature maps by a factor of 2. Therefore, after $N$ consecutive downsampling blocks, we ended up with a feature space of $(\frac{h}{2^N} \times \frac{w}{2^N} \times N_f * 2^N)$.

### 3.1.3 Decoder

In our approach, we have a single upstream branch in the decoder for generating the image $\hat{I}_B^j$ with the target pose. Starting with a feature space of dimension $(\frac{h}{2^N} \times \frac{w}{2^N} \times N_f * 2^N)$, we pass the feature maps through $N$ consecutive Up2x Blocks. At each Up2x Block, we upsampled the input size by a factor of 2 while compressing the number of feature maps by a factor of 2. Therefore, after $N$ subsequent upsampling blocks, we ended up with a feature space of $(h \times w \times N_f)$. Finally, the feature maps are passed through 4 consecutive Residual Blocks followed by a point-wise 2D convolution operation Conv1x1 to project the feature space into an output of dimension $(h \times w \times 3)$. We apply the hyperbolic tangent activation function $tanh$ on the output tensor to get the normalized generated image $\hat{I}_B^j$.

### 3.1.4 Attention Mechanism

We construct attention links between downstream and upstream branches at each resolution level. For resolution level $k$, we compute the attention mask $M_k$ by applying an element-wise $sigmoid$ activation function $\sigma$ on the encoded feature maps of the *pose branch* $H_k^{\mathcal{E}}$. The image feature maps $I_k$ at resolution level $k$ are updated by performing an element-wise product with the attention mask $M_k$. The updated image feature maps $I_k$ act as the input to the upsampling block of the decoder at resolution level $k$. We repeat these operations sequentially up to the highest resolution which results $N$ such attention links. Mathematically, at the lowest resolution level where, $k = N$,

$$I_{N-1}^{\mathcal{D}} = \mathcal{D}_N^{Up2x}(I_N^{\mathcal{E}} \odot \sigma(H_N^{\mathcal{E}}))$$

and for each subsequent higher resolution level where, $k = \{1, ..., N-1\}$,

$$I_{k-1}^{\mathcal{D}} = \mathcal{D}_k^{Up2x}(I_k^{\mathcal{D}} \odot \sigma(H_k^{\mathcal{E}}))$$

Here, we denote downstream encoding as $\mathcal{E}$ and upstream decoding as $\mathcal{D}$.

## 3.2 Discriminator

We use a Markovian PatchGAN discriminator [13] to estimate high-frequency correctness in the generated images. The discriminator complements to the L1 loss which is restricted to low-frequency details only. Such a discriminator
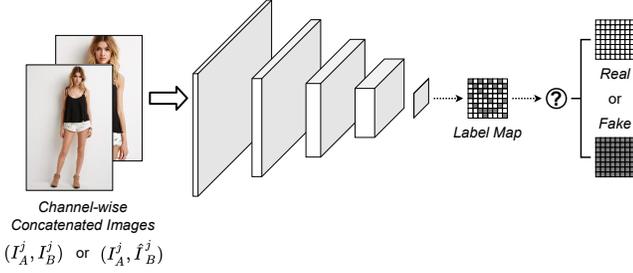
Fig. 3. Architecture of the PatchGAN discriminator. The discriminator takes two channel-wise concatenated images, either $(I_A^j, I_B^j)$ or $(I_A^j, \hat{I}_B^j)$, as input and estimates a label map where, each label corresponds to the binary class probability of an input patch.

operates on $S \times S$ image patches by classifying each patch as real or fake. As explained by the authors [13], PatchGAN functions as a style/texture loss by modeling the image as a Markov random field, assuming independence between pixels separated by more than a patch diameter.

In our approach, we enforce adversarial discrimination on the image transition rather than the image itself. We do this by depth-wise concatenating the condition image $I_A^j$ either with the target image $I_B^j$ or with the generated image $\hat{I}_B^j$ where, $(I_A^j, I_B^j)$ is labeled as real and $(I_A^j, \hat{I}_B^j)$ is labeled as fake. We adopt an identical network architecture as [13] that effectively operates on a $70 \times 70$ receptive field (patch) of the input. In Fig. 3, we illustrate the architecture of the PatchGAN discriminator.

## 3.3 Training

We train our network in an adversarial scheme where, we have two competing objective functions for the generator $G$ and the discriminator $D$. During training, each objective tries to minimize the penalty incurred by itself while trying to maximize the penalty for the other one.

### 3.3.1 Generator Objective

To ensure the low-frequency correctness in the generated images, we compute $L_1$ loss between the target image $I_B^j$ and the generated image $\hat{I}_B^j$. Mathematically,

$$\mathcal{L}_1^G = \|\hat{I}_B^j - I_B^j\|_1$$

To estimate the high-frequency correctness in the generated images, we compute Binary Cross-Entropy (BCE) loss using the PatchGAN discriminator. Mathematically,

$$\mathcal{L}_{GAN}^G = \mathcal{L}_{BCE}(D(I_A^j, \hat{I}_B^j), 1)$$

To improve the visual fidelity of the generated images, we also include a perceptual loss [10] in our generator objective. Mathematically,

$$\mathcal{L}_{P_\rho}^G = \frac{1}{h_\rho w_\rho c_\rho} \sum_{x=1}^{h_\rho} \sum_{y=1}^{w_\rho} \sum_{z=1}^{c_\rho} \|\phi_\rho(\hat{I}_B^j) - \phi_\rho(I_B^j)\|_1$$

where, $\mathcal{L}_{P_\rho}^G$ denotes the perceptual loss computed from the output of the $\rho^{\text{th}}$ layer of a pre-trained VGG19 model [29],

$\phi_\rho$ denotes the output of the $\rho^{\text{th}}$ layer having a feature space of dimension $(h_\rho \times w_\rho \times c_\rho)$. In our approach, we compute the perceptual loss at two different layers ($4^{\text{th}}$ and $9^{\text{th}}$) of a VGG19 model pre-trained on the ImageNet dataset [30] to capture both coarse and fine details in the generated images.

