

VITON-HD: High-Resolution Virtual Try-On via Misalignment-Aware Normalization

Seunghwan Choi* Sunghyun Park* Minsoo Lee* Jaegul Choo
KAIST, Daejeon, South Korea

{shadow2496, psh01087, alstn2022, jchoo}@kaist.ac.kr

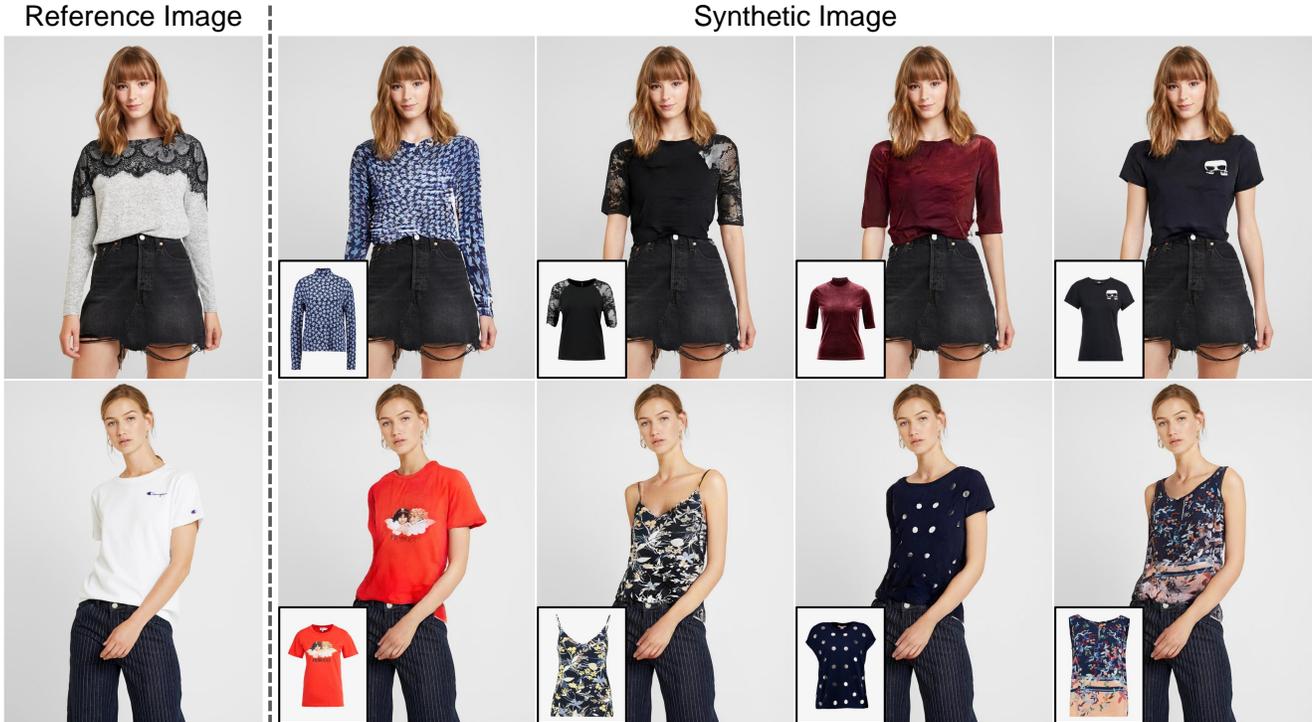


Figure 1: Given a pair of a reference image (containing a person) and a target clothing image, our method successfully synthesizes 1024×768 virtual try-on images.

Abstract

The task of image-based virtual try-on aims to transfer a target clothing item onto the corresponding region of a person, which is commonly tackled by fitting the item to the desired body part and fusing the warped item with the person. While an increasing number of studies have been conducted, the resolution of synthesized images is still limited to low (e.g., 256×192), which acts as the critical limitation against satisfying online consumers. We argue that the limitation stems from several challenges: as the resolution increases, the artifacts in the misaligned areas between the warped clothes and the desired clothing regions become noticeable in the final results; the architectures used in ex-

isting methods have low performance in generating high-quality body parts and maintaining the texture sharpness of the clothes. To address the challenges, we propose a novel virtual try-on method called VITON-HD that successfully synthesizes 1024×768 virtual try-on images. Specifically, we first prepare the segmentation map to guide our virtual try-on synthesis, and then roughly fit the target clothing item to a given person's body. Next, we propose ALIGNment-Aware Segment (ALIAS) normalization and ALIAS generator to handle the misaligned areas and preserve the details of 1024×768 inputs. Through rigorous comparison with existing methods, we demonstrate that VITON-HD highly surpasses the baselines in terms of synthesized image quality both qualitatively and quantitatively. Code is available at <https://github.com/shadow2496/VITON-HD>.

* These authors contributed equally.

1. Introduction

Image-based virtual try-on refers to the image generation task of changing the clothing item on a person into a different item, given in a separate product image. With a growing trend toward online shopping, virtually wearing the clothes can enrich a customer’s experience, as it gives an idea about how these items would look on them.

Virtual try-on is similar to image synthesis, but it has unique and challenging aspects. Given images of a person and a clothing product, the synthetic image should meet the following criteria: (1) The person’s pose, body shape, and identity should be preserved. (2) The clothing product should be naturally deformed to the desired clothing region of the given person, by reflecting his/her pose and body shape. (3) Details of the clothing product should be kept intact. (4) The body parts initially occluded by the person’s clothes in the original image should be properly rendered. Since the given clothing image is not initially fitted to the person image, fulfilling these requirements is challenging, which leaves the development of virtual try-on still far behind the expectations of online consumers. In particular, the resolution of virtual try-on images is low compared to the one of normal pictures on online shopping websites.

After Han *et al.* [10] proposed VITON, various image-based virtual try-on methods have been proposed [31, 36, 35, 6]. These methods follow two processes in common: (1) warping the clothing image initially to fit the human body; (2) fusing the warped clothing image and the image of the person that includes pixel-level refinement. Also, several recent methods [9, 36, 35] add a module that generates segmentation maps and determine the person’s layout from the final image in advance.

However, the resolution of the synthetic images from the previous methods is low (*e.g.*, 256×192) due to the following reasons. First, the misalignment between the warped clothes and a person’s body results in the artifacts in the misaligned regions, which become noticeable as the image size increases. It is difficult to warp clothing images to fit the body perfectly, so the misalignment occurs as shown in Fig. 2. Most of previous approaches utilize the thin-plate spline (TPS) transformation to deform clothing images. To accurately deform clothes, ClothFlow [9] predicts the optical flow maps of the clothes and the desired clothing regions. However, the optical flow maps does not remove the misalignment completely on account of the regularization. In addition, the process requires more computational costs than other methods due to the need of predicting the movement of clothes at a pixel level. (The detailed analysis of ClothFlow is included in the supplementary.) Second, a simple U-Net architecture [25] used in existing approaches is insufficient in synthesizing initially occluded body parts in final high-resolution (*e.g.*, 1024×768) images. As noted in Wang *et al.* [32], applying a simple U-Net-based archi-



Figure 2: An example of misaligned regions.

ture to generate high-resolution images leads to unstable training as well as unsatisfactory quality of generated images. Also, refining the images once at the pixel level is insufficient in preserving the details of high-resolution clothing images.

To address the above-mentioned challenges, we propose a novel high-resolution virtual try-on method, called VITON-HD. In particular, we introduce a new clothing-agnostic person representation that leverages the pose information and the segmentation map so that the clothing information is eliminated thoroughly. Afterwards, we feed the segmentation map and the clothing item deformed to fit the given human body to the model. Using the additional information, our novel ALIgment-Aware Segment (ALIAS) normalization removes information irrelevant to the clothing texture in the misaligned regions and propagates the semantic information throughout the network. The normalization separately standardizes the activations corresponding to the misaligned regions and the other regions, and modulates the standardized activations using the segmentation map. Our ALIAS generator employing ALIAS normalization synthesizes the person image wearing the target product while filling the misaligned regions with the clothing texture and preserving the details of the clothing item through the multi-scale refinement at a feature level. To validate the performance of our framework, we collected a 1024×768 dataset that consists of pairs of a person and a clothing item for our research purpose. Our experiments demonstrate that VITON-HD significantly outperforms the existing methods in generating 1024×768 images, both quantitatively and qualitatively. We also confirm the superior capability of our novel ALIAS normalization module in dealing with the misaligned regions.

We summarize our contributions as follows:

- We propose a novel image-based virtual try-on approach called VITON-HD, which is, to the best of our knowledge, the first model to successfully synthesize 1024×768 images.
- We introduce a clothing-agnostic person representation that allows our model to remove the dependency on the clothing item originally worn by the person.

- To address the misalignment between the warped clothes and the desired clothing regions, we propose ALIAS normalization and ALIAS generator, which is effective in maintaining the details of clothes.
- We demonstrate the superior performance of our method through experiments with baselines on the newly collected dataset.

2. Related Work

Conditional Image Synthesis. Conditional generative adversarial networks (cGANs) utilize additional information, such as class labels [19, 2], text [24, 34], and attributes [28], to steer the image generation process. Since the emergence of pix2pix [14], numerous cGANs conditioned on input images have been proposed to generate high-resolution images in a stable manner [32, 1, 21]. However, these methods tend to generate blurry images when handling a large spatial deformation between the input image and the target image. In this paper, we propose a method that can address the spatial deformation of input images and properly generate 1024×768 images.

