

HumanGAN: A Generative Model of Human Images

Kripasindhu Sarkar

Lingjie Liu

Vladislav Golyanik

Christian Theobalt

Max Planck Institute for Informatics, SIC



Figure 1: Exemplary samples synthesized by *HumanGAN*. Our pose-guided generative model can sample random appearances conditioned on a given pose (*left*), create persistent appearance and identity across different poses (*middle*), and sample body parts (*right*). All the nine images shown in this figure are synthesized by our generative model.

Abstract

Generative adversarial networks achieve great performance in photorealistic image synthesis in various domains, including human images. However, they usually employ latent vectors that encode the sampled outputs globally. This does not allow convenient control of semantically-relevant individual parts of the image, and cannot draw samples that only differ in partial aspects, such as clothing style. We address these limitations and present a generative model for images of dressed humans offering control over pose, local body part appearance and garment style. This is the first method to solve various aspects of human image generation, such as global appearance sampling, pose transfer, parts and garment transfer, and part sampling jointly in a unified framework. As our model encodes part-based latent appearance vectors in a normalized pose-independent space and warps them to different poses, it preserves body and clothing appearance under varying posture. Experiments show that our flexible and general generative method outperforms task-specific baselines for pose-conditioned image generation, pose transfer and part sampling in terms of realism and output resolution.

1. Introduction

Algorithms to generate images of clothed humans find many applications in such fields as virtual and augmented reality, data generation and augmentation for neural net-

work training. For content creation, it is often desired to have full control over semantic properties of the generated images (e.g., pose, body appearance and garment style). One way of achieving this is to use computer graphics with precise control over the image rendering. However, creating just a single photo-realistic image in this way is challenging and tedious and requires expert knowledge in 3D modeling, animation and rendering algorithms.

Recently, generative adversarial networks (GANs) have made significant progress in generating photo-realistic images, which can also be applied to synthesize imagery of humans [3, 12, 13, 14]. These methods learn a mapping from an easy-to-sample latent space that can be sampled to the image domain and is able to synthesize photo-realistic imagery without the need to resort to complex compute graphics style modeling and light transport simulation. A limitation of these generative models from a content creation perspective is that, in contrast to computer graphics synthesis, they do not easily permit control over semantic attributes of the output imagery. Recently, Tewari *et al.* [41, 42] study StyleGAN [13] to achieve a rig-like control over 3D interpretable face parameters, such as face pose, expression and illumination. However, these methods only work well for faces with limited 3D poses [41], and extending such a method to humans—that relies on differentiable rendering of a 3D morphable model—is not straightforward.

In a setting related to ours, there has been significant progress in conditional image generation with explicit inputs of specific control variables [11, 45]. Conditional

models of this type that use conditioning inputs from a parametric human body model have shown impressive results in applications such as pose and garment transfer [37, 10, 38, 32]. Unfortunately, their underlying translation network does not constitute a full generative model that can be sampled from, as they are designed to produce a single output deterministically.

In this work, we present *HumanGAN*, *i.e.*, a novel generative model for full-body images of clothed humans, which enables control of body pose, as well as independent control and sampling of appearance and clothing style on a body part level¹. Our method combines the advantages of both worlds, *i.e.*, the latent-space-based GANs and controllable translation networks, in the framework of conditional variational autoencoder [17]. We encode the true posterior probability of the latent variable from a space of pose-independent appearance and reconstruct a photo-realistic image of the human with a high-fidelity generator using the encoded latent vector. To disentangle a pose from local part appearance, we propose a novel strategy where we condition the latent vectors on body parts and warp them to a different pose before performing the reconstruction. This permits pose control under persistent appearance and clothing style, and appearance sampling on a localized body part level, without affecting the pose, see Fig. 1. To summarize, our **contributions** are as follows:

- HumanGAN, *i.e.*, a new state-of-the-art generalized model for human image generation that can perform global appearance sampling, pose transfer, parts and garment transfer, as well as part sampling. For the first time, a *single method* can support *all these tasks* (Sec. 3).
- A novel strategy of part-based encoding and part-specific spatial warping in a variational framework that disentangles pose and appearance over the body parts.

In our experiments (Sec. 4), we significantly outperform the state of the art (and other tested methods) for human appearance sampling in realism, diversity and output resolution (512×512). Furthermore, our general model shows commendable results for *pose transfer* that are on par with the state-of-the-art methods developed for this task.

2. Related Work

Deep Generative Models have made remarkable achievements on image generation. As the original GAN model [9] was only able to synthesize low-resolution images, the follow-ups improved it with multiple discriminators [7, 31, 6], self-attention mechanism [48, 3] and progressive training strategy [12]. These methods use a single latent vector \mathbf{z} to resemble the latent factor distribution of training data,

¹project webpage: gvv.mpi-inf.mpg.de/projects/HumanGAN/

which leads to unavoidable entanglements and limited control over image synthesis. StyleGANs [13, 14] approach this problem by mapping \mathbf{z} to an intermediate latent space \mathbf{w} , which is then fed into the generator to change different levels of attributes. Although it enables more control on image synthesis, it does not disentangle different feature factors. Recent works [41, 42] extend StyleGAN to synthesize images of faces with a rig-like control over 3D interpretable face parameters such as face pose, expression and scene illumination. Compared to faces, synthesizing the full human appearance with control of 3D body pose and human appearance is a much more difficult problem due to more severe 3D pose and appearance changes. We propose the first method to this problem allowing photo-realistic image synthesis of a full human body with controls of the 3D pose as well as the appearance of each body part.

Conditional GAN (cGAN) uses conditional information for the generator and discriminator. cGAN is useful for applications such as class conditional image generation [29, 30, 33] and image-to-image translation [11, 45]. Many works [29, 11, 45, 34, 44] require paired data for fully-supervised training. Pix2Pix [11] and Pix2PixHD [45] learn the mapping from input images to output images. Some works [51, 46, 23, 5, 2] tackle a harder problem of learning the mapping between two domains based on unpaired data. cGAN is a deterministic model which produces a single output. Our approach also applies a conditional GAN to map from a warped noise image to output images. However, unlike cGAN methods, we are able to randomly sample noise from a normal distribution for each body part for synthesizing different output images.

Pose Transfer refers to the problem of transferring person appearance from one pose to another [27]. Most approaches formulate it as an image-to-image mapping problem, *i.e.*, given a reference image of target person, mapping the body pose in the format of renderings of a skeleton [4, 39, 35, 18, 53], dense mesh [22, 44, 21, 37, 32, 10], dense labeled pose maps [?] or joint position heatmaps [27, 1, 28] to real images. Though our method is not specifically designed for pose transfer, it can also be applied to this problem with high-quality results, as our generated samples retain identity across different poses. We demonstrate this in Sec. 4.

Variational Autoencoders (VAE) are likelihood-based models which can effectively model stochasticity [17]. Larsen *et al.* [19] first introduced a combined VAE and GAN to achieve higher quality than vanilla VAE. VUNet [8] combined VAE with UNet for pose-guided image generation. Lassner *et al.* [20] present a two-stage VAE framework using a parametric model of human meshes [26] as pose and shape conditioners. Our method builds on conditional VAE-GAN. Unlike the existing methods, it generates images of higher resolution and quality, and offers more control over part sampling and garment transfer.

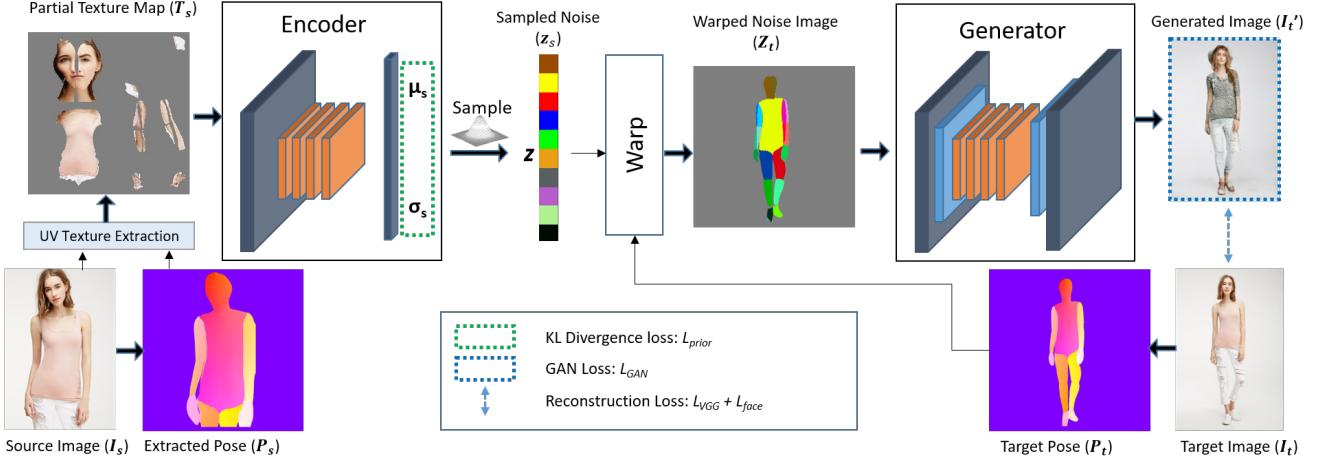


Figure 2: Overview of our method. Given a source image I_s , we extract a UV texture map of the human appearance T_s . The encoder then encodes T_s to part-specific latent vectors z_s . The target pose P_t is used to warp and broadcast the latent vectors to the corresponding parts in the target image to create a noise image Z_t . Finally, the generator converts Z_t to a realistic image. At testing, we generate random samples by controlling the latent vectors $z \in \mathcal{N}(0, 1)$ and the target pose P_t .

