# StableVITON: Learning Semantic Correspondence
# with Latent Diffusion Model for Virtual Try-On

Jeongho Kim    Gyojung Gu    Minho Park    Sunghyun Park    Jaegul Choo

KAIST, Daejeon, South Korea

{rlawjdghek, gyojung.gu, m.park, psh01087, jchoo}@kaist.ac.kr

Figure 1. Generated results of StableVITON: VITON-HD (the first row), SHHQ-1.0 (the first two images in the second row), and web-crawled images (the last two images in the second row). All results are generated using StableVITON trained on VITON-HD dataset.

## Abstract

*Given a clothing image and a person image, an image-based virtual try-on aims to generate a customized image that appears natural and accurately reflects the characteristics of the clothing image. In this work, we aim to expand the applicability of the pre-trained diffusion model so that it can be utilized independently for the virtual try-on task. The main challenge is to preserve the clothing details while effectively utilizing the robust generative capability of the pre-trained model. In order to tackle these issues, we propose StableVITON, learning the semantic correspondence between the clothing and the human body within the latent space of the pre-trained diffusion model in an end-to-end manner. Our proposed zero cross-attention blocks not only preserve the clothing details by learning the semantic correspondence but also generate high-fidelity images by utilizing the inherent knowledge of the pre-trained model in the warping process. Through our proposed novel attention total variation loss and applying augmentation, we achieve the sharp attention map, resulting in a more precise representation of clothing details. StableVITON outperforms the baselines in qualitative and quantitative evaluation, showing promising quality in arbitrary person images. Our code is available at https://github.com/rlawjdghek/StableVITON.*

# 1. Introduction

The objective of an image-based virtual try-on is to dress a given clothing image on a target person image. Most of the previous virtual try-on approaches [3, 7, 10, 14, 15, 30, 32, 34] leverage paired datasets consisting of clothing images and person images wearing those garments for training purposes. These methods typically include two modules: (1) a warping network to learn the semantic correspondence between the clothing and the human body, and (2) a generator that fuses the warped clothing and the person image.

Despite achieving significant advancements, previous methods [3, 8, 15, 32] still have limitations in achieving generalizability, particularly in maintaining the complex background in an arbitrary person image. The nature of matching clothing and individuals in the virtual try-on dataset [3, 10, 19] makes it challenging to collect data in diverse environments [21], which in turn leads to limitations in the generator's generative capability.

Meanwhile, recent advancements in large-scale pretrained diffusion models [24, 25, 28] have led to the emergence of downstream tasks [6, 8, 16, 20, 27, 35, 37] that control the pre-trained diffusion models for task-specific image generation. Thanks to the powerful generative ability, several works [18, 35] have succeeded in synthesizing high-fidelity human images using the prior knowledge of the pre-trained models, which signifies the potential for extension to the virtual try-on task.

In this paper, we aim to expand the applicability of the pre-trained diffusion model to provide a standalone model for the virtual try-on task. In the effort to adapt the pre-trained diffusion model for virtual try-on, a significant challenge is to preserve the clothing details while harnessing the knowledge of the pre-trained diffusion model. This can be achieved by learning the semantic correspondence between clothing and the human body using the provided dataset. Recent research [8, 20] that has employed pretrained diffusion models in virtual try-on has shown limitations due to the following two issues: (1) insufficient spatial information available for learning the semantic correspondence [20], and (2) the pre-trained diffusion model not being fully utilized, as it pastes the warped clothing in the RGB space, relying on external warping networks as previous approaches [3, 7, 15, 32, 34] for aligning the input condition.

To overcome these issues, we propose StableVITON, which learns the semantic correspondence between the clothing and the human body within the latent space of the pre-trained diffusion model. To incorporate the spatial information of the clothing for learning semantic correspondence, we introduce an encoder [35] that takes clothing as input and conditions the U-Net with the intermediate features of the encoder via zero cross-attention blocks. Warping through the zero cross-attention block in a pre-trained



Figure 2. Visualization of the semantic correspondence learned by our StableVITON. We overlay the attention map for the clothing regions onto the generated images for visualization.

diffusion model has the following two advantages: (1) preserving the clothing details by learning the semantic correspondence; (2) synthesizing high-fidelity images by leveraging the pre-trained models' inherent knowledge about humans in the warping process. As shown in Fig. 2, the attention mechanism in the latent space performs patch-wise warping by activating each token corresponding to clothing alignment within the generation region.

To further sharpen attention maps, we propose a novel attention total variation loss and apply the augmentation, which yields improved preservation of clothing details. By not impairing the pre-trained diffusion model, this architecture generates high-quality images even when images with complex backgrounds are provided, only using an existing virtual try-on dataset. Our extensive experiments show that StableVITON outperforms the existing virtual try-on method by a large margin. In summary, our contributions are as follows:

- Our proposed StableVITON, to the best of our knowledge, is the first end-to-end virtual try-on method finetuned on the pre-trained diffusion model without an independent warping process.
- We propose a zero cross-attention block, which learns semantic correspondence between the clothing and the human body, to condition the intermediate features from the spatial encoder.
- We propose a novel attention total variation loss and apply augmentation for further precise semantic correspondence learning.
- StableVITON shows state-of-the-art performance over existing virtual try-on models in both qualitative and quantitative results. Moreover, through the evaluation of a trained model on multiple datasets, StableVITON demonstrates its promising quality in a real-world setting.