Normalization Layers. Normalization layers [13, 30] have been widely applied in modern deep neural networks. Normalization layers, whose affine parameters are estimated with external data, are called conditional normalization layers. Conditional batch normalization [4] and adaptive instance normalization [12] are such conditional normalization techniques and have been used in style transfer tasks. SPADE [20] and SEAN [39] utilize segmentation maps to apply spatially varying affine transformations. Using the misalignment mask as external data, our proposed normalization layer computes the means and the variances of the misaligned area and the other area within an instance separately. After standardization, we modulate standardized activation maps with affine parameters inferred from human-parsing maps to preserve semantic information.

Virtual Try-On Approaches. There are two main categories for virtual try-on approaches: 3D model-based approaches [8, 27, 23, 22] and 2D image-based approaches [10, 31, 9, 36, 35, 5]. 3D model-based approaches can accurately simulate the clothes but are not widely applicable due to their dependency on 3D measurement data.

2D image-based approaches do not rely on any 3D information, thus being computationally efficient and appropriate for practical use. Jetchev and Bergmann [15] proposed CAGAN, which first introduced the task of swapping fashion articles on human images. VITON [10] addressed the same problem by proposing a coarse-to-fine synthesis framework that involves TPS transformation of clothes. Most existing virtual try-on methods tackle different aspects of VITON to synthesize perceptually convincing photo-realistic images. CP-VTON [31] adopted a geometric matching module to learn the parameters of TPS transfor-

mation, which improves the accuracy of deformation. VT-NFP [36] and ACGPN [35] predicted the human-parsing maps of a person wearing the target clothes in advance to guide the try-on image synthesis. Even though the image quality at high resolution is an essential factor in evaluating the practicality of the generated images, none of the methods listed above could generate such photo-realistic images at high resolution.

3. Proposed Method

Model Overview. As described in Fig. 3, given a reference image $I \in \mathbb{R}^{3 \times H \times W}$ of a person and a clothing image $c \in \mathbb{R}^{3 \times H \times W}$ (H and W denote the image height and width, respectively), the goal of VITON-HD is to generate a synthetic image $\hat{I} \in \mathbb{R}^{3 \times H \times W}$ of the same person wearing the target clothes c , where the pose and body shape of I and the details of c are preserved. While training the model with (I, c, \hat{I}) triplets is straightforward, construction of such dataset is costly. Instead, we use (I, c, I) where the person in the reference image I is already wearing c .

Since directly training on (I, c, I) can harm the model’s generalization ability at test time, we first compose a clothing-agnostic person representation that leaves out the clothing information in I and use it as an input. Our new clothing-agnostic person representation uses both the pose map and the segmentation map of the person to eliminate the clothing information in I (Section 3.1). The model generates the segmentation map from the clothing-agnostic person representation to help the generation of \hat{I} (Section 3.2). We then deform c to roughly align it to the human body (Section 3.3). Lastly, we propose the ALIgnment-Aware Segment (ALIAS) normalization that removes the misleading information in the misaligned area after deforming c . ALIAS generator fills the misaligned area with the clothing texture and maintains the clothing details (Section 3.4).

3.1. Clothing-Agnostic Person Representation

To train the model with pairs of c and I already wearing c , a person representation without the clothing information in I has been utilized in the virtual try-on task. Such representations have to satisfy the following conditions: (1) the original clothing item to be replaced should be deleted; (2) sufficient information to predict the pose and the body shape of the person should be maintained; (3) the regions to be preserved (*e.g.*, face and hands) should be kept to maintain the person’s identity.

Problems of Existing Person Representations. In order to maintain the person’s shape, several approaches [10, 31, 36] provide a coarse body shape mask as a cue to synthesize the image, but fail to reproduce the body parts elaborately (*e.g.*, hands). To tackle this issue, ACGPN [35] employs the detailed body shape mask as the input, and the neural network attempts to discard the clothing informa-

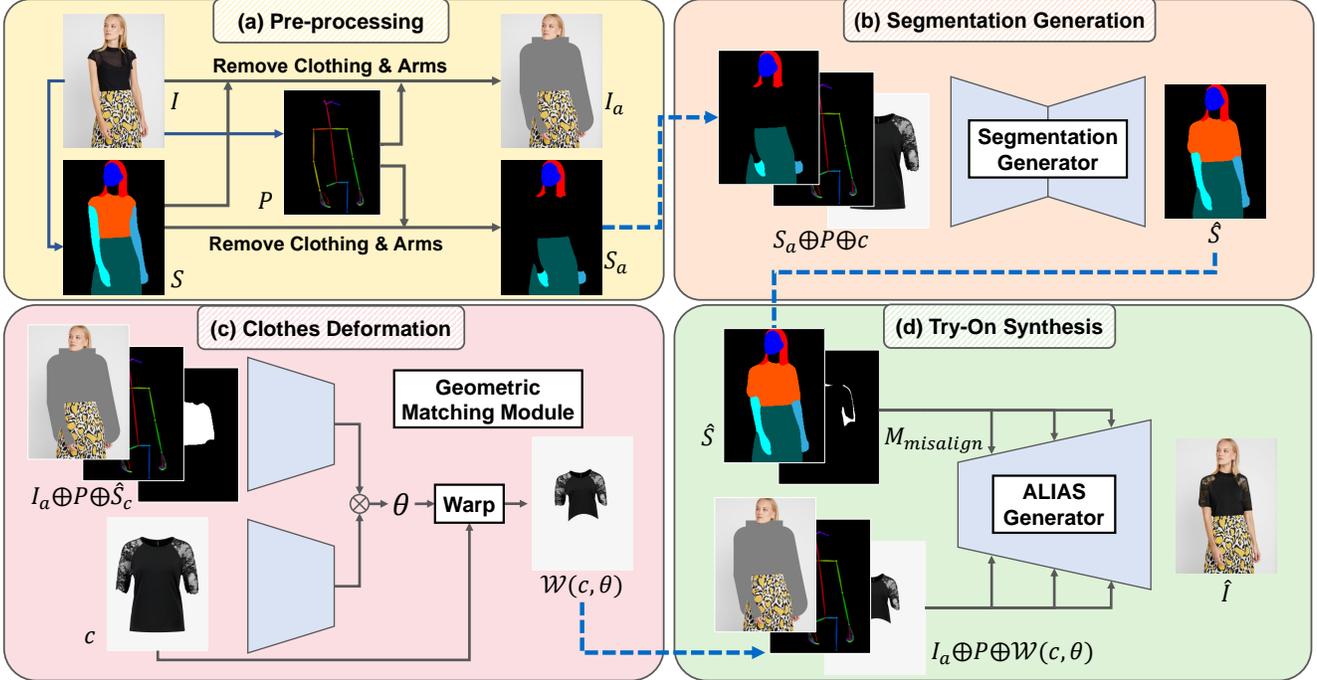


Figure 3: Overview of a VITON-HD. (a) First, given a reference image I containing a target person, we predict the segmentation map S and the pose map P , and utilize them to pre-process I and S as a clothing-agnostic person image I_a and segmentation S_a . (b) Segmentation generator produces the synthetic segmentation \hat{S} from (S_a, P, c) . (c) Geometric matching module deforms the clothing image c according to the predicted clothing segmentation \hat{S}_c extracted from \hat{S} . (d) Finally, ALIAS generator synthesizes the final output image \hat{I} based on the outputs from the previous stages via our ALIAS normalization.

tion to be replaced. However, since the body shape mask includes the shape of the clothing item, neither the coarse body shape mask nor the neural network could perfectly eliminate the clothing information. As a result, the original clothing item that is not completely removed causes problems in the test phase.

Clothing-Agnostic Person Representation. We propose a clothing-agnostic image I_a and a clothing-agnostic segmentation map S_a as inputs of each stage, which truly eliminate the shape of clothing item and preserve the body parts that need to be reproduced. We first predict the segmentation map $S \in \mathbb{L}^{H \times W}$ and the pose map $P \in \mathbb{R}^{3 \times H \times W}$ of the image I by utilizing the pre-trained networks [7, 3], where \mathbb{L} is a set of integers indicating the semantic labels. The segmentation map S is used to remove the clothing region to be replaced and preserve the rest of the image. The pose map P is utilized to remove the arms, but not the hands, as they are difficult to reproduce. Based on S and P , we generate the clothing-agnostic image I_a and the clothing-agnostic segmentation map S_a , which allow the model to remove the original clothing information thoroughly, and preserve the rest of the image. In addition, unlike other previous work, which adopts the pose heatmap with each channel corresponded to one keypoint, we con-

catenate I_a or S_a to the RGB pose map P representing a skeletal structure that improves generation quality.

3.2. Segmentation Generation

Given the clothing-agnostic person representation (S_a, P) , and the target clothing item c , the segmentation generator G_S predicts the segmentation map $\hat{S} \in \mathbb{L}^{H \times W}$ of the person in the reference image wearing c . We train G_S to learn the mapping between S and (S_a, P, c) , in which the original clothing item information is completely removed. As the architecture of G_S , we adopt U-Net [25], and the total loss \mathcal{L}_S of the segmentation generator are written as

$$\mathcal{L}_S = \mathcal{L}_{cGAN} + \lambda_{CE} \mathcal{L}_{CE}, \quad (1)$$

where \mathcal{L}_{CE} and \mathcal{L}_{cGAN} denote the pixel-wise cross-entropy loss and conditional adversarial loss between \hat{S} and S , respectively. λ_{CE} is the hyperparameter corresponding to the relative importance between two losses.