3. Method

Our goal is to learn a generative model of human images, which is conditioned on body pose and a low dimensional latent vector encapsulating the appearance of different body parts. We use DensePose [36] to represent the human pose. Our task can be then formulated as a deterministic mapping $G : (\mathbf{P}, \mathbf{z}) \rightarrow \mathbf{I}$. Here, $\mathbf{P} \in \mathbb{R}^{H \times W \times 3}$ is a three-channel DensePose image representing the conditioning pose, $\mathbf{z} \in \mathbb{R}^{M \times N}$ is the latent vector comprised of M human body parts and N part-specific latent vector dimensions, and $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$ is the generated image.

For learning such function, we build our method on variational autoencoders (VAE) [17]. As in any latent vector model, the observed variable of images \mathbf{I} is assumed to be dependent on the unobserved hidden latent variable \mathbf{z} that encodes semantically meaningful information about \mathbf{I} . The goal is then to find their joint distribution $p(\mathbf{I}, \mathbf{z})$ by maximizing the evidence lower bound (ELBO) – by jointly optimizing the encoder $q(\mathbf{z}|\mathbf{I})$ and the decoder $p(\mathbf{I}|\mathbf{z})$. In this work, our *key assumption* is that the latent variable \mathbf{z} depends *only* on the appearance of the subject in the image \mathbf{I} . To enforce that, we first extract a pose-independent human appearance \mathbf{T} from \mathbf{I} . Our encoder $q(\mathbf{z}|\mathbf{T})$ is then conditioned on \mathbf{T} (which is actually a function of \mathbf{I}), to encode to part-specific latent vector \mathbf{z} . Furthermore, the encoded appearance \mathbf{z} is warped by a target pose P_t different from the pose in \mathbf{I} , which is subsequently used by a generator. We next describe our method in detail.

3.1. Our Architecture

In the training stage, we take pairs of images (I_s, I_t) of the same person (but in different poses) as input. Our

method performs in four steps. In the first step, we extract SMPL UV texture map T_s from the input image I_s using the DensePose correspondences. In the second step, we use an encoding function E to map the human appearance T_s of the source image to the parameters of the distribution of the latent vector. In the third step, we sample z from the estimated distribution of the source appearance. Given a target pose P_t , we warp the encoded latent vector z to a noise image Z_t . In the fourth step, we decode the warped Z_t to a realistic image I'_t by a high-fidelity generator network. Our method is summarized in Fig. 2.

Extracting Appearance. We use a UV texture map of the SMPL surface model [26] to represent the subject’s appearance in the input image. The pixels of the input image I_s are transformed into the UV space through a mapping predicted by DensePose RCNN [36]. The pretrained network trained on COCO-DensePose dataset predicts 24 body segments and their part-specific U, V coordinates of SMPL model. For easier mapping, the 24 part-specific UV maps are combined to form a single normalized UV texture map T_s in the format provided in SURREAL dataset [43]. This normalized (partial) texture map provides us with a pose-independent appearance encoding of the subject that is located spatially according to the body parts. The 24 part segments in the texture map also provide us the placeholder for part-based noise sampling, *i.e.*, in our case, the number of body parts $M = 24$.

Encoding Appearance. As with VAE, we assume the distribution $q(\mathbf{z}|T_s)$ of the latent code \mathbf{z} , given the appearance T_s to be Gaussian, $q(\mathbf{z}|T_s) \equiv \mathcal{N}(\mu_s, \sigma_s)$. We use a convolutional neural network $E(\cdot)$ that takes the partial texture T_s as input and predicts the parameters (μ_s, σ_s) of the Gaussian distribution, $\mu_s, \sigma_s \in \mathbb{R}^{M \times N}$. The encoder comprises

a convolutional layer, five residual blocks, an average pooling layer, and finally, a fully connected layer that produces the final output.

Warping Latent Space. In the next step, we sample a latent code from the predicted distribution of the encoded appearance, $\mathbf{z}_s \sim E(\mathbf{T}_s) \equiv \mathcal{N}(\boldsymbol{\mu}_s, \boldsymbol{\sigma}_s)$. Given the noise vector \mathbf{z}_s and a target pose \mathbf{P}_t , we intend to reconstruct a realistic image \mathbf{I}'_t with the appearance encoded in $\mathbf{z}_s \in \mathbb{R}^{M \times N}$ and pose from \mathbf{P}_t . We also want the latent code for a specific body part to have direct influence on the same body part in the generated image. We enforce this by warping and broadcasting the part-based latent code to the corresponding part location in the target image and create a noise image $\mathbf{Z}_t \in \mathbb{R}^{H \times W \times N}$. Here, for each body part k , $\forall_{i,j \in k} \mathbf{Z}_t[i,j] \leftarrow \mathbf{z}[k]$ (see the warping module in Fig. 2). This operation can be easily implemented by differentiable sampling $W(\cdot)$ given the DensePose image \mathbf{P}_t , i.e., $\mathbf{Z}_t = W(\mathbf{z}_s, \mathbf{P}_t)$. The design enables us to perform part-based sampling during the test time. Other straightforward ways of using \mathbf{z}_s , such as sampling noise in the UV texture space or tiling of a single noise vector in the entire spatial dimension, did not give us the required result. See Sec. 4.1 and 4.3 for a detailed analysis.

Decoding to a Photo-Realistic Image. The warped noise image in the target pose \mathbf{Z}_t with the noise vectors correctly aligned with the body parts in the target pose, is used as an input to a generator network $G(\cdot)$. The generator and the warping module act as the conditional decoder of our pipeline. We use the high-fidelity generator from Pix2PixHD [45] that comprises three down-sampling blocks, six residual blocks and three up-sampling blocks.

3.2. Training Details

Our entire training pipeline can be summarized by the following equations:

$$\mathbf{z}_s \sim E(\mathbf{T}_s), \quad \mathbf{Z}_t = W(\mathbf{z}_s, \mathbf{P}_t), \quad \mathbf{I}'_t = G(\mathbf{Z}_t), \quad (1)$$

With the re-parameterization trick [17] for sampling \mathbf{z}_s , the entire pipeline becomes differentiable, allowing direct back-propagation to the parameters of $E(\cdot)$ and $G(\cdot)$. The pipeline is trained with an objective \mathcal{L}_{total} derived from conditional VAE-GAN:

$$\mathcal{L}_{total} = \mathcal{L}_{rec} + \mathcal{L}_{GAN} + \mathcal{L}_{prior}. \quad (2)$$

We describe in the following all three loss terms of (2).

Reconstruction Loss. The reconstruction loss $\mathcal{L}_{rec} = \lambda_{VGG} L_{VGG} + \lambda_{face} L_{face}$ quantifies the dissimilarity between the generated image $\tilde{\mathbf{I}}_t$ and the ground-truth image \mathbf{I}_t . It comprises 1) *Perceptual Loss* L_{VGG} which is the difference between the activations on different layers of the pre-trained VGG network [40] applied on the generated image and the ground-truth image; and 2) *Face Identity Loss*

L_{face} which is the difference between features of the pre-trained SphereFaceNet [24] on the cropped face of the generated image and the ground-truth image.

GAN Loss. The GAN loss $\mathcal{L}_{GAN} = \lambda_D L_D + \lambda_{FM} L_{FM}$ pushes the generator to generate realistic images. We directly use the two-scale discriminator architecture D from Pix2PixHD [45] for implementing the GAN loss. The network D is conditioned on both generated image and warped noise image at different scales. The total GAN loss comprises of multiscale adversarial loss L_D and discriminator feature matching loss L_{FM} . See [45] for more details.