The complete generator objective is calculated as a weighted linear combination of $L_1$ loss, GAN loss and perceptual loss. Mathematically,

$$\mathcal{L}^G = \arg \min_G \max_D \quad \lambda_1 \mathcal{L}_1^G + \lambda_2 \mathcal{L}_{GAN}^G + \lambda_3 (\mathcal{L}_{P_4}^G + \mathcal{L}_{P_9}^G)$$

where, $\lambda_1$, $\lambda_2$ and $\lambda_3$ denote the weights for the corresponding loss functions.

### 3.3.2 Discriminator Objective

Our discriminator objective has a single GAN loss component which is calculated as the average BCE loss over a real image transition $(I_A^j, I_B^j)$ and a fake image transition $(I_A^j, \hat{I}_B^j)$. Mathematically,

$$\mathcal{L}_{GAN}^D = \frac{1}{2} \left[ \mathcal{L}_{BCE}(D(I_A^j, I_B^j), 1) + \mathcal{L}_{BCE}(D(I_A^j, \hat{I}_B^j), 0) \right]$$

where, we assume real label as 1 and fake label as 0. The complete discriminator objective is given by,

$$\mathcal{L}^D = \arg \min_D \max_G \quad \mathcal{L}_{GAN}^D$$

## 3.4 Implementation

In our implementation, the generator is constructed with 4 downsampling blocks in both the encoders and consequently 4 upsampling blocks in the decoder ($N = 4$). The initial number of feature maps is set to 64 ($N_f = 64$). In the generator objective, we set the weights as, $\lambda_1 = 5$, $\lambda_2 = 1$ and $\lambda_3 = 5$. We initialize the parameters of both generator and discriminator before training by sampling from a normal distribution of 0 mean and 0.02 standard deviation. We optimize both generator and discriminator using the stochastic Adam optimizer [31] with learning rate $\eta = 1e^{-3}$, $\beta_1 = 0.5$, $\beta_2 = 0.999$, $\epsilon = 1e^{-8}$ and weight decay = 0. We train the network on a single NVIDIA TITAN X GPU for 270K iterations with a batch size of 8 and serialize the network weights after every 500 iterations. During inference, we select the checkpoint with best evaluation metrics among 70 most recent checkpoints.

## 4 EXPERIMENTS

To evaluate the performance of the proposed architecture, we carry out extensive qualitative and quantitative comparisons against a number of previous major pose transfer methods [1], [3], [4], [5] on a fixed dataset [6]. In most of the evaluation metrics, our method outperforms the previous methods. We also perform user studies for a subjective visual quality assessment.

**Dataset:** In our experiments, we use DeepFashion [6] *In-shop Clothes Retrieval* dataset for performance evaluation and comparison with previous methods. The dataset contains $176 \times 256$ person images centered on $256 \times 256$

square grids. It features a wide variation of outfit and pose which makes this dataset a relevant choice for evaluating and comparing pose transfer algorithms. For direct and fair comparison, we adopt the exact train-test split provided by Zhu *et al.* [5] where, 101,966 image pairs are randomly selected for training and 8,570 image pairs for testing. Also for a better generalization, the image pairs are selected such that the person identities of the training set do not overlap with those of the testing set.

**Metrics:** At present, a quantifiable generalized metric for visual quality assessment of images is an open problem in computer vision. However, in the previous pose transfer algorithms [1], [3], [4], [5], authors estimate a few widely used evaluation metrics for quantifying visual quality. This includes Structural Similarity Index (SSIM) [32], Inception Score (IS) [33], Detection Score (DS) [34] and PCKh [35]. SSIM measures the perceived quality of generated images by comparing with respective real images and considering image degradation as perceived change in structural information. IS uses Inception architecture [36] as an image classifier to estimate the KL divergence [37] between the label distribution and the marginal distribution for a large set of images. DS uses an object detector to estimate the target class recognition confidence of the object detection model as a measure of the perceptual quality. PCKh aims to quantify the shape consistency between generated and real person images by estimating the percentage of correctly aligned keypoints. We also evaluate the Learned Perceptual Image Patch Similarity (LPIPS) [38] metric which is a more modern standard for perceptual image quality assessment.

### 4.1 Qualitative and Quantitative Comparison

In Fig. 4, we show a qualitative comparison among previously proposed major pose transfer algorithms [1], [3], [4], [5] and our method. For a direct and fair comparison we use the same image pairs $(I_A^j, I_B^j)$ as [5] and expanded their visual analysis by incorporating results of our method. It can be observed that images generated by our method are able to retain better skin color, hair style, facial hair and limb structure. Our method also preserves dress texture better than other methods which is more prominent in 2nd and 4th row in Fig. 4. From a generic visual inspection, generated images by our method apparently look more realistic compared to other methods.