3.3. Clothing Image Deformation

In this stage, we deform the target clothing item c to align it with \hat{S}_c , which is the clothing area of \hat{S} . We employ the geometric matching module proposed in CP-VTON [31] with the clothing-agnostic person representa-

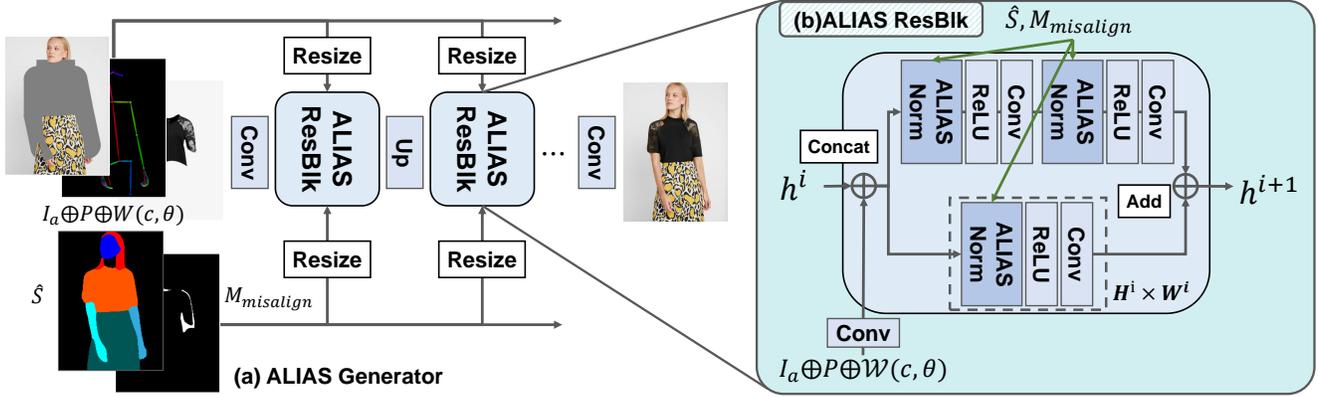


Figure 4: ALIAS generator. (a) The ALIAS generator is composed of a series of ALIAS residual blocks, along with up-sampling layers. The input $(I_a, P, \mathcal{W}(c, \theta))$ is resized and injected into each layer of the generator. (b) A detailed view of a ALIAS residual block. Resized $(I_a, P, \mathcal{W}(c, \theta))$ is concatenated to h^i after passing through a convolution layer. Each ALIAS normalization layer leverages resized \hat{S} and $M_{misalign}$ to normalize the activation.

tion (I_a, P) and \hat{S}_c as inputs. A correlation matrix between the features extracted from (I_a, P) and c is first calculated. With the correlation matrix as an input, the regression network predicts the TPS transformation parameters $\theta \in \mathbb{R}^{2 \times 5 \times 5}$, and then c is warped by θ . In the training phase, the model takes S_c extracted from S instead of \hat{S}_c . The module is trained with the L1 loss between the warped clothes and the clothes I_c that is extracted from I . In addition, the second-order difference constraint [35] is adopted to reduce obvious distortions in the warped clothing images from deformation. The overall objective function to warp the clothes to fit the human body is written as

$$\mathcal{L}_{warp} = \|I_c - \mathcal{W}(c, \theta)\|_{1,1} + \lambda_{const} \mathcal{L}_{const}, \quad (2)$$

where \mathcal{W} is the function that deforms c using θ , \mathcal{L}_{const} is a second-order difference constraint, and λ_{const} is the hyper-parameter for \mathcal{L}_{const} .

3.4. Try-On Synthesis via ALIAS Normalization

We aim to generate the final synthetic image \hat{I} based on the outputs from the previous stages. Overall, we fuse the clothing-agnostic person representation (I_a, P) and the warped clothing image $\mathcal{W}(c, \theta)$, guided by \hat{S} . $(I_a, P, \mathcal{W}(c, \theta))$ is injected into each layer of the generator. For \hat{S} , we propose a new conditional normalization method called the ALIgnment-Aware Segment (ALIAS) normalization. ALIAS normalization enables the preservation of semantic information, and the removal of misleading information from the misaligned regions by leveraging \hat{S} and the mask of these regions.

Alignment-Aware Segment Normalization. Let us denote $h^i \in \mathbb{R}^{N \times C^i \times H^i \times W^i}$ as the activation of the i -th layer of a network for a batch of N samples, where H^i , W^i , and C^i indicate the height, width, and the number of channels of h^i , respectively. ALIAS normalization has two inputs:

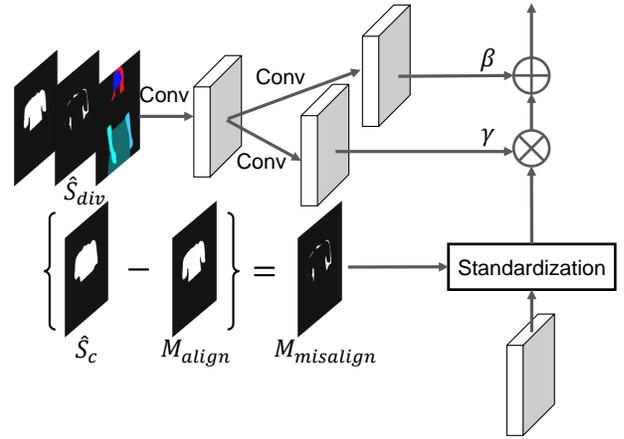


Figure 5: ALIAS normalization. First, the activation is separately standardized according to the regions divided by $M_{misalign}$, which can be obtained from the difference between \hat{S}_c and M_{align} . Next, \hat{S}_{div} is convolved to create the modulation parameters γ and β , and then the standardized activation is modulated with the parameters γ and β .

(1) the synthetic segmentation map \hat{S} ; (2) the misalignment binary mask $M_{misalign} \in \mathbb{L}^{H \times W}$, which excludes the warped mask of the target clothing image $\mathcal{W}(M_c, \theta)$ from \hat{S}_c (M_c denotes the target clothing mask), i.e.,

$$M_{align} = \hat{S}_c \cap \mathcal{W}(M_c, \theta) \quad (3)$$

$$M_{misalign} = \hat{S}_c - M_{align}. \quad (4)$$

Fig. 5 illustrates the workflow of the ALIAS normalization. We first obtain M_{align} and $M_{misalign}$ from Eq. (3) and Eq. (4). We define the modified version of \hat{S} as \hat{S}_{div} , where \hat{S}_c in \hat{S} separates into M_{align} and $M_{misalign}$. ALIAS normalization standardizes the regions of $M_{misalign}$ and the other regions in h^i separately, and then modulates the standardized activation using affine transformation parameters

	256 × 192			512 × 384			1024 × 768		
	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow
CP-VTON	0.739	0.159	56.23	0.791	0.141	31.96	0.786	0.158	43.28
ACGPN	0.842	0.064	26.45	0.863	0.067	15.22	0.856	0.102	43.39
VITON-HD*	-	-	-	-	-	-	0.893	0.054	12.47
VITON-HD	0.844	0.062	27.83	0.870	0.052	14.05	0.895	0.053	11.74

Table 1: Quantitative comparison with baselines across different resolutions. VITON-HD* is a VITON-HD variant where the standardization in ALIAS normalization is replaced by channel-wise standardization as in the original instance normalization. For the SSIM, higher is better. For the LPIPS and the FID, lower is better.

inferred from \hat{S}_{div} . The activation value at site $(n \in N, k \in C^i, y \in H^i, x \in W^i)$ is calculated by

$$\gamma_{k,y,x}^i(\hat{S}_{div}) \frac{h_{n,k,y,x}^i - \mu_{n,k}^{i,m}}{\sigma_{n,k}^{i,m}} + \beta_{k,y,x}^i(\hat{S}_{div}), \quad (5)$$

where $h_{n,k,y,x}^i$ is the activation at the site before normalization and $\gamma_{k,y,x}^i$ and $\beta_{k,y,x}^i$ are the functions that convert \hat{S}_{div} to modulation parameters of the normalization layer. $\mu_{n,k}^{i,m}$ and $\sigma_{n,k}^{i,m}$ are the mean and standard deviation of the activation in sample n and channel k . $\mu_{n,k}^{i,m}$ and $\sigma_{n,k}^{i,m}$ are calculated by

$$\mu_{n,k}^{i,m} = \frac{1}{|\Omega_n^{i,m}|} \sum_{(y,x) \in \Omega_n^{i,m}} h_{n,k,y,x}^i \quad (6)$$

$$\sigma_{n,k}^{i,m} = \sqrt{\frac{1}{|\Omega_n^{i,m}|} \sum_{(y,x) \in \Omega_n^{i,m}} (h_{n,k,y,x}^i - \mu_{n,k}^{i,m})^2}, \quad (7)$$

where $\Omega_n^{i,m}$ denotes the set of pixels in region m , which is $M_{misalign}$ or the other region, and $|\Omega_n^{i,m}|$ is the number of pixels in $\Omega_n^{i,m}$. Similar to instance normalization [30], the activation is standardized per channel. However, ALIAS normalization divides the activation in channel k into the activation in the misaligned region and the other region.