## 2. Related Work

**GAN-based Virtual Try-On.** To properly try-on the given clothing image to the target person, existing approaches [3, 7, 15, 32] based on generative adversarial network (GAN) have attempted to address the virtual try-on problem using a two-stage strategy: (1) deforming the clothing to the proposal region and (2) fusing the warped clothing via try-on generator based on GAN. In order to achieve precise clothing deformation, previous methods [1, 7, 11, 15, 32] leverage a trainable network that estimates a dense flow map [38] to deform the clothing to the human body. At the same time, several approaches [3, 7, 14, 15, 32, 34] have been attempted to address the misalignment between the warped clothing and the human body, such as using a normalization [3] or distillation [7, 14]. However, the existing approaches still are not generalized well, leading to significant performance degradation in arbitrary person images with complex backgrounds. In this paper, we effectively address such issues by proposing a method that leverages the powerful generation ability of the pre-trained model.

**Diffusion-based Virtual Try-On.** Due to the remarkable generative capabilities, research on virtual try-on has extensively discussed the application of the diffusion models. While TryOnDiffusion [39] introduces an architecture for try-on using two U-Nets, this method requires a large-scale and challenging-to-collect dataset, consisting of image pairs of the same person wearing the same clothing in two different poses. Therefore, much recent research has shifted their focus towards using the prior of a large-scale pre-trained diffusion models [13, 23, 25, 33] in the virtual try-on task. LADI-VTON [20] represents the clothing as pseudo-words, and DCI-VTON [8] applies a warping network to input the clothing as conditions for the pre-trained diffusion models. While both models deal with background-related issues, they suffer from preserving high-frequency details due to the excessive loss of spatial information from the CLIP encoder [20] and drawbacks such as incorrectly warped clothing inherited from the independent warping network [8]. On the other hand, we propose to condition the intermediate feature maps of a spatial encoder through zero cross-attention block, which allows for using the prior knowledge of the pre-trained model in the warping process.

## 3. Preliminary

**Stable Diffusion Model.** Stable Diffusion model [25] is a large-scale diffusion model trained on LAION dataset [29], built upon the Latent Diffusion model (LDM) [25], which performs a denoising process in the latent space of an autoencoder. With a fixed encoder ($\mathcal{E}$), an input image $\mathbf{x}$ is first transformed to latent feature $\mathbf{z}_0 = \mathcal{E}(\mathbf{x})$. Given a predefined variance schedule $\beta_t$, we can define a forward diffusion process in the latent space following denoising diffu-

sion probabilistic models [13]:

$$q(\mathbf{z}_t|\mathbf{z}_0) = \mathcal{N}(\mathbf{z}_t; \sqrt{\bar{\alpha}_t}\mathbf{z}_0, (1 - \bar{\alpha}_t)\mathbf{I}), \qquad (1)$$

where $t \in \{1, ..., T\}$, $T$ represents the number of steps in the forward diffusion process, $\alpha_t := 1 - \beta_t$, and $\bar{\alpha}_t := \Pi_{s=1}^{t}\alpha_s$. As a training loss, Stable Diffusion model employs the simplified objective function from LDM [15]:

$$\mathcal{L}_{LDM} = \mathbb{E}_{\mathcal{E}(\mathbf{x}),\mathbf{y},\epsilon\sim\mathcal{N}(0,1),t}\left[\|\epsilon - \epsilon_\theta(\mathbf{z}_t, t, \tau_\theta(\mathbf{y}))\|_2^2\right], \tag{2}$$