The apparent visual superiority of our method is further reflected in quantitative evaluation of multiple perceptual metrics. In Table 1, we show an analytical comparison among different pose transfer methods by evaluating SSIM, IS, DS, PCKh and LPIPS scores. To ensure cycle consistency, we generate images in both directions – $\hat{I}_B^j$ from the image pairs $(I_A^j, I_B^j)$ followed by $\hat{I}_A^j$ from the reverse image pairs $(I_B^j, I_A^j)$. We evaluate each metric on both $\hat{I}_B^j$ and $\hat{I}_A^j$ followed by estimating their mean as the final score. For SSIM and IS, we obtain slightly lower scores than the best results. For DS, PCKh and LPIPS, our method is able to achieve best scores. We improve the PCKh score over PATN [5] by 2% indicating a superior shape consistency with better alignment among keypoints. Similar to [39], we evaluate LPIPS score for both VGG16 [29] and SqueezeNet [40]. In
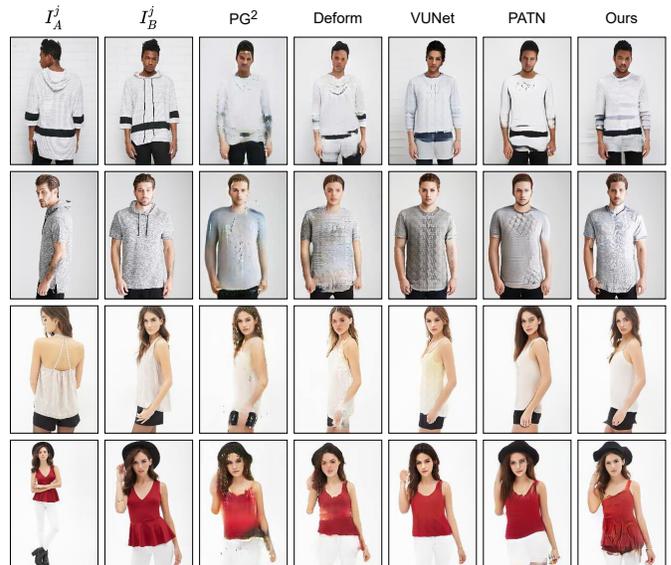


Fig. 4. Qualitative comparison among different pose transfer methods. $I_A^j$ denotes the condition image, $I_B^j$ denotes the target image and subsequent columns show the generated images by PG$^2$ [1], Deformable GANs [3], VUNet [4], PATN [5] and our method.

TABLE 1
Quantitative comparison among different pose transfer methods.

| Method | SSIM | IS | DS | PCKh | LPIPS (VGG) | LPIPS (SqzNet) |
|---|---|---|---|---|---|---|
| PG$^2$ [1] | 0.773 | 3.163 | 0.951 | 0.89 | 0.523 | 0.416 |
| Deform [3] | 0.760 | 3.362 | 0.967 | 0.94 | - | - |
| VUNet [4] | 0.763 | **3.440** | 0.972 | 0.93 | - | - |
| PATN [5] | **0.773** | 3.209 | 0.976 | 0.96 | 0.299 | 0.170 |
| Ours | 0.769 | 3.379 | **0.976** | **0.98** | **0.200** | **0.111** |
| Real Data | 1.000 | 3.864 | 0.974 | 1.00 | 0.000 | 0.000 |

each case, our method achieves significantly better score indicating improved perceptual quality over other methods.

### 4.2 User Study

Although the evaluation metrics discussed so far are widely used for quantifying visual quality, a sufficiently large number of outliers show the limitations of these metrics to be able to generalize the subjective nature of such assessment. Therefore, human perception is perhaps the most reliable way to assess visual quality. For this reason, we perform an opinion based user study to analyze the visual quality of the generated images as perceived by humans. We include two tracks in this study – a constrained test followed by an unconstrained test. In the constrained test, users need to perform discrimination between real and fake images within a fixed amount of time. In the unconstrained test, users need to perform similar discrimination but without any time limit. We follow a similar protocol as [1], [3], [5] for the constrained test except the allowed exposure time for each image. We increase the exposure time from 1 second to 5 second which allows users a better observation before making a decision. This also puts our method in a more challenging position compared to [1], [3], [5].

For our user study, we select 130 real and 130 generated images. 10 images from each set are used as the fixed prac-

TABLE 2
Evaluation scores of user study for different pose transfer methods.

| Method | Exposure Time (second) | R2G (%) | G2R (%) | Accuracy (%) |
|---|---|---|---|---|
| PG$^2$ [1] | 1.0 | 9.20 | 14.90 | 87.95 |
| Deform [3] | 1.0 | 12.42 | 24.61 | 81.49 |
| PATN [5] | 1.0 | 19.14 | 31.78 | 74.54 |
| Ours | 5.0 | 25.90 | **54.26** | **59.92** |
| Ours | $\infty$ | **30.37** | 46.30 | 61.67 |

tice samples. A user is shown 20 images during a test where, 10 images are randomly drawn from each set of the remaining 120 images. For every anonymous user submission, the fraction of real images identified as generated (**R2G**) and the fraction of generated images identified as real (**G2R**) are recorded. The final score is calculated as the mean of global aggregation of all submissions. In Table 2, we show the evaluation scores of our user study conducted with 61 individuals along with the scores reported in previous works [1], [3], [5]. Even with higher exposure time for better visual observation, our method exhibits significantly higher R2G and G2R scores and consequently, much lower recognition accuracy by the users. These results further imply that the images generated by our method are visually more realistic than other methods leading to higher confusion among the users during discrimination.