The rationale behind this strategy is to remove the misleading information in the misaligned regions. Specifically, the misaligned regions in the warped clothing image match the background that is irrelevant to the clothing texture. Performing a standardization separately on these regions leads to a removal of the background information that causes the artifacts in the final results. In modulation, affine parameters inferred from the segmentation map modulate the standardized activation. Due to injecting semantic information at each ALIAS normalization layer, the layout of the human-parsing map in the final result is preserved.

ALIAS Generator. Fig. 4 describes the overview of the ALIAS generator, which adopts the simplified architecture that discards the encoder part of an encoder-decoder network. The generator employs a series of residual blocks (ResBlk) with upsampling layers. Each ALIAS ResBlk

consists of three convolutional layers and three ALIAS normalization layers. Due to the different resolutions that ResBlks operate at, we resize the inputs of the normalization layers, \hat{S} and $M_{misalign}$, before injecting them into each layer. Similarly, the input of the generator, $(I_a, P, \mathcal{W}(c, \theta))$, is resized to different resolutions. Before each ResBlk, the resized inputs $(I_a, P, \mathcal{W}(c, \theta))$ are concatenated to the activation of the previous layer after passing through a convolution layer, and each ResBlk utilizes the concatenated inputs to refine the activation. In this manner, the network performs the multi-scale refinement at a feature level that better preserves the clothing details than a single refinement at the pixel level. We train the ALIAS generator with the conditional adversarial loss, the feature matching loss, and the perceptual loss following SPADE [20] and pix2pixHD [32]. Details of the model architecture, hyperparameters, and the loss function are described in the supplementary.

4. Experiments

4.1. Experiment Setup

Dataset. We collected 1024×768 virtual try-on dataset for our research purpose, since the resolution of images on the dataset provided by Hanet *al.* [10] was low. Specifically, we crawled 13,679 frontal-view woman and top clothing image pairs on an online shopping mall website. The pairs are split into a training and a test set with 11,647 and 2,032 pairs, respectively. We use the pairs of a person and a clothing image to evaluate a paired setting, and we shuffle the clothing images for an unpaired setting. The paired setting is to reconstruct the person image with the original clothing item, and the unpaired setting is to change the clothing item on the person image with a different item.

Training and Inference. With the goal of reconstructing I from (I_a, c) , the training of each stage proceeds individually. During the training of the geometric matching module and the ALIAS generator, we use S instead of \hat{S} . While we aim to generate 1024×768 try-on images, we train the segmentation generator and the geometric matching module at 256×192. In the inference phase, after being predicted by the segmentation generator at 256×192, the segmentation map is upsampled to 1024×768 and passed

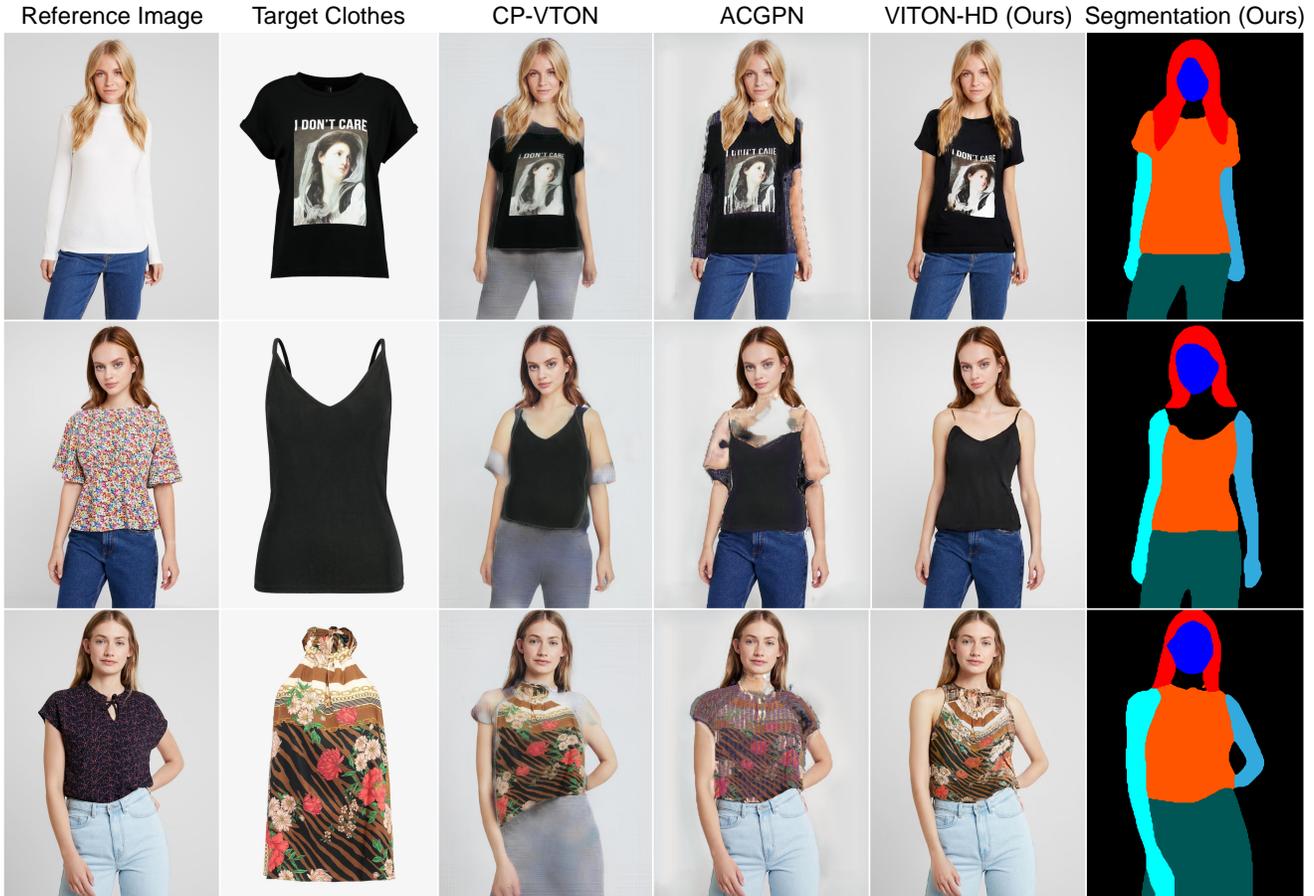


Figure 6: Qualitative comparison of the baselines.

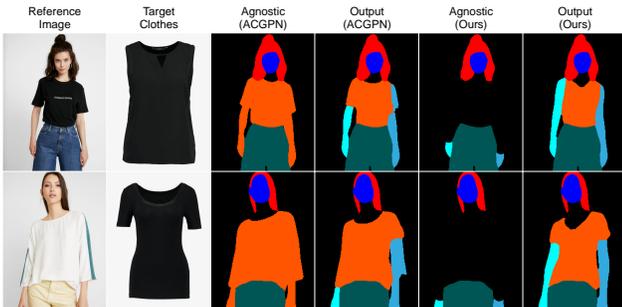


Figure 7: Qualitative comparison of the segmentation generator of ACGPN and VITON-HD. The clothing-agnostic segmentation map used by each model is also reported.

to subsequent stages. Similarly, the geometric matching module predicts the TPS parameters θ at 256×192 , and the 1024×768 clothing image deformed by the parameters θ is used in the ALIAS generator. We empirically found that this approach makes these two modules perform better with a lower memory cost than those trained at 1024×768 . Details of the model architecture and hyperparameters are described in the supplementary.

4.2. Qualitative Analysis

We compare VITON-HD with CP-VTON [31] and ACGPN [35], whose codes are publicly available. Following the training and inference procedure of our model, segmentation generators and geometric matching modules of the baselines are trained at 256×192 , and the outputs from the modules are upscaled to 1024×768 during the inference.

Comparison with Baselines. Fig. 6 demonstrates that VITON-HD generates more perceptually convincing 1024×768 images compared to the baselines. Our model clearly preserves the details of the target clothes, such as the logos and the clothing textures, due to the multi-scale refinement at a feature level. In addition, regardless of what clothes the person is wearing in the reference image, our model synthesizes the body shape naturally. As shown in Fig. 7, the shape of the original clothing item remains in the synthetic segmentation map generated by ACGPN. On the other hand, the segmentation generator in VITON-HD successfully predicts the segmentation map regardless of the original clothing item, due to our newly proposed clothing-agnostic person representation. Although our model surpasses the baselines qualitatively, there are a few limitations



Figure 8: Effects of ALIAS normalization. The orange colored areas in the enlarged images indicate the misaligned regions.

to VITON-HD, which are reported in the supplementary with the additional qualitative results.

Effectiveness of the ALIAS Normalization. We study the effectiveness of ALIAS normalization by comparing our model to VITON-HD*, where the standardization in ALIAS normalization is replaced by channel-wise standardization, as in the original instance normalization [30]. Fig. 8 shows that ALIAS normalization has the capability to fill the misaligned areas with the target clothing texture by removing the misleading information. On the other hand, without utilizing ALIAS normalization, the artifacts are produced in the misaligned areas, because the background information in the warped clothing image is not removed as described in Section 3.4. ALIAS normalization, however, can handle the misaligned regions properly.