where the denoising network $\epsilon_\theta(\cdot)$ is implemented with a U-Net architecture [26] and $\tau_\theta(\cdot)$ is the CLIP [23] text encoder to condition the text prompt $y$.
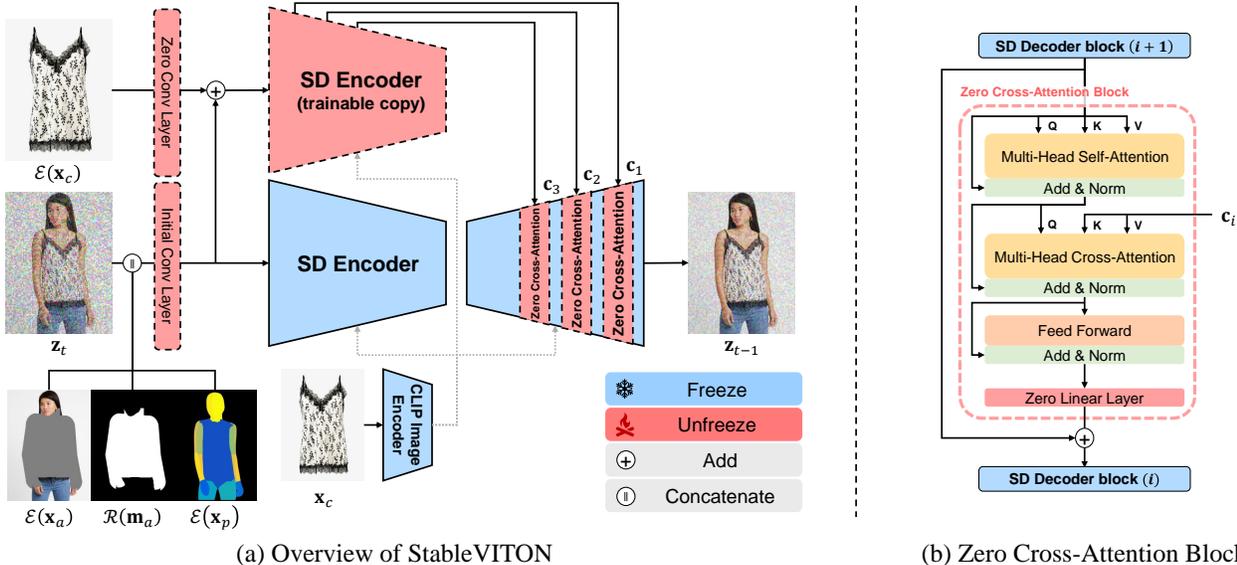
## 4. Method

### 4.1. Model Overview

An overview of the StableVITON is presented in Fig. 3(a). Given a person image $\mathbf{x} \in \mathbb{R}^{H \times W \times 3}$, the clothing-agnostic person representation $\mathbf{x}_a \in \mathbb{R}^{H \times W \times 3}$ (we call it as 'agnostic map') [3] is proposed to eliminate any clothing information in $\mathbf{x}$. In this work, we approach the virtual try-on as an exemplar-based image inpainting problem [33] to fill the agnostic map $\mathbf{x}_a$ with the clothing image $\mathbf{x}_c$. As the input of the U-Net, we concatenate four components: (1) the noisy image ($\mathbf{z}_t$), (2) latent agnostic map ($\mathcal{E}(\mathbf{x}_a)$), (3) the resized clothing-agnostic mask ($\mathbf{x}_{m_a}$), (4) latent dense pose condition ($\mathcal{E}(\mathbf{x}_p)$) [9] to preserve the person's pose. To align the input channels, we expand the initial convolution layer of the U-Net to 13 (*i.e.,* 4+4+1+4=13) channels with a convolution layer initialized with zero weights. For exemplar conditioning, we input the $\mathbf{x}_c$ to the CLIP image encoder [33].

To preserve the fine details of the clothing, we introduce a spatial encoder, which takes latent clothing ($\mathcal{E}(\mathbf{x}_c)$) as input. This spatial encoder copies the weight of the pre-trained U-Net [35] and conditions the intermediate feature maps of the encoder to U-Net via zero cross-attention blocks. During training, we apply augmentation and further finetune the model with our proposed attention total variation loss, which makes the attention region on the clothing sharper. The detailed model architecture is described in the supplementary material.

### 4.2. StableVITON

**Zero Cross-Attention Block.** We aim to condition the intermediate feature maps of the clothing to U-Net, properly aligning with the human body. The operation of adding the unaligned clothing feature map to the human feature map is insufficient to preserve clothing details due to misalignment between the human body and the clothing. Therefore, we proposed a zero cross-attention block to be a flexible operation by applying an attention mechanism for conditioning.

(a) Overview of StableVITON

(b) Zero Cross-Attention Block

Figure 3. For the virtual try-on task, StableVITON additionally takes three conditions: agnostic map, agnostic mask, and dense pose, as the input of the pre-trained U-Net, which serves as the query (Q) for the cross-attention. The feature map of the clothing is used as the key (K) and value (V) for the cross-attention and is conditioned on the UNet, as depicted in (b).

Specifically, as shown in Fig. 3(b), the feature map of the U-Net decoder block inputs to self-attention, followed by the cross-attention layer where the query (Q) comes from the previous self-attention layer and the spatial encoder's feature map serves as the key (K) and value (V). To eliminate harmful noise, we introduce a linear layer initialized with zero weight after the feed-forward operation [35].

To successfully align the clothing to the human body part via cross-attention, it is crucial to ensure semantic correspondence between the key tokens (clothing) and the query tokens (human body). For instance, when dealing with a query token related to the right shoulder, the corresponding key tokens should exhibit higher attention scores in the corresponding right shoulder area of the clothing. In Fig. 4(a), we averaged the attention maps of the resolution of $32 \times 24$ across the head dimension and arranged them flatly. For clear visualization, we downsample the generated image to a resolution of $32 \times 24$ and then resize it back to $32^2 \times 24^2$. Subsequently, we overlay this generated image with the attention map corresponding to each query token. Zooming in on the upper and middle of the generated clothing area, we observe that the key tokens unrelated to the corresponding query token, such as the bottom of the clothing, are activated in the attention map. This indicates that the cross-attention layer fails to learn the exact semantic correspondence between query and key tokens, combining the several key tokens of the clothing to generate the color corresponding to the query token during training. Therefore, as shown in Fig. 4, the stripes on the clothing are not distinctly visible.