### 4.3 Ablation Study

The core design characteristic of the proposed network architecture is focused around the attention links at different underlying resolution levels of the generator as shown in Fig. 2. To study the efficacy of such architecture, we perform an ablation study with 4 variants of the generator. In the first variant, denoted as **A0**, we remove every attention link from the generator and thereby making it a generic encoder-decoder network. In the second variant, denoted as **A1-LR**, we keep only one attention link at the lowest resolution level. For the third variant, denoted as **A1-HR**, we keep only one attention link at the highest resolution level. For the fourth variant, denoted as **FULL**, we retain every attention link present in the generator as shown in Fig. 2. We train each network on the same training data for 125K iterations keeping the discriminator, training mechanism and all other implementation conditions same as discussed in Sec. 3.

In Fig. 5, we show a qualitative comparison among the generated images by different model variants. Model A0 generates blurry and visually inconsistent images. Model A1-LR generates visually and structurally consistent images but lacks realistic detail in the face and limbs. Model A1-HR performs very poorly resulting in blurry and incomplete images. Model FULL performs significantly better than other variants producing realistic images with finer detail while preserving visual and structural consistency.

The subjective visual analysis is further supported by a quantitative comparison among the model variants as shown in Table 3 where, we evaluate SSIM, IS, DS, PCKh and LPIPS scores for each model. For most of the metrics, model FULL achieves best evaluation scores as expected. Interestingly, model A1-LR achieves marginally better IS and DS scores than model FULL even with visually inferior



Fig. 5. Ablation study – Qualitative comparison among model variants. $I_A^j$ denotes the condition image, $I_B^j$ denotes the target image and each subsequent column shows the generated images by respective model variant.

TABLE 3
Ablation study – Quantitative comparison among model variants.

| Model | SSIM | IS | DS | PCKh | LPIPS (VGG) | LPIPS (SqzNet) |
|---|---|---|---|---|---|---|
| A0 | 0.760 | 3.055 | 0.969 | 0.97 | 0.221 | 0.124 |
| A1-LR | 0.758 | **3.211** | **0.976** | 0.97 | 0.210 | 0.117 |
| A1-HR | 0.755 | 2.859 | 0.965 | 0.95 | 0.229 | 0.134 |
| FULL | **0.764** | 3.171 | 0.975 | **0.97** | **0.204** | **0.113** |
| Real Data | 1.000 | 3.864 | 0.974 | 1.00 | 0.000 | 0.000 |

generation results as evident from Fig. 5. This anomaly further indicates the limitations of present perceptual metrics to be able to properly generalize human perception. For this reason, the outcome of such analysis is based on multiple metrics rather than relying on a single one.

From the visual and analytical ablation study, we conclude that the dense attention links between every underlying resolution levels of encoder and decoder in the generator network result in significant improvements of visual quality of the generated images. Consequently, this leads to the architecture of the proposed generator as shown in Fig. 2.

### 4.4 Failure Cases

Similar to the previous keypoint based pose transfer methods [1], [2], [3], [5], performance of the proposed method is directly dependent on the success of an HPE such as [25] in estimating the keypoints correctly. An erroneously estimated pose directly impacts the generation performance leading to inconsistency in the generated images. The supervised learning approach during optimization is also associated with the usual problem of data imbalance. In some occasions, the visual quality of the generated image is inconsistent for a distinctly uncommon pose or outfit with

Fig. 6. Failure cases. For each example, $I_A^j$ denotes the condition image, $I_B^j$ denotes the target image and $\hat{I}_B^j$ denotes the generated image.

very few image samples in the training dataset. In Fig. 6, we show some examples where, the proposed method fails to maintain consistency in the generated images.

## 4.5 Extended Applications

To show the efficacy of the proposed network architecture across multiple application domains, we extend our experiments to address a few different problems. An important point to note that we do not modify the network architecture in these experiments. While domain specific modifications of the architecture may further improve generation quality, those analyses are beyond the scope of this paper. However, we show that even without any modification, the proposed network works remarkably well in many different practical scenarios.

### 4.5.1 Conditional Semantic Reconstruction

In novel view synthesis, one typical solution is image reconstruction from the crude semantic map. Unlike an unconditional synthesis, here we use a condition image as reference to transfer various image attributes to the generated image. In our experiments, we use semantic maps provided with the DeepFashion [6] *In-shop Clothes Retrieval* dataset. Here, a semantic map is represented as an aggregation of 16



Fig. 7. Qualitative results of conditional semantic reconstruction using the proposed method.

different parts where, each part is annotated with a unique color. Similar to the keypoint based approach, as discussed in Sec. 3, we represent a semantic map as a 16 channel sparse heatmap where, each channel corresponds to the binary mask of one specific part. We use 5,161 image pairs for training and 795 image pairs for testing. The network is trained for 65K iterations keeping all other implementation conditions same as discussed in Sec. 3. Our method is able to successfully reconstruct the semantic maps by transferring visual attributes from the reference image. In Fig. 7, we show a few qualitative results of conditional semantic reconstruction using the proposed method.

### 4.5.2 Virtual Try-On

Although, the proposed network is not originally intended for virtual try-on applications, the ability of conditional semantic reconstruction can be further extended to address such problems to a significant extent. In this application, the main objective is to replace selected parts of the attire of a target person with that of a reference person. We achieve this in two sequential steps. Initially, we perform a crude reconstruction of the semantic map of the target person by conditioning on the reference person image. The crude reconstructed image is then refined by bitwise operations involving the target person image and the binary mask with selected parts of the attire. Mathematically,

$$I_{fine} = [M_{parts} \odot I_{crude}] \oplus [(1 - M_{parts}) \odot I_{target}]$$

where, $I_{fine}$ denotes the refined target person image with replaced attire, $I_{crude}$ denotes the reconstructed target person image, $I_{target}$ denotes the target person image with original attire and $M_{parts}$ denotes a binary mask of selected parts of the attire. In Fig. 8, we show a few qualitative results of virtual try-on using this approach. Our method is able to successfully replace local parts of the target person's attire.