4.3. Quantitative Analysis

We perform the quantitative experiments in both a paired and an unpaired settings, in which a person wears the original clothes or the new clothes, respectively. We evaluate our method using three metrics widely used in virtual try-on. The structural similarity (SSIM) [33] and the learned perceptual image patch similarity (LPIPS) [38] are used in the paired setting, and the frechet inception distance (FID) [11] score is adopted in the unpaired setting. The inception score [26] is not included in the experiments, since it cannot reflect whether the details of the clothing image are maintained [10]. The input of the each model contains different amount of information that offers advantages in reconstructing the segmentation maps, thus we use the segmentation maps from the test set instead of the synthetic segmentation maps in the paired setting.

Comparison across Different Resolutions. We compare the baselines quantitatively across different resolutions (256×192 , 512×384 , and 1024×768) as shown in Table 1. Our model outperforms the baselines for SSIM and LPIPS across all resolutions. For FID score, our model significantly surpasses CP-VTON, regardless of the resolutions. The FID score in ACGPN is slightly lower than that of our model at the 256×192 resolution. However, at the

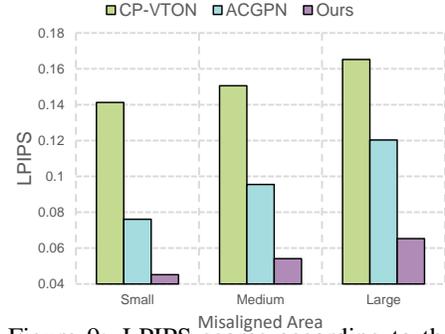


Figure 9: LPIPS scores according to the degree of misalignment.

1024×768 resolution, our model achieves a lower FID score than ACGPN with a large margin. The results indicate that the baselines cannot handle 1024×768 images, while our model is trained in a stable manner, even at a high resolution. This may be due to the limited capability of the U-Net architecture employed in the baseline models.

Comparison According to the Degree of Misalignment. To verify the ability of filling the misaligned areas with the clothing texture, we perform experiments in the paired setting according to the degree of the misalignment. According to the number of pixels in the misaligned areas, we divide the test dataset in three types: small, medium, and large. For a fair comparison, each model uses the same segmentation maps and the same warped clothes as inputs to match the misaligned regions. We evaluate LPIPS to measure the semantic distances between the reference images and the reconstructed images. As shown in Fig. 9, the wider the misaligned areas, the worse the performance of models, which means that the misalignment hinders the models from generating photo-realistic virtual try-on images. Compared to the baselines, our model consistently performs better, and the performance of our model decreases less as the degree of misalignment increases.

5. Conclusions

We propose the VITON-HD that synthesizes photo-realistic 1024×768 virtual try-on images. The proposed ALIAS normalization can properly handle the misaligned areas and propagate the semantic information throughout the ALIAS generator, which preserves the details of the clothes via the multi-scale refinement. Qualitative and quantitative experiments demonstrate that VITON-HD surpasses existing virtual try-on methods with a large margin.

Acknowledgments. This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korean government (MSIT) (No. NRF-2019R1A2C4070420) and Seoul R&BD Program (CD200024) through the Seoul Business Agency (SBA) funded by the Seoul Metropolitan Government.

References

- [1] Ivan Anokhin, Pavel Solovev, Denis Korzhenkov, Alexey Kharlamov, Taras Khakhulin, Aleksei Silvestrov, Sergey Nikolenko, Victor Lempitsky, and Gleb Sterkin. High-resolution daytime translation without domain labels. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 7488–7497, 2020. [3](#)
- [2] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. In *Proc. the International Conference on Learning Representations (ICLR)*, 2018. [3](#)
- [3] Z Cao, T Simon, SE Wei, YA Sheikh, et al. Openpose: Realtime multi-person 2d pose estimation using part affinity fields. *The IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2019. [4](#)
- [4] Harm De Vries, Florian Strub, Jérémie Mary, Hugo Larochelle, Olivier Pietquin, and Aaron C Courville. Modulating early visual processing by language. In *Proc. the Advances in Neural Information Processing Systems (NeurIPS)*, pages 6594–6604, 2017. [3](#)
- [5] Haoye Dong, Xiaodan Liang, Xiaohui Shen, Bochao Wang, Hanjiang Lai, Jia Zhu, Zhiting Hu, and Jian Yin. Towards multi-pose guided virtual try-on network. In *Proc. of the IEEE international conference on computer vision (ICCV)*, pages 9026–9035, 2019. [3](#)
- [6] Haoye Dong, Xiaodan Liang, Xiaohui Shen, Bowen Wu, Bing-Cheng Chen, and Jian Yin. Fw-gan: Flow-navigated warping gan for video virtual try-on. In *Proc. of the IEEE international conference on computer vision (ICCV)*, pages 1161–1170, 2019. [2](#)
- [7] Ke Gong, Xiaodan Liang, Yicheng Li, Yimin Chen, Ming Yang, and Liang Lin. Instance-level human parsing via part grouping network. In *Proc. of the European Conference on Computer Vision (ECCV)*, pages 770–785, 2018. [4](#)
- [8] Peng Guan, Loretta Reiss, David A Hirshberg, Alexander Weiss, and Michael J Black. Drape: Dressing any person. *ACM Transactions on Graphics (TOG)*, 31(4):1–10, 2012. [3](#)
- [9] Xintong Han, Xiaojun Hu, Weilin Huang, and Matthew R Scott. Clothflow: A flow-based model for clothed person generation. In *Proc. of the IEEE international conference on computer vision (ICCV)*, pages 10471–10480, 2019. [2](#), [3](#), [13](#)
- [10] Xintong Han, Zuxuan Wu, Zhe Wu, Ruichi Yu, and Larry S Davis. Viton: An image-based virtual try-on network. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 7543–7552, 2018. [2](#), [3](#), [6](#), [8](#)
- [11] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Proc. the Advances in Neural Information Processing Systems (NeurIPS)*, pages 6629–6640, 2017. [8](#)
- [12] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proc. of the IEEE international conference on computer vision (ICCV)*, pages 1501–1510, 2017. [3](#)
- [13] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proc. the International Conference on Machine Learning (ICML)*, pages 448–456, 2015. [3](#)
- [14] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 1125–1134, 2017. [3](#)
- [15] Nikolay Jetchev and Urs Bergmann. The conditional analogy gan: Swapping fashion articles on people images. In *Proc. of the IEEE international conference on computer vision workshop (ICCVW)*, pages 2287–2292, 2017. [3](#)
- [16] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014. [13](#)
- [17] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *Proc. of the IEEE international conference on computer vision (ICCV)*, pages 2794–2802, 2017. [12](#)
- [18] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *Proc. the International Conference on Learning Representations (ICLR)*, 2018. [11](#)
- [19] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier gans. In *Proc. the International Conference on Machine Learning (ICML)*, pages 2642–2651, 2017. [3](#)
- [20] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 2337–2346, 2019. [3](#), [6](#), [12](#)
- [21] Taesung Park, Jun-Yan Zhu, Oliver Wang, Jingwan Lu, Eli Shechtman, Alexei A. Efros, and Richard Zhang. Swapping autoencoder for deep image manipulation. In *Proc. the Advances in Neural Information Processing Systems (NeurIPS)*, 2020. [3](#)
- [22] Chaitanya Patel, Zhouyingcheng Liao, and Gerard Pons-Moll. Tailornet: Predicting clothing in 3d as a function of human pose, shape and garment style. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 7365–7375, 2020. [3](#)
- [23] Gerard Pons-Moll, Sergi Pujades, Sonny Hu, and Michael J Black. Clothcap: Seamless 4d clothing capture and retargeting. *ACM Transactions on Graphics (TOG)*, 36(4):1–15, 2017. [3](#)
- [24] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. In *Proc. the International Conference on Machine Learning (ICML)*, pages 1060–1069, 2016. [3](#)
- [25] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer Assisted Intervention*, pages 234–241, 2015. [2](#), [4](#), [11](#)
- [26] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques

- for training gans. In *Proc. the Advances in Neural Information Processing Systems (NeurIPS)*, pages 2234–2242, 2016. 8
- [27] Masahiro Sekine, Kaoru Sugita, Frank Perbet, Björn Stenger, and Masashi Nishiyama. Virtual fitting by single-shot body shape estimation. In *International Conference on 3D Body Scanning Technologies*, pages 406–413, 2014. 3
- [28] Wei Shen and Rujie Liu. Learning residual images for face attribute manipulation. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 4030–4038, 2017. 3
- [29] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 13
- [30] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016. 3, 6, 8
- [31] Bochao Wang, Huabin Zheng, Xiaodan Liang, Yimin Chen, Liang Lin, and Meng Yang. Toward characteristic-preserving image-based virtual try-on network. In *Proc. of the European Conference on Computer Vision (ECCV)*, pages 589–604, 2018. 2, 3, 4, 7, 13
- [32] 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 *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 8798–8807, 2018. 2, 3, 6, 11, 12
- [33] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004. 8
- [34] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, XiaoLei Huang, and Xiaodong He. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 1316–1324, 2018. 3
- [35] Han Yang, Ruimao Zhang, Xiaobao Guo, Wei Liu, Wangmeng Zuo, and Ping Luo. Towards photo-realistic virtual try-on by adaptively generating-preserving image content. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 7850–7859, 2020. 2, 3, 5, 7, 12, 13
- [36] Ruiyun Yu, Xiaoqi Wang, and Xiaohui Xie. Vtnfp: An image-based virtual try-on network with body and clothing feature preservation. In *Proc. of the IEEE international conference on computer vision (ICCV)*, pages 10511–10520, 2019. 2, 3
- [37] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *Proc. the International Conference on Machine Learning (ICML)*, pages 7354–7363, 2019. 13
- [38] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2018. 8
- [39] Peihao Zhu, Rameen Abdal, Yipeng Qin, and Peter Wonka. Sean: Image synthesis with semantic region-adaptive normalization. In *Proc. of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 5104–5113, 2020. 3