**Augmentation.** To mitigate such issues of key tokens unre-

lated to query tokens being attended to, we alter the feature map by applying augmentation, including random shifts to input conditions. Detailed settings of augmentation are described in supplementary material. Along with the augmented input conditions, we train our model with the objective function defined as follows:

$$\mathcal{L}_{LDM} = \mathbb{E}_{\zeta, \mathbf{x}_c, \mathcal{E}(\mathbf{x}_c), \epsilon, t} \left[ \|\epsilon - \epsilon_\theta(\zeta, t, \tau_\phi(\mathbf{x}_c), \mathcal{E}(\mathbf{x}_c))\|_2^2 \right],$$
(3)

where $\zeta = [\mathbf{z}_t; \mathcal{E}(\mathbf{x}_a); \mathbf{x}_{m_a}; \mathcal{E}(\mathbf{x}_p)]$, and $\tau_\phi$ is the CLIP image encoder. Note that we do not update the parameters of the original blocks, as depicted in Fig. 3(a).

The rationale behind this strategy is to force the model to learn fine-grained semantic correspondence using augmentation, instead of just moderately injecting the clothes at similar positions. As shown in Fig. 4(b), we can confirm that key tokens related to query tokens have high attention scores, signifying that the cross-attention layer has learned the high semantic correspondence between the clothing-agnostic region and clothing.

**Attention Total Variation Loss.** While the cross-attention layer successfully aligns the clothing to the agnostic map, the points with high attention scores appear in dispersed positions, as shown in the attention map of Fig. 4(b). This causes inaccurate details in generated images, such as color discrepancies.

To address such an issue, we propose attention total variation loss. As the attention scores are the weight for the output, we calculate the center coordinates as a weighted sum of the attention map and the grid. Therefore, given the $H_q W_q$ query tokens and $h_k w_k$ key tokens, we calculate
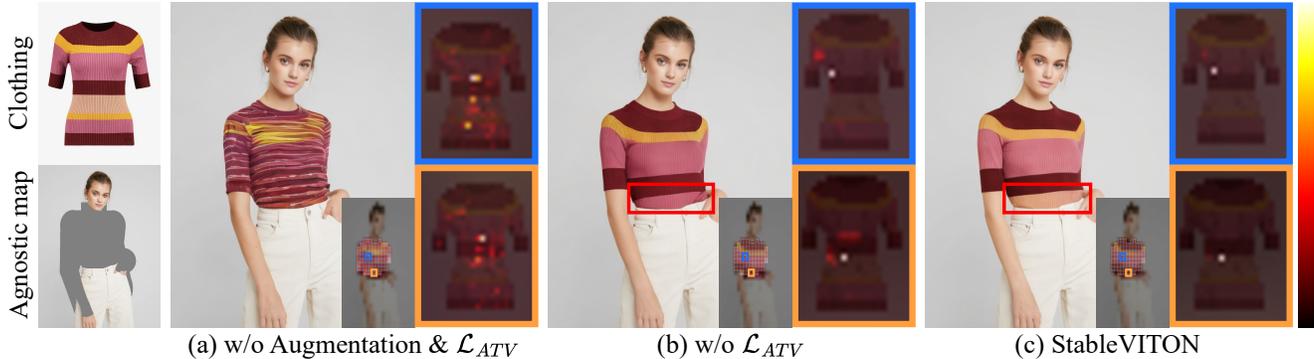
4

Figure 4. Visualization of attention map from a zero cross-attention block of 32 resolution.

center coordinate map $F \in \mathbb{R}^{H_q \times W_q \times 2}$ as follows:

$$F_{ijn} = \frac{1}{h_k w_k} \sum_{k=1}^{h_k} \sum_{l=1}^{w_k} (A_{ijkl} \odot G_{kln}), \qquad (4)$$

where we average the attention map over the head dimension and reshape it as $A \in \mathbb{R}^{H_q \times W_q \times h_k \times w_k}$, and $G \in [-1, 1]^{h_k \times w_k \times 2}$ is a 2D normalized coordinate. $\odot$ indicates element-wise multiplication operation.

For each query token in each clothing-agnostic region, the center coordinates should be evenly distributed, and the attention total variation loss $\mathcal{L}_{ATV}$ is defined as follows:

$$\mathcal{L}_{ATV} = \| \nabla(F \odot M) \|_1, \qquad (5)$$

where $M \in \{0, 1\}^{H_q \times W_q}$ is the ground truth clothing mask to only affect the clothing region. The attention total variation loss $\mathcal{L}_{ATV}$ is designed to enforce the center coordinates on the attention map uniformly distributed, thereby alleviating interference among attention scores located at dispersed positions. As illustrated in Figure (c), this leads to the generation of a sharper attention map, thereby more accurately reflecting the color of the clothing.