### 4.5.3 Skeleton Guided Style Transfer

As a part of our earlier work on a scene text editing framework, STEFANN [41], we propose two independent network architectures for preserving structural and color consistencies in a generative character-to-character transformation. In this cross-domain application, we show that the performance of such a two-stage style transfer technique can be significantly improved by leveraging the proposed end-to-end architecture with skeleton supervision. The structural information of the source character is provided as a single channel binary image tensor of its skeleton. However, in this case, the geometric structure of the target character is not known before inference. For this reason, we use the target character skeleton from a fixed font as a crude structural approximation of the target character. In our experiments, we train the network with a much smaller subset of the dataset of STEFANN. We randomly select 200 image pairs of uppercase characters from both training and testing sets. This leads to 203000 image pairs from 1015 fonts for training and 60000 image pairs from 300 fonts for testing, where a random color is applied for each font. The skeleton is estimated from the binarized character image by applying Gaussian blur with a $3 \times 3$ kernel followed by a parallel
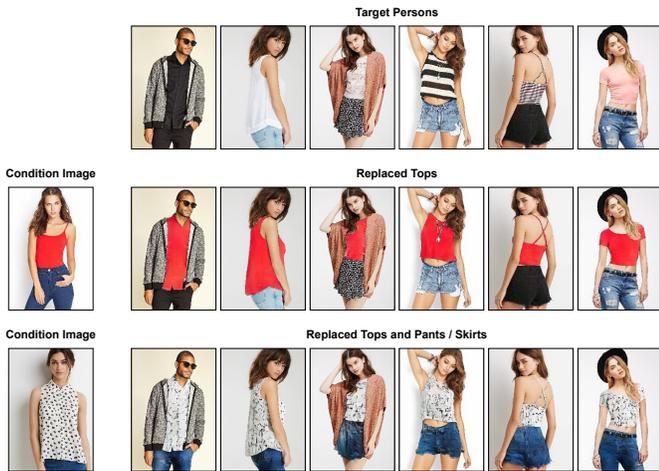
Fig. 8. Qualitative results of virtual try-on using the proposed method.

thinning algorithm [42]. We train the network for 65K iterations keeping all other implementation conditions same as discussed in Sec. 3.
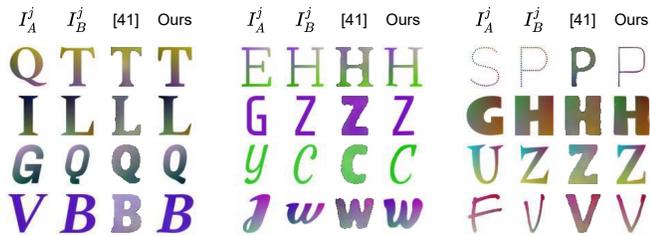


Fig. 9. Qualitative comparison of different methods of style transfer. $I_A^j$ denotes the condition image, $I_B^j$ denotes the target image and subsequent columns show the generated images by STEFANN [41] and our method using skeleton supervision.

TABLE 4
Quantitative comparison of different methods of style transfer.

| Method | SSIM | PSNR (dB) | LPIPS (VGG) | LPIPS (SqzNet) |
|---|---|---|---|---|
| STEFANN [41] | 0.450 | 13.347 | 0.397 | 0.273 |
| Ours | **0.638** | **16.876** | **0.208** | **0.089** |
| Real Data | 1.000 | $\infty$ | 0.000 | 0.000 |

In Fig. 9, we show a qualitative comparison of different methods of style transfer. Even with over 70% reduction in the training data, the end-to-end skeleton guided approach can preserve structural and color consistencies significantly better than the two-stage pipeline. This visual analysis is further reflected in multiple perceptual metrics including SSIM, PSNR and LPIPS as shown in Table 4.

## 5 CONCLUSION

In this paper, we propose a keypoint based generative method for human pose transfer. The key contribution of our work is an improved network architecture by introducing attention links at every underlying resolution level in the encoding and decoding streams of the generator. Our method achieves state-of-the-art results on the DeepFashion dataset in both qualitative and quantitative benchmarks. The visual superiority of the generated images by our method is further supported by an opinion based subjective user study. We also show the efficacy of the proposed network across multiple application domains including conditional semantic reconstruction, virtual try-on and skeleton guided style transfer. These extended applications immediately encourage future works to explore domain specific modifications of the proposed end-to-end architecture.
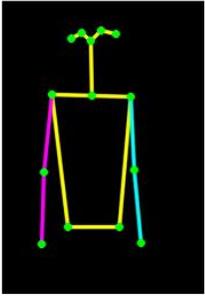
## REFERENCES

[1] L. Ma, X. Jia, Q. Sun, B. Schiele, T. Tuytelaars, and L. Van Gool, "Pose guided person image generation," in *The Conference on Neural Information Processing Systems (NIPS)*, 2017.

[2] L. Ma, Q. Sun, S. Georgoulis, L. Van Gool, B. Schiele, and M. Fritz, "Disentangled person image generation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[3] A. Siarohin, E. Sangineto, S. Lathuilière, and N. Sebe, "Deformable GANs for pose-based human image generation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[4] P. Esser, E. Sutter, and B. Ommer, "A variational u-net for conditional appearance and shape generation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[5] Z. Zhu, T. Huang, B. Shi, M. Yu, B. Wang, and X. Bai, "Progressive pose attention transfer for person image generation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[6] Z. Liu, P. Luo, S. Qiu, X. Wang, and X. Tang, "DeepFashion: Powering robust clothes recognition and retrieval with rich annotations," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[7] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *The Conference on Neural Information Processing Systems (NIPS)*, 2014.