Prior Loss \mathcal{L}_{prior} . To enable sampling at inference time, the encoding $E(\mathbf{T}_s)$ is encouraged to be close to a standard Gaussian distribution, i.e., the prior distribution on the \mathbf{z} vector is assumed to be $\mathcal{N}(0, I)$. Therefore, we employ the prior loss $\mathcal{L}_{prior} = \lambda_{KL} \mathcal{D}_{KL}(E(\mathbf{T}_s) || \mathcal{N}(0, I))$, where $\mathcal{D}_{KL}(p||q)$ is the Kullback-Leibler divergence between the probability distributions $p(x)$ and $q(x)$.

With reparameterization trick on sampling \mathbf{z}_s , we train the system end-to-end and optimize the parameters of the networks E , G and D . The final objective \mathcal{L}_{total} in Eq. (2) is minimized with respect to the generator G and the encoder E , while maximized with respect to the the discriminator D . We use Adam optimiser [16] for our optimization with an initial learning rate of 2×10^{-4} , β_1 as 0.5 and no weight decay. The loss weights are set empirically to $\lambda_{VGG} = 10$, $\lambda_{face} = 5$, $\lambda_D = 1$, $\lambda_{FM} = 10$, $\lambda_{KL} = 0.01$. For speed, we pre-compute DensePose on the images and directly read them as input.

3.3. Inference: Sampling Poses and Body Parts

During testing, we *sample* the appearance vector \mathbf{z} from the prior distribution. We warp \mathbf{z} with the conditioning pose \mathbf{P} and feed the resulting noise image to the trained generator G to get a generated image $\mathbf{I}_{z,P}$:

$$\mathbf{z} \sim \mathcal{N}(0, I), \quad \mathbf{Z}_P = W(\mathbf{z}, \mathbf{P}), \quad \mathbf{I}_{z,P} = G(\mathbf{Z}_P), \quad (3)$$

The appearance \mathbf{z} can also be encoded from an input image by using the encoder on its partial texture map \mathbf{T} , i.e., $\mathbf{z} = \boldsymbol{\mu}$, where $\boldsymbol{\mu}, \boldsymbol{\sigma} = E(\mathbf{T})$. Keeping \mathbf{z} fixed and varying the the pose \mathbf{P} , we can perform *pose transfer* [37, 10, 38, 32], i.e., re-rendering a subject with different pose and viewpoint. We can also perform *parts-based sampling* and *garment transfer* by only varying the vector $\mathbf{z}[k]$ corresponding to the part k . See Sec. 4 for all possible applications of our system.

4. Experimental Results

Dataset. We use the *In-shop Clothes Retrieval Benchmark* of DeepFashion dataset [25] for our main experiments. The dataset comprises of around 52K high-resolution images of fashion models with 13K different clothing items in different poses. Training and testing splits are also provided. To

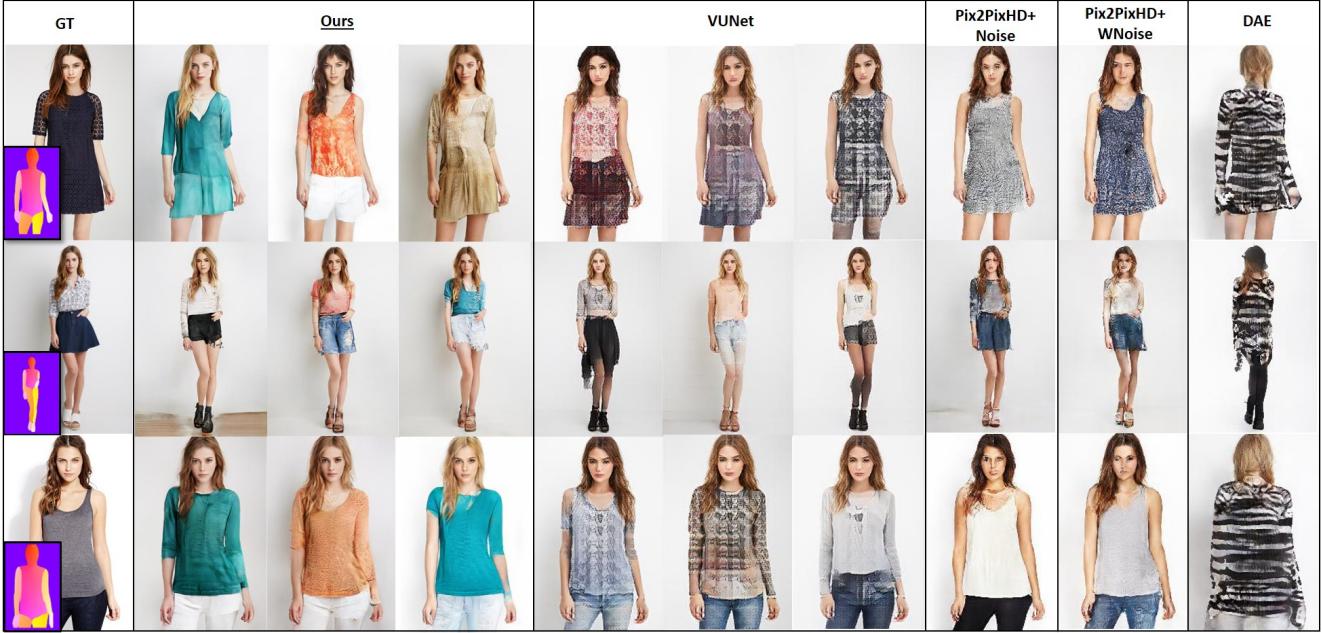


Figure 3: Results of our HumanGAN approach for pose-guided image generation, and its comparison with VUNet [8], Pix2PixHD+Noise, Pix2PixHD+WNoise and Deterministic Auto Encoder (DAE). The conditioning pose is shown in the left corner. Our approach produces samples of higher quality than the baseline methods.

filter non-human images, we discard all the images where we could not compute DensePose, resulting in 38K training images and 3K testing images. We train our system with the resulting training split and use the testing split for conditioning poses. We also show qualitative results of our method with Fashion dataset [47] that has 500 training and 100 test videos, each containing roughly 350 frames.

Experimental Setup. We train our model for the resolution of 512×512 with the training procedure described in Sec. 3.2. The resolution of the partial texture map T is chosen to be 256×256 , number of body parts $M = 24$, and the latent vector dimension $N = 16$. To evaluate our model for the ability to generate diverse and realistic images, preserve the identity of the generated output across different poses, and perform part-based sampling, we use the *same trained model* for *all the experiments* in the following subsections. All the ablation experiments use the same setting in terms of output resolution. When the result of the comparison methods do not have a trained model for 512 resolution (especially for the SOA methods on pose transfer), we resize our image before performing the quantitative evaluation.

4.1. Appearance Sampling

We next evaluate the ability of our system to generate diverse and realistic images of humans given a conditioning pose. Given a *fixed pose*, we *randomly generate samples* from the latent vector $\mathbf{z} \sim \mathcal{N}(0^{MN}, I^{MN})$ and compare our method with the following baseline methods.

	Diversity LPIPS Distance \uparrow	Realism FID \downarrow
VUnet [8]	0.182	50.0
Pix2PixHD+Noise	4.6e-6	109.4
Pix2PixHD+WNoise	0.008	101.9
DAE+WNoise	0.083	187.4
Ours	0.219	24.9
ground truth	0.44	0.0

Table 1: Diversity vs realism. We use LPIPS distance [49] between the randomly generated samples of the same pose to measure Diversity, and FID to measure the realism of the generated samples. \uparrow (\downarrow) means higher (lower) is better.

Pix2PixHD + Noise. Some generators, such as Pix2Pix, Pix2PixHD produces a single output given a conditional input. Randomly drawn noise from a prior distribution can be added to the input of conditional generators to induce stochasticity. We sample noise $\mathbf{z} \in \mathbb{R}^N$ from the standard Gaussian distribution and tile it across the 3 channel condition DensePose image to produce $3 + N$ channel input vector and train Pix2PixHD for a pose-conditioned multimodal human generator. We optimize the conditional generator G and discriminator D with the GAN loss: $\max_D \min_G \mathbb{E}_{\mathbf{P}, \mathbf{T} \sim p(\mathbf{P}, \mathbf{T})} [\log(D(\mathbf{P}, \mathbf{T}))] + \mathbb{E}_{\mathbf{P} \sim p(\mathbf{P}), \mathbf{z} \sim p(\mathbf{z})} [\log(1 - D(\mathbf{P}, G(\mathbf{P}, \mathbf{z})))]$, and the reconstruction loss between $G(\mathbf{P}, \mathbf{z})$ and \mathbf{I} .

Pix2PixHD + Warped Noise. Isola *et al.* [11] and Zhu



Figure 4: Generated images with interpolated appearance encodings. Note that each of the intermediate images is a coherently dressed human image.