Supplementary Material

A. Implementation Details

A.1. Pre-processing Details

This section introduces the details of generating our clothing-agnostic person representation. To remove the dependency on the clothing item originally worn by a person, regions that can provide any original clothing information, such as the arms that hint at the sleeve length, should be eliminated. Therefore, when generating a clothing-agnostic image I_a , we remove the arms from the reference image I . For the same reason, legs should be removed if the pants are the target clothing items. We mask the regions with a gray color, so that the masked pixels of the normalized image would have a value of 0. We add padding to the masks to thoroughly remove these regions, and the width of the padding is empirically determined.

A.2. Model Architectures

This section introduces the architectures of the segmentation generator, the geometric matching module, and ALIAS generator in detail.

Segmentation Generator. The segmentation generator has the structure of U-Net [25], which consists of convolutional layers, downsampling layers, and upsampling layers. Two multi-scale discriminators [32] are employed for the conditional adversarial loss. The details of the segmentation generator architecture are shown in Fig. 10.

Geometric Matching Module. The geometric matching module consists of two feature extractors and a regression network. A correlation matrix is calculated from the two extracted features, and the regression network predicts the TPS parameter θ with the correlation matrix. The feature extractor is composed of a series of convolutional layers, and the regression network consists of a series of convolutional layers followed by a fully connected layer. The details are shown in Fig. 11.

ALIAS Generator. The architecture of the ALIAS generator consists of a series of ALIAS ResBlks with nearest-neighbor upsampling layers. We employ two multi-scale discriminators with instance normalization. Spectral normalization [18] is applied to all the convolutional layers. Note that we separately standardize the activation based on the misalignment mask $M_{misalign}$ only in the first five

ALIAS ResBlks. The details of the ALIAS generator architecture is shown in Fig. 12.

A.3. Training Details

This section introduces the losses and the hyperparameters for the segmentation generator, the geometric matching module, and the ALIAS generator.

Segmentation Generator. The segmentation generator G_S uses the clothing-agnostic segmentation map S_a , the pose map P , and the clothing item c as inputs ($\hat{S} = G_S(S_a, P, c)$) to predict the segmentation map \hat{S} of the person in the reference image wearing the target clothing item. The segmentation generator is trained with the

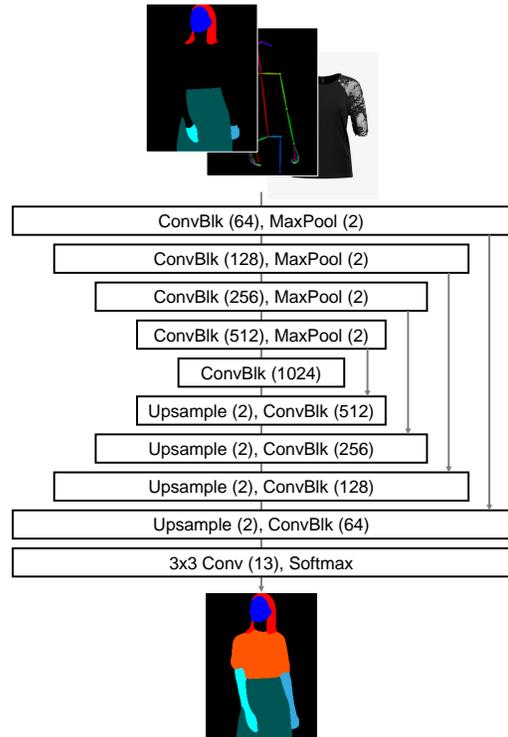


Figure 10: Segmentation Generator. $k \times k$ Conv (x) denotes a convolutional layer where the kernel size is k and the output channel is x . Also, ConvBlk (x) denotes a block, which consists of two series of 3×3 convolutional layer, instance normalization, and ReLU activation.

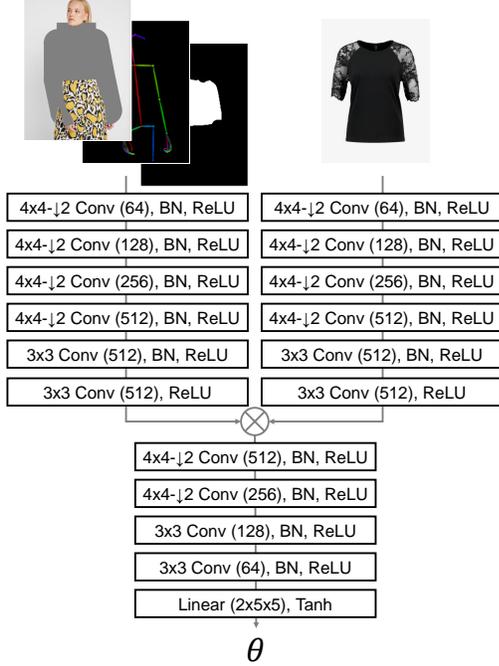


Figure 11: Geometric Matching Module. $k \times k \downarrow 2 \text{ Conv}(x)$ denotes a convolutional layer where the kernel size is k , the stride is 2, and the output channel is x .

cross-entropy loss \mathcal{L}_{CE} and the conditional adversarial loss \mathcal{L}_{cGAN} , which is LSGAN loss [17]. The full loss \mathcal{L}_S for the segmentation generator are written as

$$\mathcal{L}_S = \mathcal{L}_{cGAN} + \lambda_{CE} \mathcal{L}_{CE} \quad (8)$$

$$\mathcal{L}_{CE} = -\frac{1}{HW} \sum_{k \in C, y \in H, x \in W} S_{k,y,x} \log(\hat{S}_{k,y,x}) \quad (9)$$

$$\mathcal{L}_{cGAN} = \mathbb{E}_{(X,S)}[\log(D(X,S))] + \mathbb{E}_X[1 - \log(D(X,\hat{S}))], \quad (10)$$

where λ_{CE} is the hyperparameter for the cross-entropy loss. In the experiment, λ_{CE} is set to 10. In Eq. (9), $S_{y,x}$ and $\hat{S}_{y,x}$ indicate the pixel values of the segmentation map of the reference image S and \hat{S} corresponding to the coordinates (x, y) in channel k . The symbols H , W and C indicate the height, width, and the number of channels of S . In Eq. (10), the symbol X indicates the inputs of the generator (S_a, P, c) , and D denotes the discriminator.

The learning rate of the generator and the discriminator is 0.0004. We adopt the Adam optimizer with $\beta_1 = 0.5$ and $\beta_2 = 0.999$. We train the segmentation generator for 200,000 iterations with the batch size of 8.

Geometric Matching Module. The inputs of the geometric matching module are c , P , clothing-agnostic image I_a , and \hat{S}_c , which is the clothing area of \hat{S} . The output is

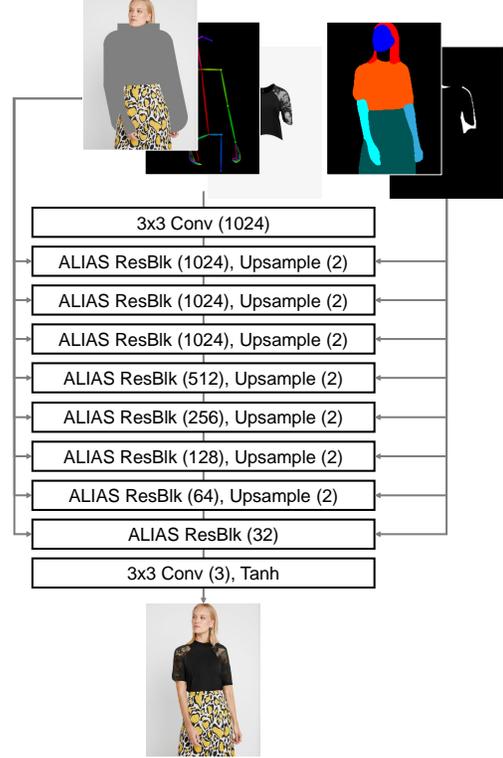


Figure 12: ALIAS Generator. The segmentation map S and the misalignment mask $M_{misalign}$ are passed to the generator through the proposed ALIAS ResBlks.

the TPS transformation parameters θ . The overall objective function is written as

$$\mathcal{L}_{warp} = \|I_c - \mathcal{W}(c, \theta)\|_{1,1} + \lambda_{const} \mathcal{L}_{const} \quad (11)$$

$$\mathcal{L}_{const} = \sum_{p \in \mathbf{P}} (|\|pp_0\|_2 - \|pp_1\|_2| + |\|pp_2\|_2 - \|pp_3\|_2|) + (|\mathcal{S}(p, p_0) - \mathcal{S}(p, p_1)| + |\mathcal{S}(p, p_2) - \mathcal{S}(p, p_3)|), \quad (12)$$

where \mathcal{W} is the function that deforms c using θ , and I_c is the clothing item extracted from the reference image I . \mathcal{L}_{const} is a second-order difference constraint [35], and λ_{const} is the hyperparameter for \mathcal{L}_{const} . In the experiment, we set λ_{const} to 0.04. In Eq. (12), the symbol p indicates a sampled TPS control point from the entire control points set \mathbf{P} , and p_0 , p_1 , p_2 , and p_3 are top, bottom, left and right point of p , respectively. The function $\mathcal{S}(p, p_i)$ denotes the slope between p and p_i .