[8] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.

[9] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," in *The International Conference on Learning Representations (ICLR)*, 2016.

[10] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *The European Conference on Computer Vision (ECCV)*, 2016.

[11] C. Lassner, G. Pons-Moll, and P. V. Gehler, "A generative model of people in clothing," in *The IEEE International Conference on Computer Vision (ICCV)*, 2017.

[12] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi, "Photo-realistic single image super-resolution using a generative adversarial network," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[13] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-Image translation with conditional adversarial networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[14] P. Sangkloy, J. Lu, C. Fang, F. Yu, and J. Hays, "Scribbler: Controlling deep image synthesis with sketch and color," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[15] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *The IEEE International Conference on Computer Vision (ICCV)*, 2017.

[16] R. A. Yeh, C. Chen, T. Yian Lim, A. G. Schwing, M. Hasegawa-Johnson, and M. N. Do, "Semantic image inpainting with deep generative models," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[17] C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2015.

[18] J. Kim, J. Kwon Lee, and K. Mu Lee, "Accurate image super-resolution using very deep convolutional networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[19] B. Zhao, X. Wu, Z.-Q. Cheng, H. Liu, Z. Jie, and J. Feng, "Multi-view image generation from a single-view," in *The ACM International Conference on Multimedia (MM)*, 2018.

[20] G. Balakrishnan, A. Zhao, A. V. Dalca, F. Durand, and J. Guttag, "Synthesizing images of humans in unseen poses," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[21] X. Han, Z. Wu, Z. Wu, R. Yu, and L. S. Davis, "VITON: An image-based virtual try-on network," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[22] B. Wang, H. Zheng, X. Liang, Y. Chen, L. Lin, and M. Yang, "Toward characteristic-preserving image-based virtual try-on network," in *The European Conference on Computer Vision (ECCV)*, 2018.

[23] M. Zanfir, A.-I. Popa, A. Zanfir, and C. Sminchisescu, "Human appearance transfer," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[24] A. Pumarola, A. Agudo, A. Sanfeliu, and F. Moreno-Noguer, "Unsupervised person image synthesis in arbitrary poses," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[25] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[26] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[27] S. Ioffe and C. Szegedy, "Batch Normalization: Accelerating deep network training by reducing internal covariate shift," in *The International Conference on Machine Learning (ICML)*, 2015.

[28] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *The International Conference on Machine Learning (ICML)*, 2010.

[29] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *The International Conference on Learning Representations (ICLR)*, 2015.

[30] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.

[31] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *The International Conference on Learning Representations (ICLR)*, 2015.

[32] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: From error visibility to structural similarity," *IEEE Transactions on Image Processing (TIP)*, 2004.

[33] T. Salimans, I. J. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training GANs," in *The Conference on Neural Information Processing Systems (NIPS)*, 2016.

[34] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single shot multibox detector," in *The European Conference on Computer Vision (ECCV)*, 2016.

[35] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele, "2D human pose estimation: New benchmark and state of the art analysis," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.

[36] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.

[37] S. Kullback and R. A. Leibler, "On information and sufficiency," *The Annals of Mathematical Statistics*, 1951.

[38] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[39] O. Gafni and L. Wolf, "Wish you were here: Context-aware human generation," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[40] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, "SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <0.5mb model size," *arXiv preprint arXiv:1602.07360*, 2016.

[41] P. Roy, S. Bhattacharya, S. Ghosh, and U. Pal, "STEFANN: Scene Text Editor using Font Adaptive Neural Network," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[42] T. Zhang and C. Y. Suen, "A fast parallel algorithm for thinning digital patterns," *Communications of the ACM*, 1984.
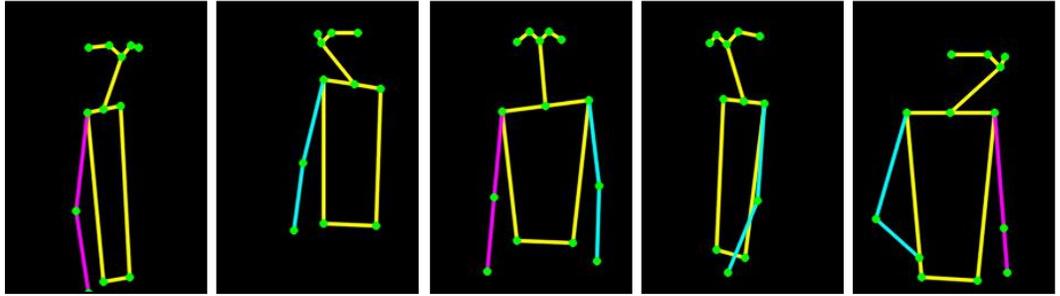
# SUPPLEMENTARY MATERIAL

In this extended text, we present more visual results as supplementary to the original paper. In Fig. S1 and S2, we show visual examples of pose transfer using our method. In Fig. S3 and S4, we show qualitative comparisons among previously proposed major pose transfer algorithms and our method. For a direct and fair comparison, we use the same image pairs $(I_A^j, I_B^j)$ as [5] and expanded their visual analysis by incorporating results of our method. Additionally, our code repository and pre-trained models are publicly available at https://github.com/prasunroy/pose-transfer.