	SSIM ↑	LPIPS ↓		SSIM ↑	LPIPS ↓
CBI [10]	0.766	0.178	DPT [32]	0.759	0.206
DSC [38]	0.750	0.214	NHRR [37]	0.768	0.164
VUNet [8]	0.739	0.202	Ours	0.777	0.187

Table 2: Quantitative results for pose transfer against several state-of-the-art methods using Structural Similarity Index (SSIM) [50] and Learned Perceptual Image Patch Similarity (LPIPS) [49]. ↑ (↓) means higher (lower) is better.

et al. [52] observed that a simple extension of a translation network with noise for the purpose multimodal generation often ends up in mode collapse. Therefore, we extend our *Pix2PixHD + Noise* baseline by sampling part specific noise vector $\mathbf{z} \in \mathbb{R}^{M \times N}$ and warping it with the input DensePose to create a noise image \mathbf{Z}_P as in Sec. 3.1. This noise image is used as conditioning input to Pix2PixHD.

Deterministic Auto Encoder with Warped Noise (DAE). This baseline serves as the deterministic version of our method. Because of the lack of constraints on the latent variable (KL divergence with the prior), and direct broadcasting of a single noise vector to multiple pixels in the warping operation, we encountered the problem of exploding gradients. To make the training process tractable, we use a UNet to encode the appearance to the UV texture space, $\mathbf{Z}_s = E(\mathbf{T}_s)$, $\mathbf{Z}_s \in \mathbb{R}^{h \times w \times N}$ (h, w are the texture map dimensions). We warp \mathbf{Z}_s from the texture coordinates to the pixel coordinates by the target pose (instead of a single vector per parts) to complete the pipeline for training. Because of the similarity of this baseline to NHRR [37], we use UNet configuration as in their work.

VUNet [8]. We compare our result to VUNet that performs disentanglement between appearance and structure and can be used for human generation. We use their publically available code and trained model on DeepFashion dataset and report their results here. Note that in contrast to our method, VUNet cannot perform part-based sampling (Sec. 4.3).

Additional Baselines We perform following additional baselines during development and training of our model: a) *NoParts* does not perform part-specific warping, but broadcasts concatenated part-encoded noise-vector to the human silhouette (see Sec. 4.3 for more details on this) b) *+DP-Cond* conditions the DensePose image \mathbf{P}_t in addition to the



Figure 5: Comparison of our reconstruction+transfer results with the state-of-the-art pose transfer methods: CBI [10], NHRR [37], DSC [38], VUNet [8], and DPT. Our approach produces more realistic renderings than the competing methods. Best viewed when zoomed digitally. More results are provided in the appendix.

noise image \mathbf{Z}_t to the generator. c) *+NoisePrior* use samples from the prior $\mathcal{N}(0, 1)$ (along with the encoded distribution) during the training, as recommended by Larsen *et al.* [19]. All the aforementioned *additional baselines* performed poorly, and we present their results in the appendix. We discuss here the main baselines that are most representative of the different methods.

The qualitative results are shown in Fig. 3. We find that *Pix2PixHD+N* produces a single realistic output on the conditioning pose and undergo full mode collapse. The baseline *Pix2PixHD+WN* produces slightly diverse output than *Pix2PixHD+N*, but the variations are still not meaningful. While the deterministic method of *DAE* is a right choice for reconstruction and pose transfer (see results of NHRR [37] in Sec. 4.2), there is no guarantee that the distribution of the encoded latent space will be close to a prior distribution, which makes the test-time sampling difficult. This makes the output far from realistic when the latent vector is sampled from $\mathcal{N}(0, I)$. VUNet produces a diverse output but lacks the quality due to its less powerful generator. In con-



Figure 6: Our results for part sampling. We change the latent noise corresponding to a specific body part.



Figure 7: Results of our method for motion transfer in a walking sequence from fashion dataset [47]. See the accompanying video for the motion results.

trast, we find that our full method produces results that are both diverse and realistic. We also find our latent space of appearance to be smooth and interpolateable. Fig. 4 shows images generated by interpolating the appearance vector \mathbf{z} between two encodings. Note how each of the intermediate images is a coherently dressed human image. More qualitative results are provided in the appendix.

To evaluate our method quantitatively, we randomly select 100 poses from the test set and generate 50 samples for each pose for all the methods. We measure diversity and realism of all the baselines by the following metrics 1) *Pairwise LPIPS distance* – we compute LPIPS distance [49] between all the generated samples for each pose, and take the mean of the distances of all such pairs. More the value of this metric, more diverse is the output. 2) *FID* – we compute the Frechét Inception Distance between the generated samples and the training split of the dataset. FID captures how close is the distribution of the generated samples, from the distribution of the ground truth in the InceptionV3 feature space, and it has been used widely in the community as a metric for quality for GANs [3, 13, 41, 42]. The quantitative results are in Table 1. Our method outperforms other baselines significantly in terms of quality of the image (FID), while maintaining diversity.

User Study. We perform a comprehensive user study to access visual fidelity and characteristics of the results between VUNet and our method for appearance sampling, as those arguably generate visually most realistic human images (see Table 1). To this end, we use 20 random test poses

and choose the results by both the methods which look most realistic, in our opinion. Users are then shown a pair of generated images from the two methods for 12 poses and asked to select the most realistic one among them. Since such comparisons on individual images can be biased, we also ask the participants to decide between image sets (two to four images) for eight poses. The purpose of those is to compensate for the possible image selection biases associated with image pairs. In the total of 36 respondents, our method is preferred over VUNet in **91.06%** of the cases.

4.2. Pose Transfer

In this section, we evaluate how our system preserves the appearance of the outputs across different poses and perform the following pose transfer experiment, *i.e.*, re-rendering of a subject under different poses and viewpoints. We encode the appearance \mathbf{T} of the input image \mathbf{I} by using the mean of the distribution predicted by the encoder, *i.e.*, $\mathbf{z} = \mu$ where $\mu, \sigma = E(\mathbf{T})$. We keep \mathbf{z} *fixed*, and use *different target poses* to generate images for pose transfer.

We compare our results with five state-of-the-art pose transfer methods, namely Coordinate Based Inpainting (CBI) [10], Deformable GAN (DSC) [38], Variational U-Net (VUNet) [8] Dense Pose Transfer (DPT) [32] and Neural Human Re-Rendering (NHRR) [37]. For both qualitative and quantitative evaluation, we use the results of 176 testing pairs that are used in the existing work [37, 10]. The qualitative results are shown in Fig. 5. Note that these methods, except for VUNet, are designed explicitly for the problem of pose transfer, while our method is designed as a generative model which is capable of retaining identity across different pose. We observe that our results show better realism than the other state-of-the-art methods, and perform comparably for the problem of pose transfer. This is confirmed by the comparable LPIPS distance of our model in comparison to the other SOA methods (Table 2). The appearance consistency is also verified by our result on the walking sequences of the fashion dataset [47] (Fig. 7).

User Study. Following the existing work on pose transfer, we perform another user study for evaluating our re-



Figure 8: The results of our method for garment transfer. By combining the appearance encoding of two images based on the body parts, we can perform garment transfer.

	Variation–Part ↑	Variation–Rest ↓
NoParts	0.45	0.44
Ours	0.37	0.11

Table 3: Quantitative evaluation for part sampling using mean pairwise L1 distance of the masked part. ↓(↑) means lower (higher) is better.

sults for this task, and compare with the method of NHRR and CBI. We adapt the user study format and question types from [37] and show the generated results of the methods from 16 source-target pairs. The user preference for *identity preservation* was found to be **ours: 65.62%**, CBI: 21.88%, NHRR: 12.5%. The user preference for *realism* was found to be **ours: 81.25%**, CBI: 6.25%, NHRR: 12.5%.

4.3. Part-Based Sampling

We next evaluate our method for part-based sampling – the ability to produce different plausible renderings of a body part (*e.g.*, head) while keeping the rest of the body same. To this end, we vary the vector $z[k] \sim \mathcal{N}(0^N, I^N)$ corresponding to the part k , and keep the rest of the noise vector $z[j] \mid j \neq k$ fixed, and perform the decoding on a given pose. When multiple elementary DensePose parts (*e.g.*, left head, right head) correspond to one logical body part for sampling (*e.g.*, head), we sample noise in all the elementary part vectors. The results are shown in Fig. 6.