The learning rate of the geometric matching module is 0.0002. We adopt the Adam optimizer with $\beta_1 = 0.5$ and $\beta_2 = 0.999$. We train the geometric matching module for 50,000 iterations with the batch size of 8.

ALIAS Generator. The loss function of ALIAS generator follows those of SPADE [20] and pix2pixHD [32], as it

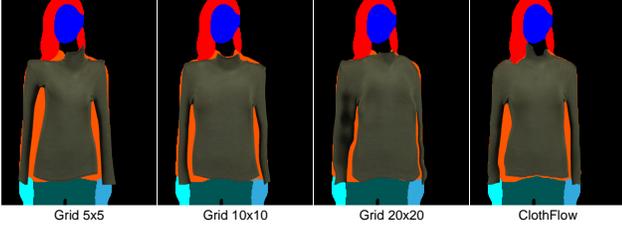


Figure 13: Qualitative comparisons of TPS transformation with various grid numbers and the flow estimation from ClothFlow.

Method	Warp-SSIM \uparrow	MACs \downarrow	Mask-SSIM \uparrow
ClothFlow	0.841*	8.13G	0.803*
VITON-HD	0.782	4.47G	0.852

Table 2: \star denotes a score taken from the ClothFlow paper, and we train VITON-HD in the same setting (e.g., dataset and resolution). We compute MACs of their warping modules at 256×192 .

contains the conditional adversarial loss \mathcal{L}_{cGAN} , the feature matching loss \mathcal{L}_{FM} , and the perceptual loss $\mathcal{L}_{percept}$. Let D_I be the discriminator, I and c be the given reference and target clothing images, and \hat{I} be the synthetic image generated by the generator. S_{div} is the modified version of the segmentation map S . The full loss \mathcal{L}_I of our generator is written as

$$\mathcal{L}_I = \mathcal{L}_{cGAN} + \lambda_{FM} \mathcal{L}_{FM} + \lambda_{percept} \mathcal{L}_{percept} \quad (13)$$

$$\begin{aligned} \mathcal{L}_{cGAN} = & \mathbb{E}_I[\log(D_I(S_{div}, I))] \\ & + \mathbb{E}_{(I,c)}[1 - \log(D_I(S_{div}, \hat{I}))] \end{aligned} \quad (14)$$

$$\mathcal{L}_{FM} = \mathbb{E}_{(I,c)} \sum_{i=1}^T \frac{1}{K_i} [\|D_I^{(i)}(S_{div}, I) - D_I^{(i)}(S_{div}, \hat{I})\|_{1,1}] \quad (15)$$

$$\mathcal{L}_{percept} = \mathbb{E}_{(I,c)} \sum_{i=1}^V \frac{1}{R_i} [\|F^{(i)}(I) - F^{(i)}(\hat{I})\|_{1,1}], \quad (16)$$

where λ_{FM} and $\lambda_{percept}$ are hyperparameters. In the experiment, both λ_{FM} and $\lambda_{percept}$ are set to 10. T is the number of layers in D_I , and $D_I^{(i)}$ and K_i are the activation and the number of elements in the i -th layer of D_I , respectively. Similarly, V is the number of layers used in the VGG network F [29], and $F^{(i)}$ and R_i are the activation and the number of elements in the i -th layer of F , respectively. We replace the standard adversarial loss with the Hinge loss [37].

The learning rate of the generator and the discriminator is 0.0001 and 0.0004, respectively. We adopt the Adam optimizer [16] with $\beta_1 = 0$ and $\beta_2 = 0.9$. We train the ALIAS generator for 200,000 iterations with the batch size of 4.

B. Additional Experiments

B.1. Comparison with ClothFlow

To demonstrate that the optical flow estimation does not solve the misalignment completely, we re-implement the flow estimation module of ClothFlow [9] based on the original paper. Fig. 13 shows that the misalignment still occurs, although both TPS with a higher grid number (e.g., a 10×10 or 20×20 grid) and the flow estimation module of ClothFlow can reduce the misaligned regions. The reason is that the regularization to avoid the artifacts (e.g., TV loss) prevents the warped clothes from fitting perfectly into the target region. In addition, we evaluate the accuracy and the computational cost of warping modules in VITON-HD and ClothFlow with Warp-SSIM [9] and MACs, respectively. We also measure how well the models reconstruct the clothing using Mask-SSIM [9]. Table 2 shows that the ClothFlow warping module has the better accuracy than ours, whereas the higher Mask-SSIM in VITON-HD proves that ALIAS normalization is more effective at solving the misalignment problem than the improved warping method. We found that the ClothFlow warping module needs a huge computational cost (MACs: 130.03G) at 1024×768 , but the cost could be reduced when predicting the optical flow map at 256×192 . Table 2 demonstrates that the ClothFlow warping module still needs more computational cost than ours, yet it is a viable option to combine the flow estimation module with ALIAS generator.

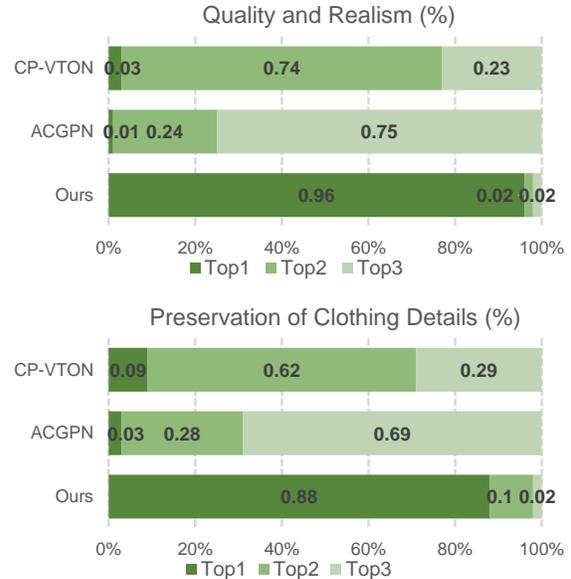


Figure 14: User study results. We compare our model with CP-VTON [31] and ACGPN [35].



Figure 15: Failure cases of VITON-HD.

B.2. User Study

We further evaluate our model and other baselines via a user study in the unpaired setting. We randomly select 30 sets of a reference image and a target clothing image from the test dataset. Given the reference images and the target clothes, the users are asked to rank the 1024×768 outputs of our model and baselines according to the following questions: (1) Which image is the most photo-realistic? (2) Which image preserves the details of the target clothing the most? As shown in Fig. 14, it can be observed that our approach achieves the rank 1 votes more than 88% for the both questions. The result demonstrates that our model generates more realistic images, and preserves the details of the clothing items compared to the baselines.

B.3. Qualitative Results

We provide additional qualitative results to demonstrate our model’s capability of handling high quality image synthesis. Fig. 16, 17, 18, and 19 show the qualitative comparison of the baselines across different resolutions. Fig. 20, 21, 22, and 23 show additional results of VITON-HD at 1024×768 resolution.

C. Failure Cases and Limitations

Fig. 15 shows the failure cases of our model caused by the inaccurately predicted segmentation map or the inner collar region indistinguishable from the other clothing region. Also, the boundaries of the clothing textures occasionally fade away.

The limitations of our model are as follows. VITON-HD is trained to preserve the bottom clothing items, limiting the presentation of the target clothes (*e.g.*, whether they are tucked in). It can be a valuable future direction to generate multiple possible outputs from a single input pair. Next, our dataset mostly consists of slim women and top clothing images, which makes VITON-HD handle only a limited range of body shapes and clothing during the inference. We believe that VITON-HD has the capability to cover more diverse cases when the images of various body shapes and clothing types are provided. Finally, existing virtual try-on methods including VITON-HD do not provide robust performance for in-the-wild images. We think generating realistic try-on images for the in-the-wild images is an interesting topic for future work.