Finally, we finetune our StableVITON by adding $\mathcal{L}_{ATV}$ to Eq. 3:

$$\mathcal{L}_{finetune} = \mathcal{L}_{LDM} + \lambda_{ATV} \mathcal{L}_{ATV}, \qquad (6)$$

where $\lambda_{ATV}$ is a weight hyper-parameter.

## 5. Experiment

**Baselines.** We compare StableVITON with three GAN-based virtual try-on methods, VITON-HD [3], HR-VITON [15], and GP-VTON [32], and two diffusion-based virtual try-on methods, LADI-VTON [20] and DCI-VTON [8]. We also evaluate a diffusion-based inpainting method, Paint-by-Example [33]. We use pre-trained weights if available; otherwise, we train the models following the official code.

**Dataset.** We conduct the experiments using two publicly available virtual try-on datasets, VITON-HD [3] and Dress-Code [19], and one human image dataset, SHHQ-1.0 [5].

We train our model with VITON-HD and upper-body images in DressCode, respectively. For the evaluation of SHHQ-1.0, we use the first 2,032 images and follow the pre-processing instruction of VITON-HD [3] to obtain the input conditions such as the agnostic maps or the dense pose.

**Evaluation.** We evaluate the performances in two test settings. Specifically, the paired setting uses a pair of a person and the original clothes for reconstruction, whereas the unpaired setting involves changing the clothing of a person image with a different clothing item. As previous work [3], training and evaluation within a single dataset are referred to as 'single dataset evaluation'. On the other hand, we extend our evaluation on other datasets (*e.g.*, SHHQ-1.0), which we refer to as a 'cross dataset evaluation'. This evaluation enables an in-depth assessment of the model's generalizability in handling arbitrary person images, demonstrating applicability in real-world scenarios. Our model is capable of training at a $1024 \times 768$ resolution, but for a fair evaluation with baselines, we used a model trained at a $512 \times 384$ resolution. More results and details about experiments are described in the supplementary material.

### 5.1. Qualitative Results

**Single Dataset Evaluation.** As shown in Fig. 5, Stable-VITON generates realistic images and effectively preserves the text and clothing textures compared to the six baseline methods. Specifically, in the first row of Fig. 5, GAN-based methods such as GP-VTON struggle to generate the arms of the target person naturally. Moreover, other diffusion-based models either fail to preserve the text (Paint-by-Example and LADI-VTON) or show an overlapped artifact between the clothing and the target person (DCI-VTON). On the other hand, despite some parts of the arm being covered by clothing, our model produces a high-fidelity result that omits the 'L' in 'Love'.

**Cross Dataset Evaluation.** We visualize the generation images of the models trained on VITON-HD for DressCode and SHHQ-1.0 datasets in Fig. 6 and Fig. 7, respectively. The results clearly demonstrate that StableVITON gener-

5

Figure 5. Qualitative comparison with baselines in a single dataset setting (VITON-HD / VITON-HD). Best viewed when zoomed in.



Figure 6. Qualitative comparison with baselines in a cross dataset setting (VITON-HD / DressCode). Best viewed when zoomed in.

| Train / Test | VITON-HD / VITON-HD | | | | D.C. Upper / D.C. Upper | | | |
| Method | SSIM | LPIPS | FID | KID | SSIM | LPIPS | FID | KID |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| **VITON-HD [3]** | 0.862 | 0.117 | 12.117 | 3.23 | - | - | - | - |
| **HR-VITON [15]** | 0.878 | 0.1045 | 11.265 | 2.73 | <u>0.936</u> | 0.0652 | 13.820 | 2.71 |
| **LADI-VTON [20]** | 0.864 | 0.0964 | 9.480 | 1.99 | 0.915 | 0.0634 | 14.262 | 3.33 |
| **Paint-by-Example [33]** | 0.802 | 0.1428 | 11.939 | 3.85 | 0.897 | 0.0775 | 15.332 | 4.64 |
| **DCI-VTON [8]** | 0.880 | <u>0.0804</u> | 8.754 | 1.10 | **0.937** | <u>0.0421</u> | 11.920 | 1.89 |
| **GP-VTON [32]** | <u>0.884</u> | 0.0814 | 9.072 | <u>0.88</u> | 0.769 | 0.2679 | 20.110 | 8.17 |
| **Ours** | 0.852 | 0.0842 | <u>8.698</u> | <u>0.88</u> | 0.911 | 0.0500 | <u>11.266</u> | <u>0.72</u> |
| **Ours (RePaint [18])** | **0.888** | **0.0732** | **8.233** | **0.49** | **0.937** | **0.0388** | **9.940** | **0.12** |

Table 1. Quantitative comparisons in single dataset settings, VITON-HD and DressCode upper-body (D.C. Upper) datasets. **Bold** and <u>underline</u> denote the best and the second best result, respectively.

ates high-fidelity images while preserving the details of the clothing. GAN-based methods especially show significant artifacts on the target person and fail to maintain background. While diffusion-based methods generate natural images, they fail to preserve clothing details or the shape of the clothing. Furthermore, even when applying the augmentation we used to DCI-VTON (denoted as DCI-VTON (Aug.)), as depicted in Fig. 7, a significant improvement in the performance of the warping network is not achieved, failing to preserve clothing details.