To explicitly see how much our design choices help for part sampling, we perform a baseline experiment **NoParts** where we encode the appearance in a single vector $z \sim E(\mathbf{T}) \mid z \in \mathbb{Z}^N$ instead of part-specific vectors. We then warp this vector using the conditioning DensePose to create a noise image for the generator as described in Sec. 3.1. We compute the following two metrics for a given part p : 1) *Variation–Part*: mean pairwise L1 distance between the samples in the masked region (by DensePose) of the part p normalized by the masked area. 2) *Variation–Rest*: mean pairwise L1 distance between the samples in the masked region of all body parts excluding p . We compute the aforementioned metrics for the following parts: “Head”, “Upper body” and “Lower body” for 2500 generated images and provide our result in Table 3. A suitable method for

part sampling should generate diverse semantically meaningful renderings of a part without changing the rest of the body. This is confirmed by high *Variation–Part* and low *Variation–Rest* in our full method in comparison to the baseline *NoParts*.

Using part-specific latent vectors allows us to naturally perform

5. Limitations and Future Work

Our method sometimes demonstrates spurious interleaving of body parts or garments in the generated images (*e.g.*, see Fig. 4, right-most images). However, this is also shared by other (less realistic) generative human models (see the result of VUNet in Fig. 3 and the baselines in Fig. 5). These artifacts could be avoided by a hierarchical generator, where the garment style is generated first. This design, however, comes with the disadvantage of not being able to perform part-based sampling. We have also found that our generated images are biased towards females. This, we hypothesize, is due to the bias in the DeepFashion dataset.

garment transfer between two images representing the body and garments. To this end, we first encode the appearance of both the body image I_b and garment image I_g , in their part-based noise vectors z_b and z_g respectively. We then construct a new noise embedding that comprises of the body parts $z_b[p] \mid p \in Body$, and garment parts $z_g[p] \mid p \in Garments$ of the two noise embeddings, and use it in the generator for the final output, see Fig. 8.

6. Conclusion

We have presented a generative model for full-body images of clothed humans, which enables control of body pose, as well as independent control and sampling of appearance and clothing style on a body part level. A framework based on variational autoencoders is used to induce stochasticity in the appearance space. To achieve the disentanglement of pose and appearance, we encode the posterior probability of the part-specific latent vectors from a space of pose-independent appearance and warp the encoded vector to a different pose before performing the reconstruction. Experiments with pose-conditioned image generation, pose

transfer, as well as parts and garment transfer, have demonstrated that the model improves over the state-of-the-arts.

Acknowledgements. This work was supported by the ERC Consolidator Grant 4DReply (770784) and Lise Meitner Postdoctoral Fellowship. We thank Dushyant Mehta for the feedback on the draft.

References

- [1] Kfir Aberman, Mingyi Shi, Jing Liao, Dani Lischinski, Baoquan Chen, and Daniel Cohen-Or. Deep video-based performance cloning. *Comput. Graph. Forum*, 38(2):219–233, 2019. [2](#)
- [2] Aayush Bansal, Shugao Ma, Deva Ramanan, and Yaser Sheikh. Recycle-gan: Unsupervised video retargeting. In *ECCV*, 2018. [2](#)
- [3] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. *CoRR*, abs/1809.11096, 2018. [1, 2, 7](#)
- [4] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. Everybody dance now. In *International Conference on Computer Vision (ICCV)*, 2019. [2](#)
- [5] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. StarGAN: Unified generative adversarial networks for multi-domain image-to-image translation. In *Computer Vision and Pattern Recognition (CVPR), 2018*, 2018. [2](#)
- [6] Thang Doan, J. Monteiro, Isabela Albuquerque, Bogdan Mazzoure, A. Durand, Joelle Pineau, and R. Devon Hjelm. Online adaptative curriculum learning for gans. In *AAAI*, 2019. [2](#)
- [7] Ishan Durugkar, Ian Gemp, and Sridhar Mahadevan. Generative multi-adversarial networks, 2017. [2](#)
- [8] Patrick Esser, Ekaterina Sutter, and Björn Ommer. A variational u-net for conditional appearance and shape generation. In *Computer Vision and Pattern Recognition (CVPR)*, pages 8857–8866, 2018. [2, 5, 6, 7, 11, 12, 13, 14, 16, 20](#)
- [9] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27, pages 2672–2680. Curran Associates, Inc., 2014. [2](#)
- [10] A. K. Grigor’ev, Artem Sevastopolsky, Alexander Vakhitov, and Victor S. Lempitsky. Coordinate-based texture inpainting for pose-guided human image generation. *Computer Vision and Pattern Recognition (CVPR)*, pages 12127–12136, 2019. [2, 4, 6, 7, 16, 21](#)
- [11] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. *CVPR*, 2017. [1, 2, 5](#)
- [12] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. In *International Conference on Learning Representations*, 2018. [1, 2](#)
- [13] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4401–4410, 2019. [1, 2, 7](#)
- [14] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *Proc. CVPR*, 2020. [1, 2](#)
- [15] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8110–8119, 2020.
- [16] Diederick P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. [4](#)
- [17] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In Yoshua Bengio and Yann LeCun, editors, *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014. [2, 3, 4](#)
- [18] Bernhard Kratzwald, Zhiwu Huang, Danda Pani Paudel, and Luc Van Gool. Towards an understanding of our world by GANing videos in the wild. arXiv:1711.11453, 2017. [2](#)
- [19] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. In *International conference on machine learning*, pages 1558–1566. PMLR, 2016. [2, 6, 11](#)
- [20] Christoph Lassner, Gerard Pons-Moll, and Peter V. Gehler. A generative model for people in clothing. In *Proceedings of the IEEE International Conference on Computer Vision*, 2017. [2](#)
- [21] Lingjie Liu, Weipeng Xu, Marc Habermann, Michael Zollhöfer, Florian Bernard, Hyeongwoo Kim, Wenping Wang, and Christian Theobalt. Neural human video rendering by learning dynamic textures and rendering-to-video translation. *IEEE Transactions on Visualization and Computer Graphics*, PP:1–1, 05 2020. [2](#)
- [22] Lingjie Liu, Weipeng Xu, Michael Zollhoefer, Hyeongwoo Kim, Florian Bernard, Marc Habermann, Wenping Wang, and Christian Theobalt. Neural rendering and reenactment of human actor videos. *ACM Transactions on Graphics (TOG)*, 2019. [2](#)
- [23] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, pages 700—708, Red Hook, NY, USA, 2017. Curran Associates Inc. [2](#)
- [24] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphereface: Deep hypersphere embedding for face recognition. In *Computer Vision and Pattern Recognition (CVPR)*, pages 212–220, 2017. [4](#)
- [25] Z. Liu, P. Luo, S. Qiu, X. Wang, and X. Tang. Deepfashion: Powering robust clothes recognition and retrieval with rich annotations. In *Computer Vision and Pattern Recognition (CVPR)*, pages 1096–1104, 2016. [4](#)
- [26] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned

- multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, Oct. 2015. 2, 3
- [27] Liqian Ma, Xu Jia, Qianru Sun, Bernt Schiele, Tinne Tuytelaars, and Luc Van Gool. Pose guided person image generation. In *Advances in Neural Information Processing Systems*, pages 405–415, 2017. 2
- [28] Liqian Ma, Qianru Sun, Stamatios Georgoulis, Luc van Gool, Bernt Schiele, and Mario Fritz. Disentangled person image generation. *Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [29] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets, 2014. 2
- [30] Takeru Miyato and Masanori Koyama. cGANs with projection discriminator. In *International Conference on Learning Representations*, 2018. 2
- [31] Gonçalo Mordido, Haojin Yang, and C. Meinel. Dropoutgan: Learning from a dynamic ensemble of discriminators. *ArXiv*, abs/1807.11346, 2018. 2
- [32] Natalia Neverova, Riza Alp Güler, and Iasonas Kokkinos. Dense pose transfer. *European Conference on Computer Vision (ECCV)*, 2018. 2, 4, 6, 7, 16
- [33] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier GANs. volume 70 of *Proceedings of Machine Learning Research*, pages 2642–2651, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR. 2
- [34] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2
- [35] Albert Pumarola, Antonio Agudo, Alberto Sanfeliu, and Francesc Moreno-Noguer. Unsupervised person image synthesis in arbitrary poses. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2
- [36] Iasonas Kokkinos Rieza Alp Gueler, Natalia Neverova. Densepose: Dense human pose estimation in the wild. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 3
- [37] Kripasindhu Sarkar, Dushyant Mehta, Weipeng Xu, Vladislav Golyanik, and Christian Theobalt. Neural re-rendering of humans from a single image. In *European Conference on Computer Vision (ECCV)*, 2020. 2, 4, 6, 7, 8, 16, 21
- [38] Aliaksandr Siarohin, Stéphane Lathuilière, Enver Sangineto, and Nicu Sebe. Appearance and pose-conditioned human image generation using deformable gans. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2019. 2, 4, 6, 7, 16
- [39] Aliaksandr Siarohin, Enver Sangineto, Stephane Lathuiliere, and Nicu Sebe. Deformable GANs for pose-based human image generation. In *CVPR 2018*, 2018. 2
- [40] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 4
- [41] Ayush Tewari, Mohamed Elgharib, Gaurav Bharaj, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zöllhofer, and Christian Theobalt. Stylerig: Rigging stylegan for 3d control over portrait images. *cvpr 2020*. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, june 2020. 1, 2, 7
- [42] Ayush Tewari, Mohamed Elgharib, Mallikarjun BR, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zöllhofer, and Christian Theobalt. Pie: Portrait image embedding for semantic control. volume 39, December 2020. 1, 2, 7
- [43] Gül Varol, Javier Romero, Xavier Martin, Naureen Mahmood, Michael J. Black, Ivan Laptev, and Cordelia Schmid. Learning from synthetic humans. In *Computer Vision and Pattern Recognition (CVPR)*, 2017. 3
- [44] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 2
- [45] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 1, 2, 4
- [46] Zili Yi, Hao Zhang, Ping Tan, and Minglun Gong. Dual-GAN: Unsupervised dual learning for image-to-image translation. In *ICCV*, pages 2868–2876, Oct. 2017. 2
- [47] Polina Zablotckaia, Aliaksandr Siarohin, Leonid Sigal, and Bo Zhao. Dwnet: Dense warp-based network for pose-guided human video generation. In *British Machine Vision Conference (BMVC)*, 2019. 5, 7
- [48] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. volume 97 of *Proceedings of Machine Learning Research*, pages 7354–7363, Long Beach, California, USA, 09–15 Jun 2019. PMLR. 2
- [49] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 5, 6, 7
- [50] Zhou 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*, 13(4):600–612, 2004. 6
- [51] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networkss. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, 2017. 2
- [52] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 465–476. Curran Associates, Inc., 2017. 6
- [53] Zhen Zhu, Tengteng Huang, Baoguang Shi, Miao Yu, Bofei Wang, and Xiang Bai. Progressive pose attention transfer for person image generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2347–2356, 2019. 2