**Condition Pose**

**Target Poses**



**Condition Image**

**Generated Images**



**Condition Pose**

**Target Poses**



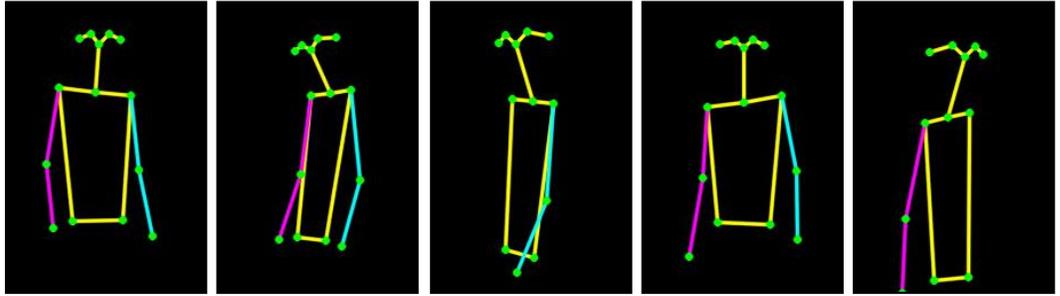**Condition Image**

**Generated Images**



Fig. S1. Qualitative results generated using the proposed method.
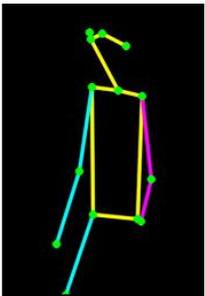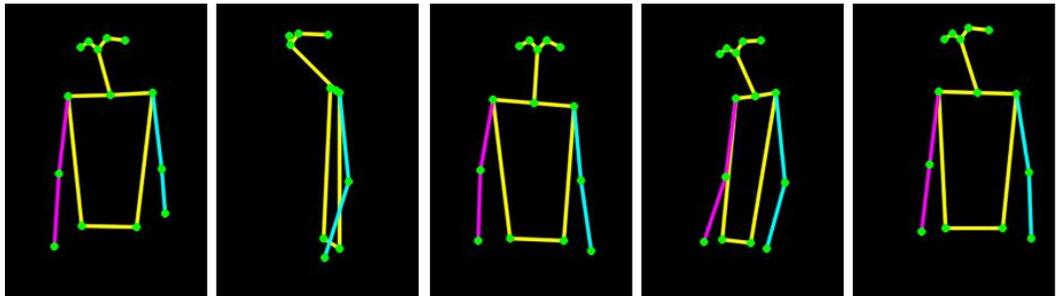
**Condition Pose**

**Target Poses**



**Condition Image**

**Generated Images**



**Condition Pose**

**Target Poses**



**Condition Image**

**Generated Images**



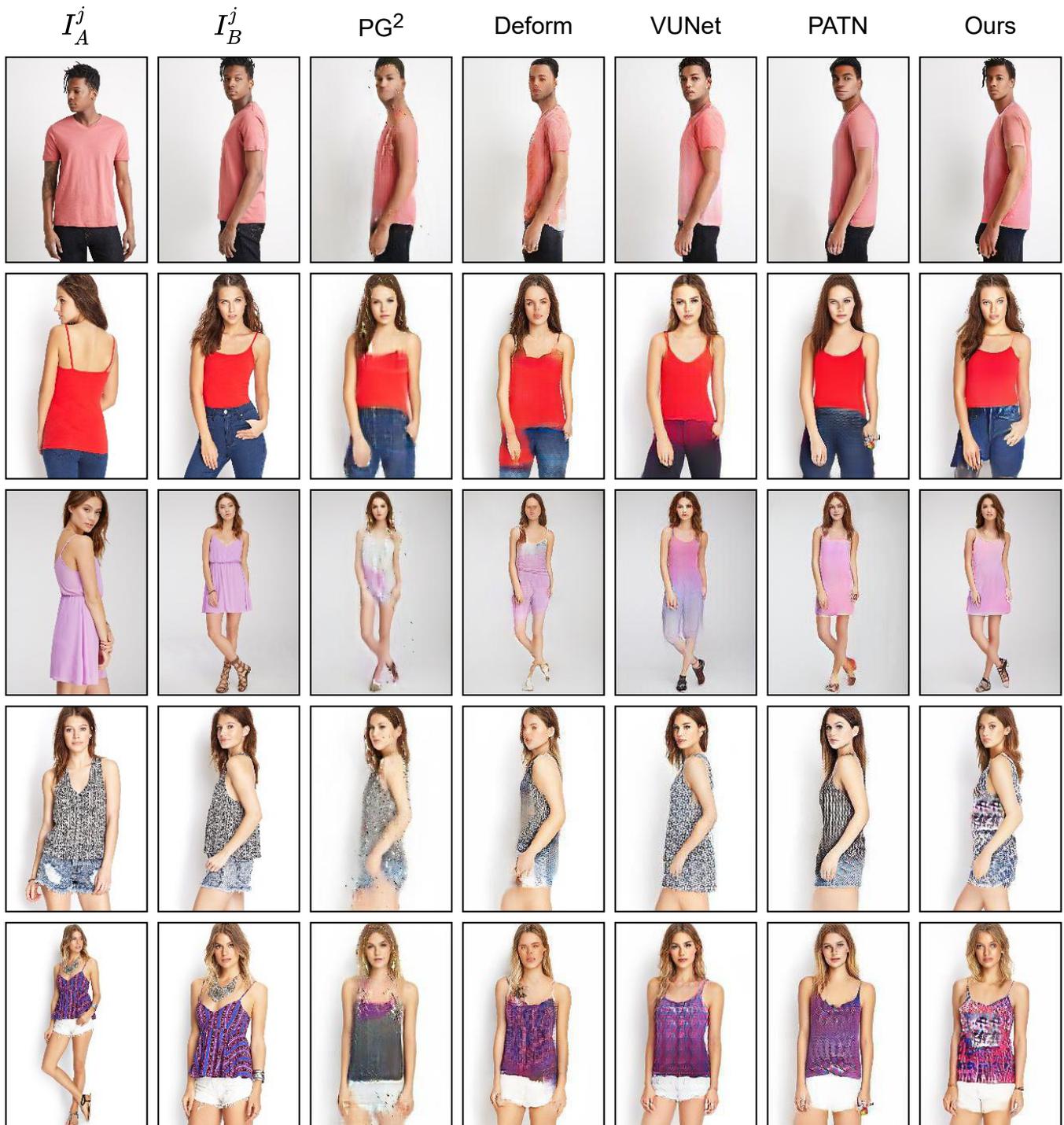Fig. S2. Qualitative results generated using the proposed method.