## 5.2. Quantitative Results

**Metrics.** For quantitative evaluation, we use SSIM [31] and LPIPS [36] in the paired setting. In an unpaired setting, we assess the realism using FID [12] and KID [2] score. We follow the evaluation paradigm [20] for the implementation [4, 22].

**Single Dataset Evaluation.** We evaluate our StableVI-TON and existing baselines on a single dataset setting and report the results in Table 1. In the unpaired setting (*i.e.,* FID and KID), StableVITON outperforms all the baselines.

Figure 7. Qualitative comparison with baselines in a cross dataset setting (VITON-HD / SHHQ-1.0). Best viewed when zoomed in.

| Train / Test | VITON-HD / D.C. Upper | | | | D.C. Upper / VITON-HD | | | | VITON-HD / SHHQ-1.0 | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Method | SSIM | LPIPS | FID | KID | SSIM | LPIPS | FID | KID | FID | KID |
| VITON-HD [3] | 0.853 | 0.1874 | 44.257 | 28.82 | - | - | - | - | 71.149 | 52.01 |
| HR-VITON [15] | 0.909 | 0.1077 | 19.970 | 7.35 | 0.811 | 0.2278 | 45.923 | 36.69 | 52.732 | 31.22 |
| LADI-VTON [20] | 0.901 | 0.1009 | 16.336 | 5.36 | 0.801 | 0.2429 | 31.790 | 23.02 | 24.904 | 6.07 |
| Paint-by-Example [33] | 0.889 | 0.0867 | 16.398 | 4.78 | 0.784 | 0.1814 | 15.625 | 7.52 | 26.274 | 9.830 |
| DCI-VTON [8] | 0.903 | 0.1217 | 23.076 | 12.03 | <u>0.825</u> | 0.1870 | 16.670 | 6.40 | 24.850 | 6.68 |
| DCI-VTON (Aug.) [8] | 0.898 | 0.1240 | 18.809 | 8.02 | - | - | - | - | 24.368 | 6.11 |
| GP-VTON [32] | 0.724 | 0.3846 | 65.711 | 66.01 | 0.804 | 0.2621 | 52.351 | 48.68 | - | - |
| Ours | <u>0.911</u> | <u>0.0603</u> | <u>12.581</u> | <u>1.70</u> | 0.817 | <u>0.1308</u> | <u>10.104</u> | <u>1.72</u> | <u>23.531</u> | <u>5.68</u> |
| Ours (RePaint [18]) | **0.938** | **0.0470** | **10.480** | **0.41** | **0.855** | **0.1173** | **9.714** | **1.35** | **21.077** | **5.10** |

Table 2. Quantitative comparisons in cross dataset settings. We train the models on VITON-HD and DressCode upper-body (D.C. Upper) datasets and evaluate them on different datasets. **Bold** and <u>underline</u> denote the best and the second best result, respectively.

We observe that the performance degradation in the paired setting occurs due to the autoencoder's reconstruction error of the agnostic map. To mitigate this issue, we adapt RePaint [18], which samples the known region (*i.e.*, agnostic map) and replaces it in each denoising steps during the inference, used in DCI-VTON. Applying RePaint, StableVITON outperforms the baselines for all evaluation metrics. Since RePaint greatly helps maintain regions unrelated to clothing, it shows notable performance improvement in the paired setting. Nevertheless, even without RePaint, our method demonstrates superior performance in terms of FID and KID in the unpaired setting compared to the baselines.

**Cross Dataset Evaluation.** Table 2 presents that our StableVITON shows state-of-the-art performance for all the evaluation metrics with a large margin. GAN-based methods fail to maintain background consistency, resulting in significantly high FID and KID scores. While diffusion-based methods exhibit better performance due to the plausible generation outcomes enabled by pre-trained diffusion models, they fail to preserve clothing details. Consequently, they exhibit lower similarity scores in the paired setting (*i.e.,* SSIM and LPIPS) compared to our method.
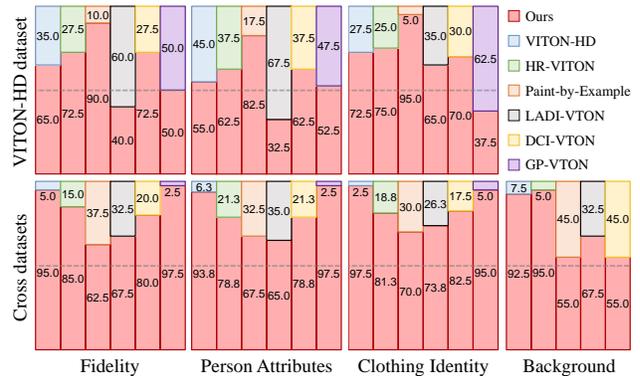


Figure 8. User study results. We compare our StableVITON with six baselines, involving a total of 40 participants.