A. Appendix

This appendix complements the main manuscript and provides more qualitative results as well as questions from the user study.

A.1. More Qualitative Results

Appearance Sampling. Figs. 9, 10, and 11 show the qualitative results of our method for pose-guided image generation and its comparison with the baselines *VUNet* [8], *Pix2PixHD+Noise*, *Pix2PixHD+WNoise* and *DAE* (Deterministic Auto Encoder). We observe that purely noise-based baselines such as *Pix2PixHD+Noise* not only failed to generate diverse output but also lacked realism. *VUNet* produces a diverse set of outputs but often shows spurious patterns. It also lacks realism because of its GAN-free architecture and less powerful generator. Our approach, in contrast, produces samples of high quality than the baseline methods.

Image Interpolation. Fig. 12 shows the resulting images generated by interpolating the appearance vector between two encodings: given two encodings z_1 and z_2 of two different human images, we generate a human image with the interpolated encoding z as

$$z = z_1 t + z_2 (1 - t), \quad t \in [0, 1]. \quad (4)$$

We find that intermediate images show coherently dressed humans that share the properties of both input images.

Pose Transfer. Fig. 13 shows our results for pose transfer and its comparison with the state-of-the-art methods. Our results show higher realism and are more visually pleasing compared to the other baselines. However, it misses some fine-scaled details in a few cases. See Figs. 17 and 18 for more samples.

Part Sampling and Garment Transfer. Fig. 14 shows our results for part sampling. Here, we only change the latent vectors representing a specific part (*e.g.*, head). We observe that the rest of the body does not change considerably with the change in the generated image parts. However, we

also observe that our method is biased towards generating realistic outputs over generating images that are highly different in the part regions but not coherent as a whole (*e.g.*, sampling the head and garments of a female will result in images with female heads). Fig. 15 shows our garment transfer results, where we use the appearance encodings of two different images corresponding to the body and garment parts.

A.2. Comparisons to Additional Baselines

Variational autoencoders are notoriously difficult to train, and we have made several observations while developing and training our model. Appending noise to the conditioning DensePose image (baseline *Pix2PixHD+Noise*) not only failed to generate diverse output but also lacked realism. Global latent vectors (instead of the part-specific) for appearance (baseline *NoParts*) lacked realism as well. This section shows the results of two additional baselines (also introduced in the main manuscript):

a) **+DPCond.** This baseline conditions the DensePose image P_t in addition to the noise image Z_t to the generator. *i.e.* we concatenate P_t and Z_t channel-wise and input the resulting tensor to the generator.

b) **+NoisePrior.** This baseline use samples from the prior $\mathcal{N}(0, 1)$ (along with the encoded distribution) during the training, as recommended by Larsen *et al.* [19].

The qualitative results are shown in Fig. 16. We observe that conditioning P_t along with the noise image Z_t to the generator, resulted in fewer variations during sampling. We assume the reason to be the overpowering of the conditioning variable — it does not let the latent vectors learn semantics. In our HumanGAN, we force the generator to produce output only from the warped noise vector, thereby enforcing semantics for sampling. We also observe that using samples from the prior $\mathcal{N}(0, 1)$ during the training leads to high variation during sampling. However, it creates highly distorted faces and other body parts.

A.3. User Studies

In Figs. 17 and 18, we show the list of questions used in the two user studies.

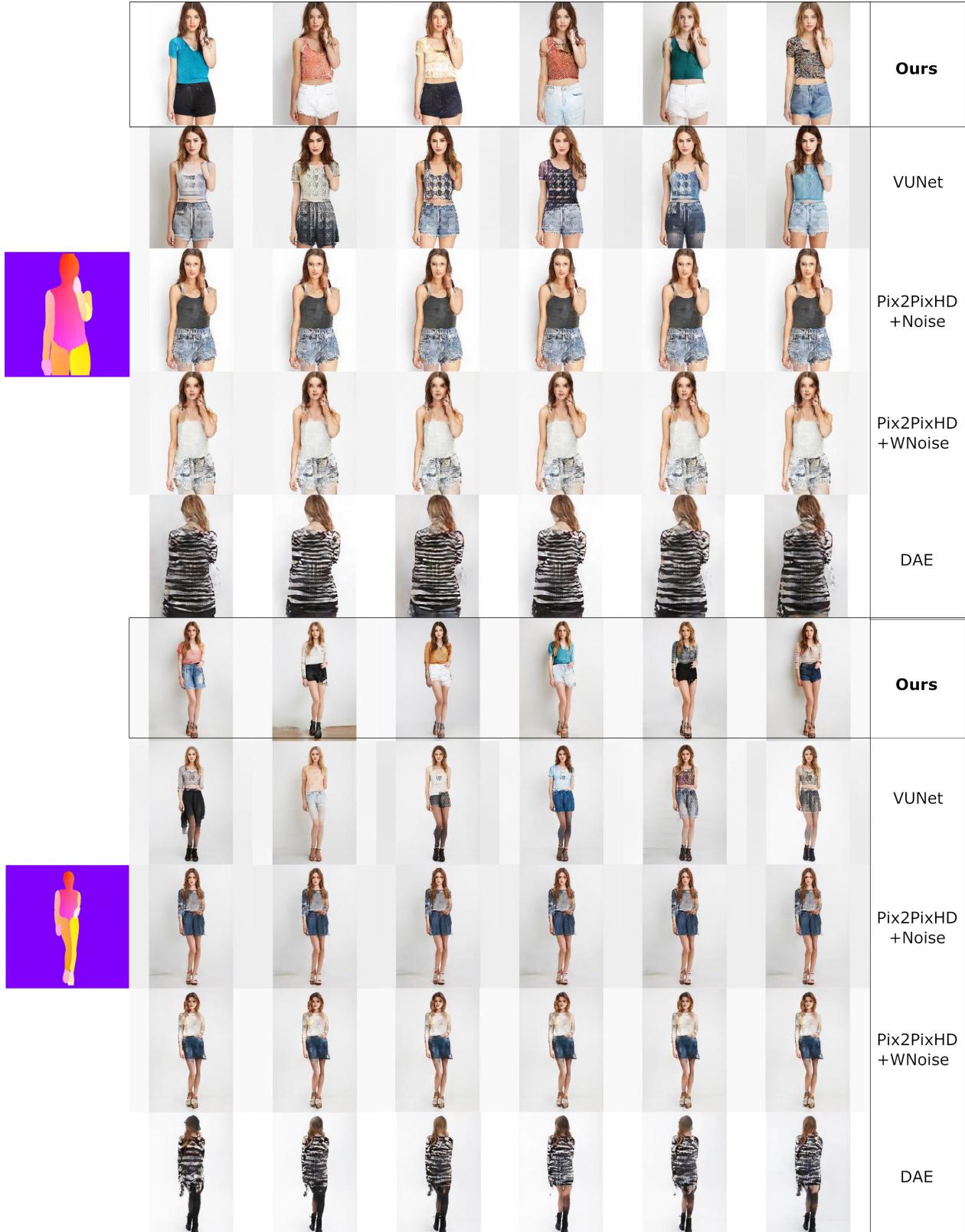


Figure 9: Results of our method for pose-guided image generation and its comparison with VUNet [8] and other baselines. The conditioning pose in the form of DensePose is shown in the left column.