Fig. S3. Qualitative comparison among different pose transfer methods. $I_A^j$ denotes the condition image, $I_B^j$ denotes the target image and subsequent columns show the generated images by PG$^2$ [1], Deformable GANs [3], VUNet [4], PATN [5] and our method.

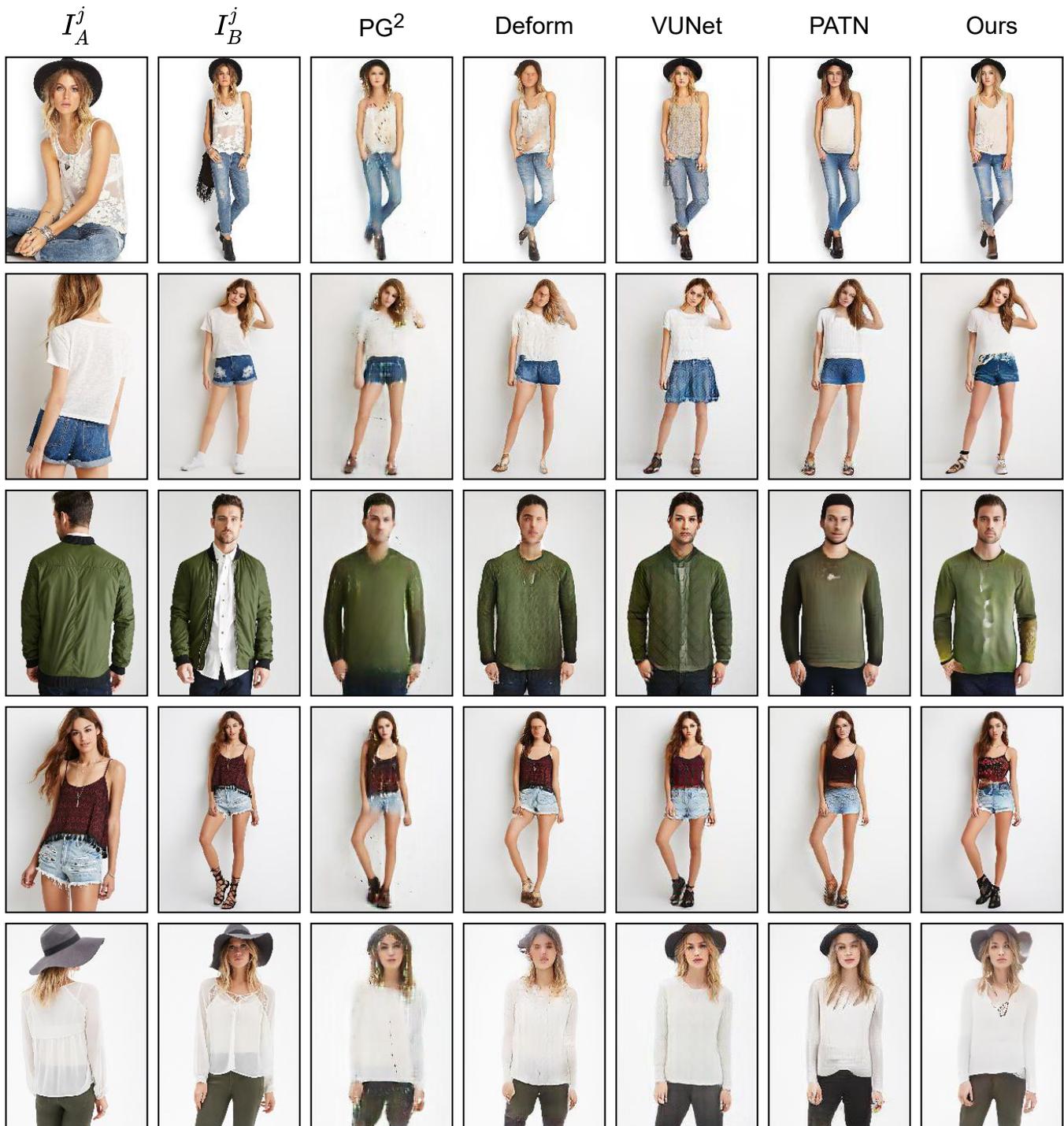| $I_A^j$ | $I_B^j$ | PG$^2$ | Deform | VUNet | PATN | Ours |
|---|---|---|---|---|---|---|



Fig. S4. Qualitative comparison among different pose transfer methods. $I_A^j$ denotes the condition image, $I_B^j$ denotes the target image and subsequent columns show the generated images by PG$^2$ [1], Deformable GANs [3], VUNet [4], PATN [5] and our method.