## 5.3. User Study

For the models trained on the VITON-HD dataset, we conducted a user study with 40 participants. Each participant was shown one image generated by the baseline and the other by our model. They were asked to choose the better image based on three criteria: (1) fidelity, (2) person attributes, and (3) clothing identity. We added a question

about (4) background quality for the cross dataset setting. Detailed questions can be found in the supplementary material. As shown in Fig. 8, StableVITON outperforms in most of the criteria, and especially in the cross dataset setting, our method is overwhelming in human evaluations. While LADI-VTON shows better preference in the evaluation of fidelity and person attributes on the VITON-HD dataset, it fails to preserve clothing details, resulting in a 35% preference in the clothing identity criteria.

| Train / Test on VITON-HD | SSIM | LPIPS | FID | KID |
|---|---|---|---|---|
| ControlNet + Aug. | 0.832 | 0.1157 | 9.81 | 1.81 |
| ControlNet-W + Aug. | 0.822 | 0.1124 | 9.66 | 1.57 |
| Zero Cross-Attention Block + Aug. | **0.850** | **0.0851** | **8.74** | **0.91** |

Table 3. Quantitative comparison results between our zero cross-attention block and ControlNet. We train ControlNet with clothing and warped clothing [15] (ControlNet-W). We apply augmentation to all the models in training.

## 5.4. Ablation Study

**Comparison with ControlNet.** To demonstrate the effectiveness of StableVITON in tackling the alignment issue compared to ControlNet [35], we train ControlNet under two different input conditions: (1) clothing, and (2) warped clothing [15] (dubbed as ControlNet-W). We apply our proposed augmentation to both models during training. Using the warping network to align the clothing helps ControlNet capture coarse features, such as the overall shape and color of the logo, as shown in Fig. 9(b) and (c). However, since the misalignment still exists between warped clothing and the human body in the training phase, ControlNet-W struggles to reflect more fine-grained clothing details to the generation results. These observations highlight that ControlNet is highly sensitive to subtle misalignment across the input, stemming from the limitations of the ControlNet's direct addition operation in conditioning. In contrast, as shown in Fig. 9(d), the zero cross-attention block, free from the alignment constraints, successfully preserves the logos and patterns and leads to a qualitative performance improvement, as shown in Table 3.

**Effect of Training Components.** We investigate the effect of the two proposed components during training: augmentation and attention total variation loss. In Fig. 10, we visualize the generated images while incrementally introducing the proposed training components one by one. Compared to Fig 10(a) and (b), we observe that detailed features such as logos and patterns of the clothing are more preserved when augmentation is applied. However, as the attention maps are not sufficiently distinct, we observe inaccuracies in the generated images, such as the 'M' in 'PUMA' being incorrectly depicted or lines blurring, as shown in Fig. 10(b). After finetuning with our proposed attention total variation loss, these finer details are significantly improved. Such visual enhancements correspond to quantitative performance improvements as demonstrated in Table 4.



(a) Target Person (b) ControlNet (c) ControlNet-W (d) Ours

Figure 9. Comparison between our StableVITON and ControlNets under the two different input conditions: (1) clothing, and (2) warped clothing [15] (ControlNet-W). Best viewed when zoomed in.



(a) w/o Aug. & $\mathcal{L}_{ATV}$ (b) w/o $\mathcal{L}_{ATV}$ (c) Ours

Figure 10. Visual comparisons according to our training components. Best viewed when zoomed in.

| Aug. | Attention TV Loss | SSIM | LPIPS | FID | FID |
|---|---|---|---|---|---|
| ✗ | ✗ | 0.847 | 0.0969 | 9.35 | 1.33 |
| ✓ | ✗ | 0.850 | 0.0851 | 8.744 | 0.91 |
| ✓ | ✓ | **0.852** | **0.042** | **8.698** | **0.88** |

Table 4. Ablation study of our proposed training components on VITON-HD dataset.

## 6. Conclusion

We propose StableVITON, a novel image-based virtual try-on method using the pre-trained diffusion model. Our proposed zero cross-attention block learns semantic correspondence between the clothing and the human body, enabling try-on in the latent feature space. A novel attention total variation loss and augmentation are designed to preserve the clothing details better. Extensive experiments, including cross dataset evaluation, clearly demonstrate that StableVITON shows the state-of-the-art performance compared to the existing methods and its promising quality in the real-world setting.