Figure 10: Results of our method for pose-guided image generation and its comparison with VUNet [8] and other baselines. The conditioning pose in the form of DensePose is shown in the left column.



Figure 11: Results of our method for pose-guided image generation and its comparison with VUNet [8] and other baselines. The conditioning pose in the form of DensePose is shown in the left column.

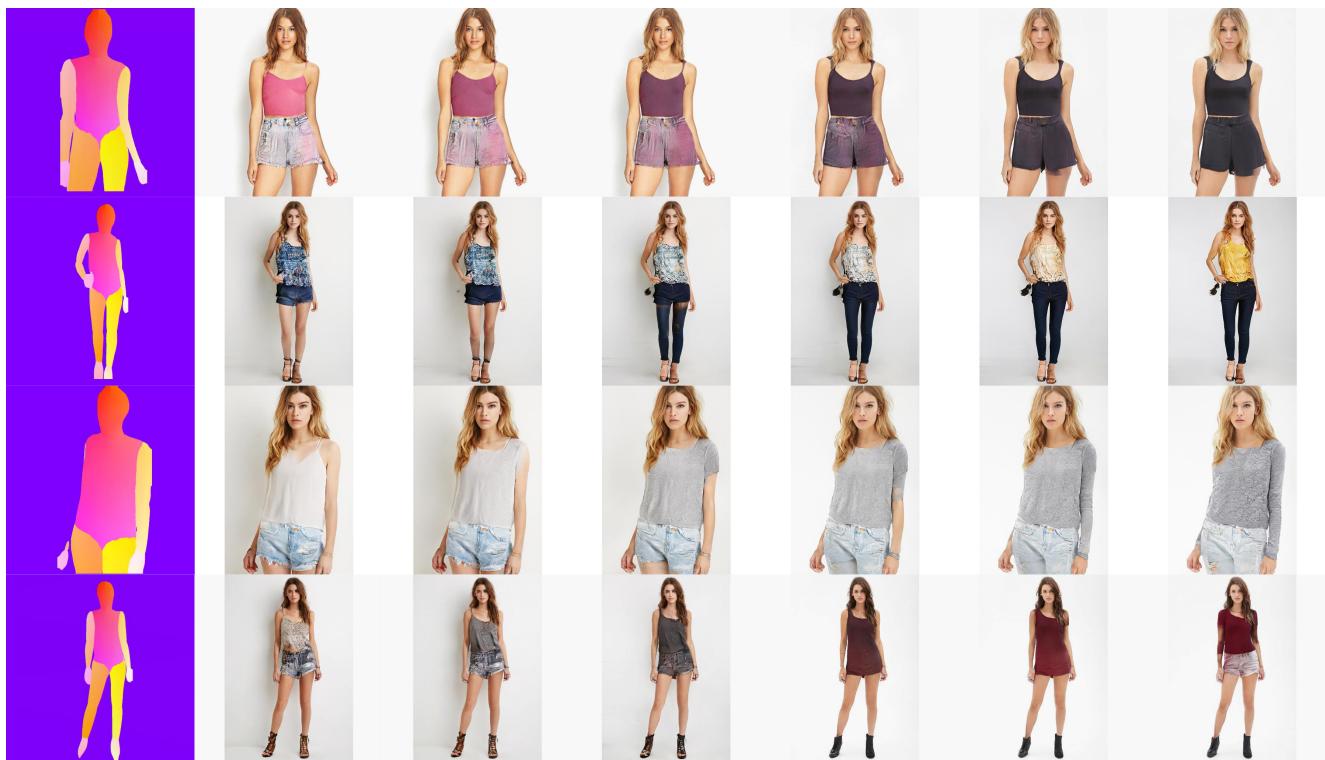


Figure 12: Generated images with interpolated appearance encodings. The conditioning pose is shown on the left.



Figure 13: Pose Transfer. Comparison of our reconstruction+transfer results with the state-of-the-art pose transfer methods CBI [10], NHRR [37], DSC [38], VUNet [8] and DPT [32]. Our HumanGAN produces more realistic renderings than the competing methods.



Figure 14: Our results for part sampling (head, lower body and upper body). Conditioning pose is shown in the left column.

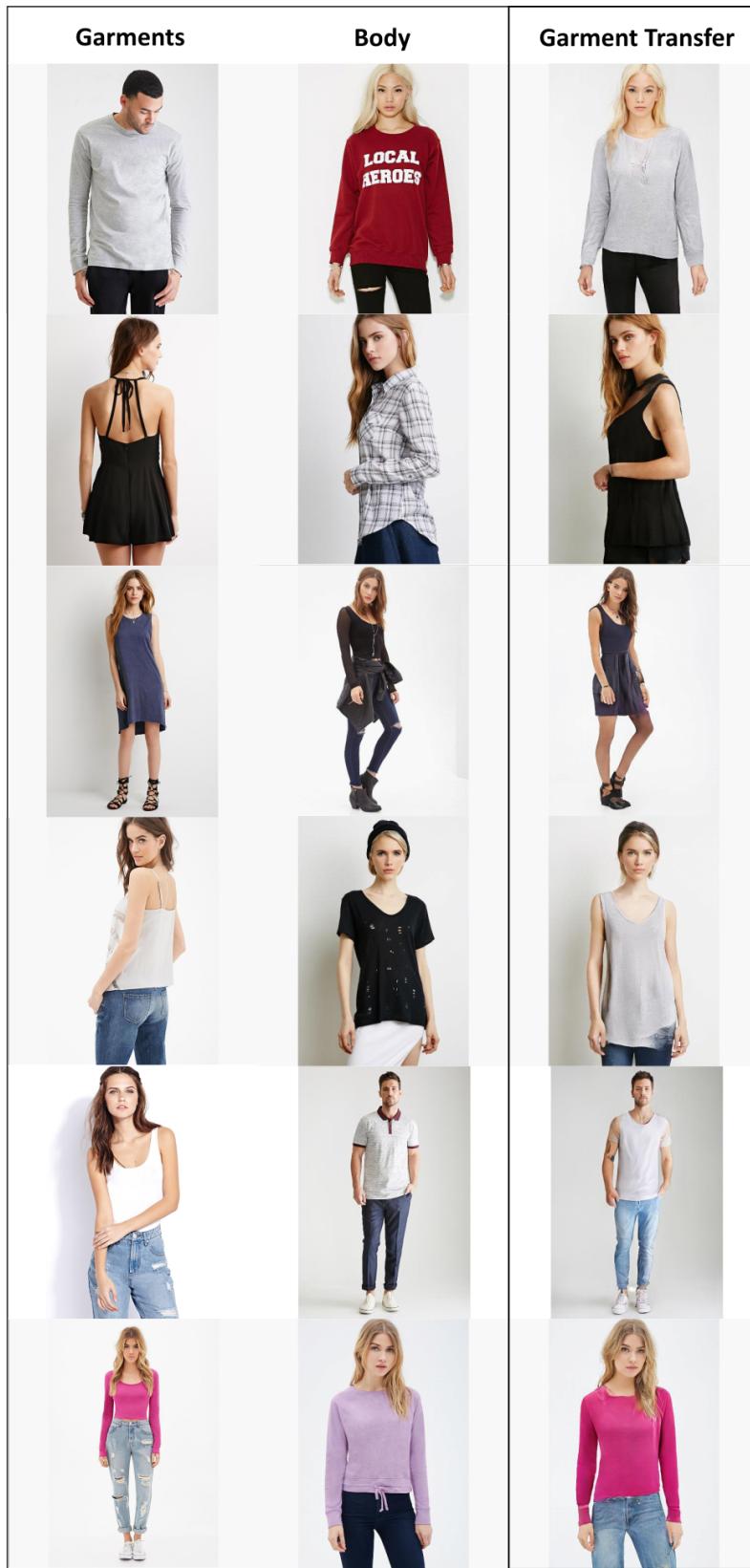


Figure 15: Our results for garment transfer.