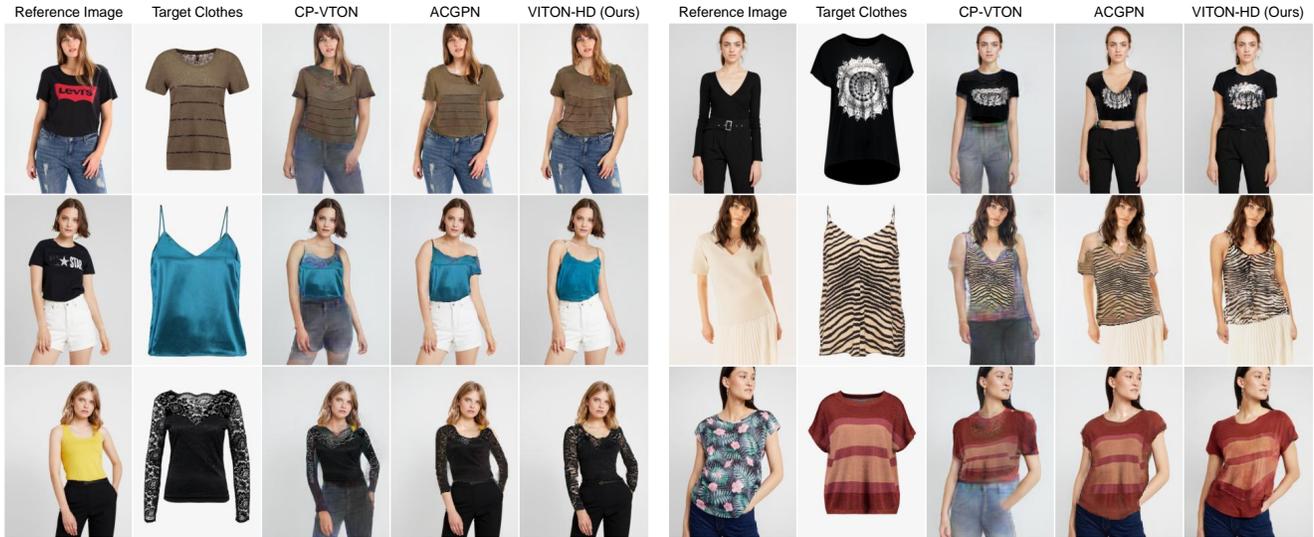


Figure 16: Qualitative comparison of the baselines (256×192).

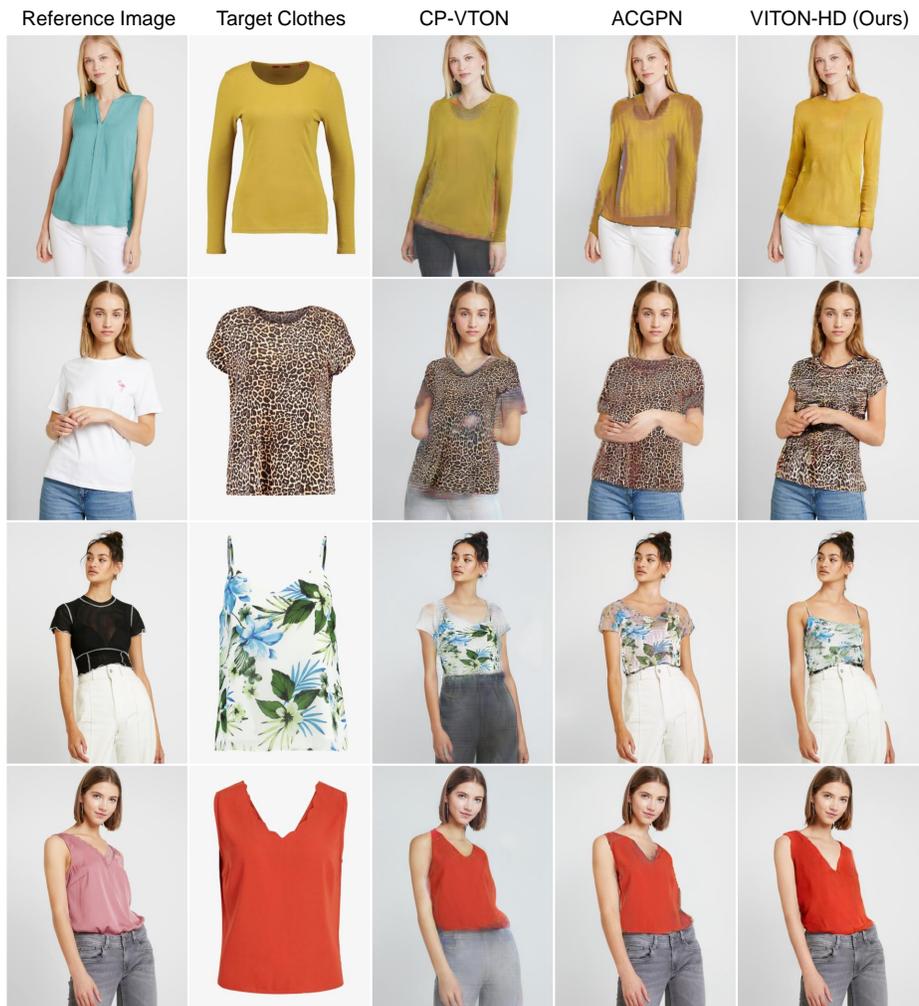


Figure 17: Qualitative comparison of the baselines (512×384).



Figure 18: Qualitative comparison of the baselines (1024×768).



Figure 19: Qualitative comparison of the baselines (1024×768).

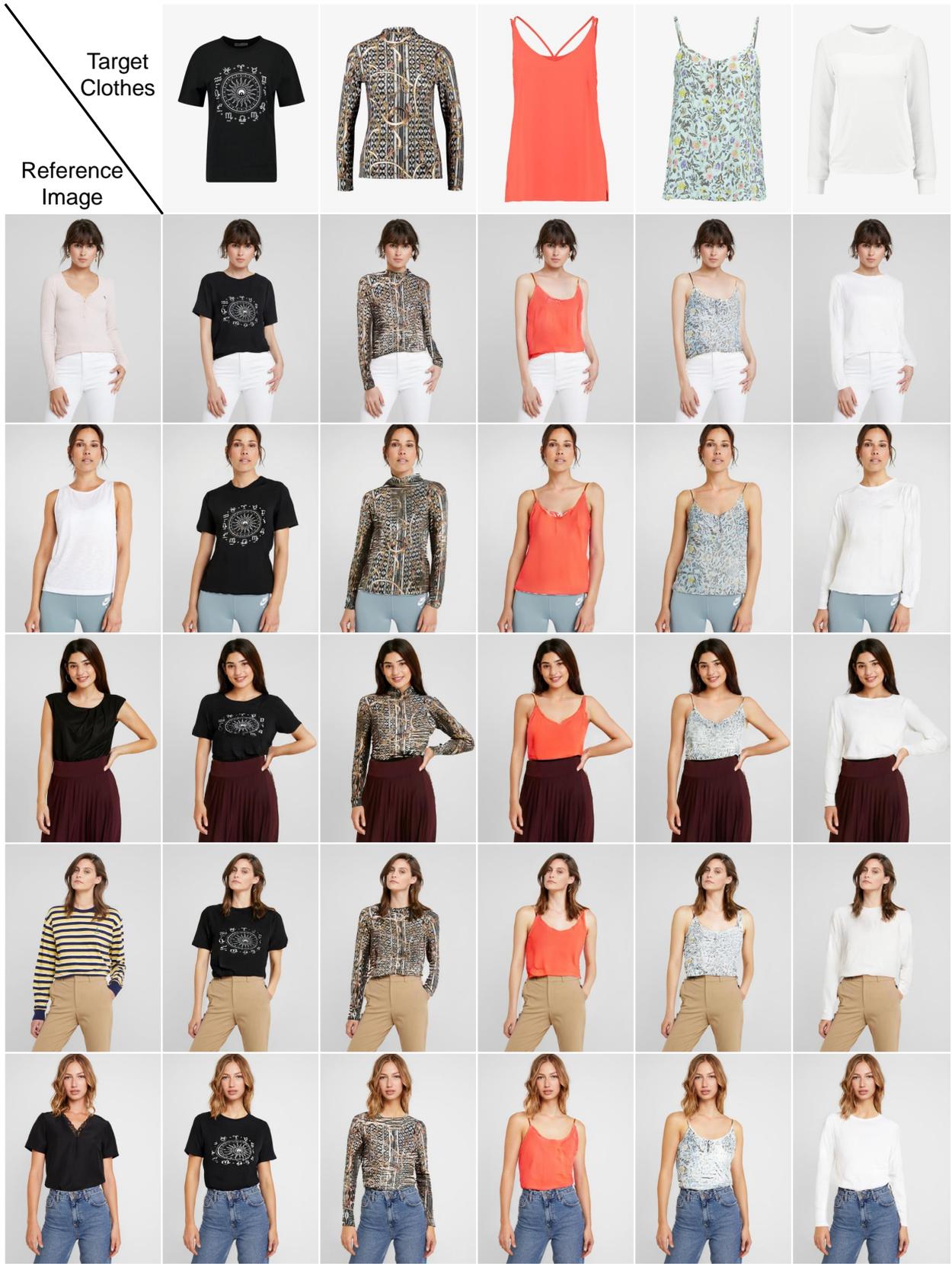


Figure 20: Additional qualitative results of VITON-HD.



Figure 21: Sample 1 of VITON-HD. (Left) The synthetic image. (Right) The reference image and the target clothing item.



Figure 22: Sample 2 of VITON-HD. (Left) The synthetic image. (Right) The reference image and the target clothing item.



Figure 23: Sample 3 of VITON-HD. (Left) The synthetic image. (Right) The reference image and the target clothing item.