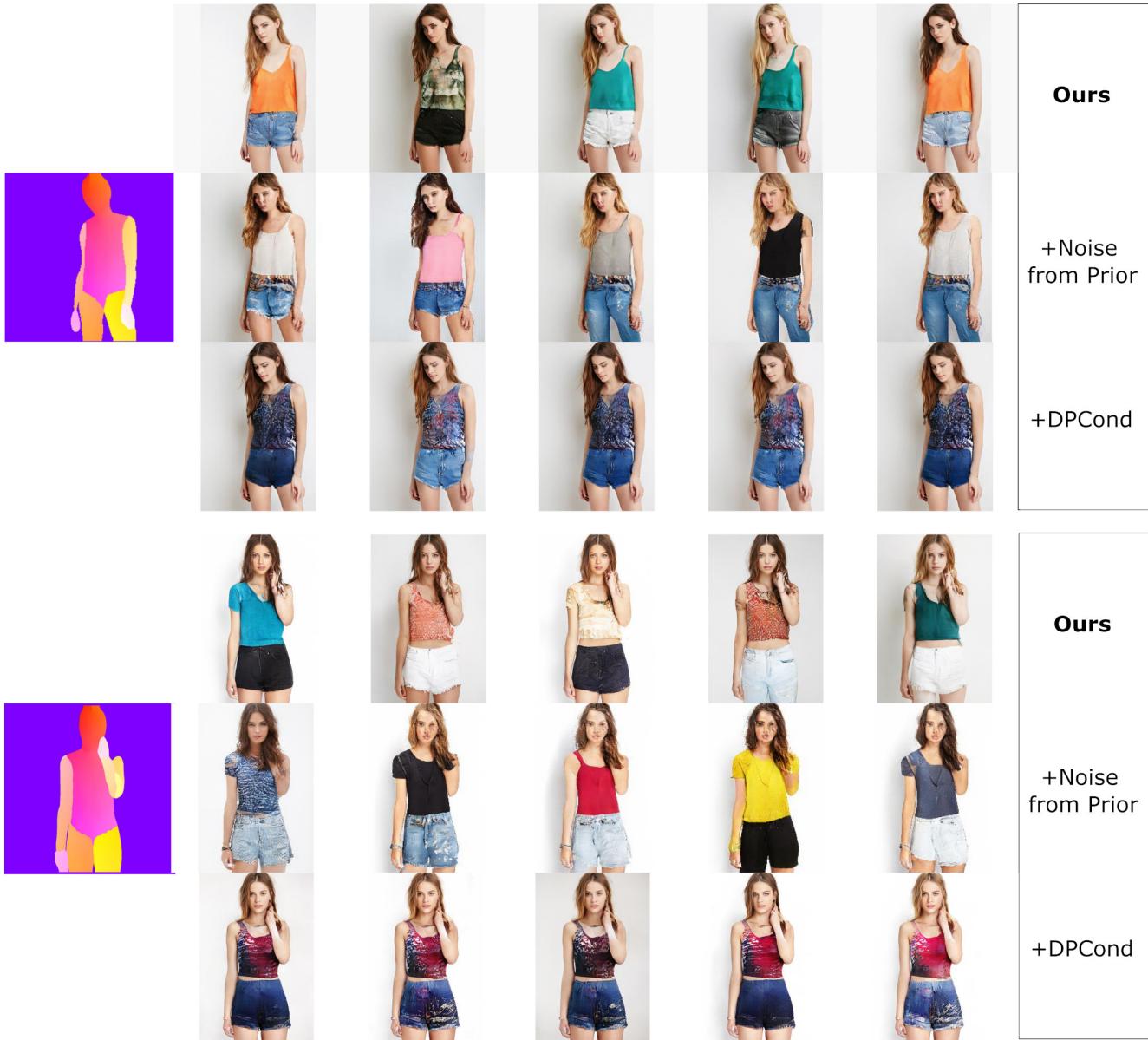


Figure 16: Comparison of our method with additional baselines. We observe that *+Noise from Prior* leads to less realistic images, while *+DPCond* leads to less variation during sampling. See Sec. A.2 for more details.

Questions 1 – 9 (image sets)	ours	VUnet	VUnet	ours	VUnet	ours
						
						
	VUnet	ours	ours	VUnet	ours	VUnet
						
						
	VUnet	ours	ours	VUnet		
						
						
Questions 10 – 20 (individual images)	VUnet	ours	ours	VUnet	VUnet	ours
						
	ours	VUnet	VUnet	ours	VUnet	ours
						
	ours	VUnet	ours	VUnet	ours	VUnet
						
	VUnet	ours	VUnet	ours	ours	VUnet
						

Figure 17: The samples and the sets used in the *first user study* where we compare our results with VUNet [8] for appearance sampling. The keys on the top left were not shown during the user study (they are replaced with A and B variants).

Questions 21 – 36 (pose transfer)

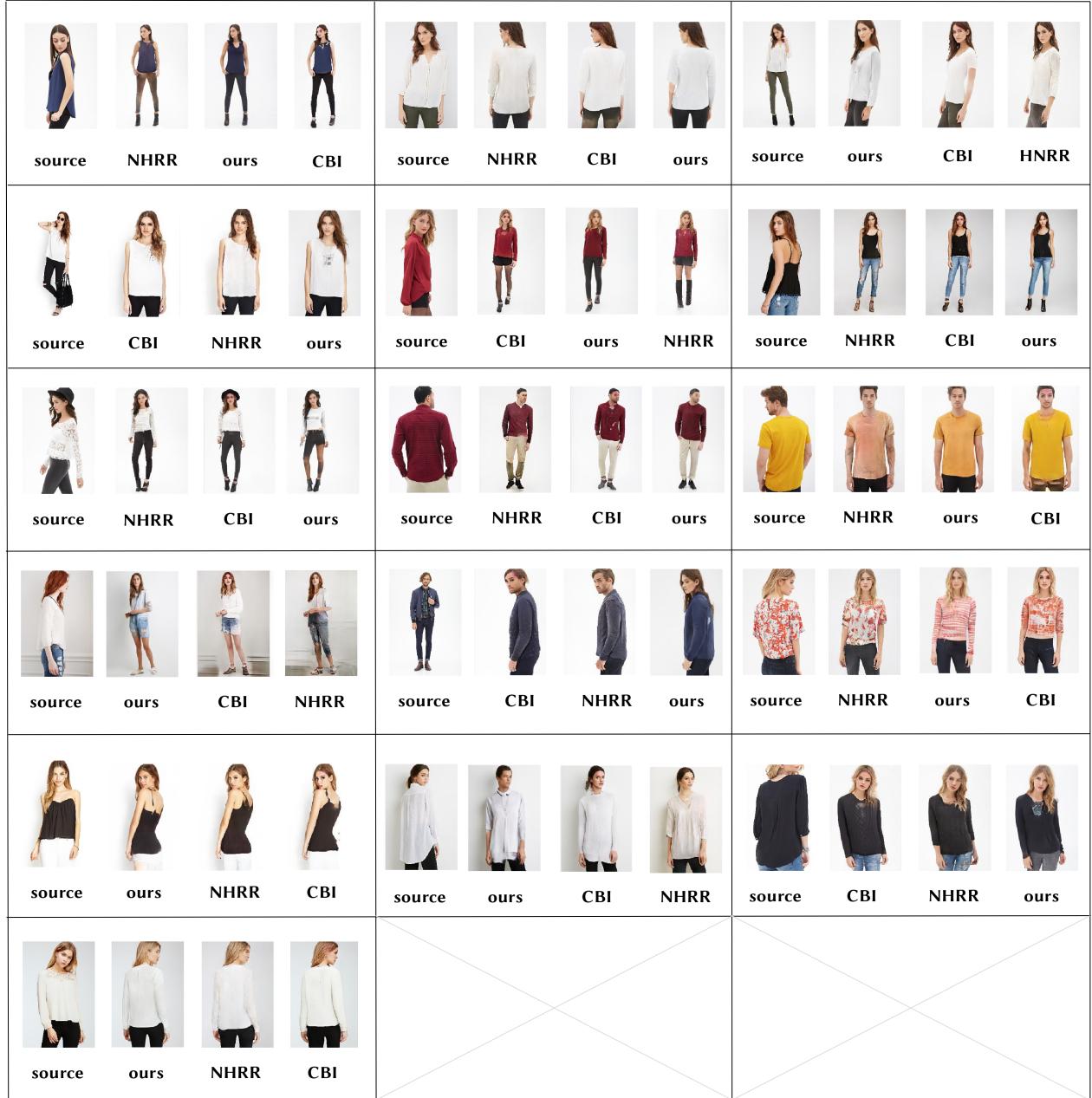


Figure 18: The samples and the sets used in the *second user study* where we compare our results with CBI [10] and NHRR [37] for pose transfer. The keys on the bottom were not shown during the user study (they are replaced with A, B and C variants).