# References

[1] Shuai Bai, Huiling Zhou, Zhikang Li, Chang Zhou, and Hongxia Yang. Single stage virtual try-on via deformable attention flows. In *European Conference on Computer Vision*, pages 409–425. Springer, 2022. 3

[2] Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. *arXiv preprint arXiv:1801.01401*, 2018. 6

[3] Seunghwan Choi, Sunghyun Park, Minsoo Lee, and Jaegul Choo. Viton-hd: High-resolution virtual try-on via misalignment-aware normalization. In *CVPR*, pages 14131–14140, 2021. 2, 3, 5, 6, 7, 11

[4] Nicki Skafte Detlefsen, Jiri Borovec, Justus Schock, Ananya Harsh Jha, Teddy Koker, Luca Di Liello, Daniel Stancl, Changsheng Quan, Maxim Grechkin, and William Falcon. Torchmetrics-measuring reproducibility in pytorch. *Journal of Open Source Software*, 7(70):4101, 2022. 6

[5] Jianglin Fu, Shikai Li, Yuming Jiang, Kwan-Yee Lin, Chen Qian, Chen-Change Loy, Wayne Wu, and Ziwei Liu. Stylegan-human: A data-centric odyssey of human generation. *arXiv preprint*, arXiv:2204.11823, 2022. 5

[6] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *European Conference on Computer Vision*, pages 89–106. Springer, 2022. 2

[7] Yuying Ge, Yibing Song, Ruimao Zhang, Chongjian Ge, Wei Liu, and Ping Luo. Parser-free virtual try-on via distilling appearance flows. In *CVPR*, pages 8485–8493, 2021. 2, 3

[8] Junhong Gou, Siyu Sun, Jianfu Zhang, Jianlou Si, Chen Qian, and Liqing Zhang. Taming the power of diffusion models for high-quality virtual try-on with appearance flow. *arXiv preprint arXiv:2308.06101*, 2023. 2, 3, 5, 6, 7

[9] Riza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation in the wild. In *CVPR*, pages 7297–7306, 2018. 3

[10] Xintong Han, Zuxuan Wu, Zhe Wu, Ruichi Yu, and Larry S Davis. Viton: An image-based virtual try-on network. In *CVPR*, pages 7543–7552, 2018. 2

[11] Xintong Han, Xiaojun Hu, Weilin Huang, and Matthew R Scott. Clothflow: A flow-based model for clothed person generation. In *ICCV*, pages 10471–10480, 2019. 3

[12] 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. 2017. 6

[13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 3

[14] Thibaut Issenhuth, Jérémie Mary, and Clément Calauzenes. Do not mask what you do not need to mask: a parser-free virtual try-on. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX 16*, pages 619–635. Springer, 2020. 2, 3

[15] Sangyun Lee, Gyojung Gu, Sunghyun Park, Seunghwan Choi, and Jaegul Choo. High-resolution virtual try-on with misalignment and occlusion-handled conditions. In *ECCV*, pages 204–219. Springer, 2022. 2, 3, 5, 6, 7, 8

[16] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22511–22521, 2023. 2

[17] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. *arXiv preprint arXiv:2202.09778*, 2022. 11

[18] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022. 2, 6, 7

[19] Davide Morelli, Matteo Fincato, Marcella Cornia, Federico Landi, Fabio Cesari, and Rita Cucchiara. Dress code: High-resolution multi-category virtual try-on. In *ECCV*, pages 2231–2235, 2022. 2, 5, 11

[20] Davide Morelli, Alberto Baldrati, Giuseppe Cartella, Marcella Cornia, Marco Bertini, and Rita Cucchiara. Ladi-vton: Latent diffusion textual-inversion enhanced virtual try-on. *arXiv preprint arXiv:2305.13501*, 2023. 2, 3, 5, 6, 7

[21] Assaf Neuberger, Eran Borenstein, Bar Hilleli, Eduard Oks, and Sharon Alpert. Image based virtual try-on network from unpaired data. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5184–5193, 2020. 2

[22] Gaurav Parmar, Richard Zhang, and Jun-Yan Zhu. On aliased resizing and surprising subtleties in gan evaluation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11410–11420, 2022. 6

[23] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 3

[24] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 2

[25] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2, 3, 11

[26] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015. 3

[27] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023. 2

[28] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 2

[29] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022. 3

[30] Bochao Wang, Huabin Zheng, Xiaodan Liang, Yimin Chen, Liang Lin, and Meng Yang. Toward characteristic-preserving image-based virtual try-on network. In *ECCV*, pages 589–604, 2018. 2

[31] 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. 6

[32] Zhenyu Xie, Zaiyu Huang, Xin Dong, Fuwei Zhao, Haoye Dong, Xijin Zhang, Feida Zhu, and Xiaodan Liang. Gp-vton: Towards general purpose virtual try-on via collaborative local-flow global-parsing learning. In *CVPR*, pages 23550–23559, 2023. 2, 3, 5, 6, 7

[33] Binxin Yang, Shuyang Gu, Bo Zhang, Ting Zhang, Xuejin Chen, Xiaoyan Sun, Dong Chen, and Fang Wen. Paint by example: Exemplar-based image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18381–18391, 2023. 3, 5, 6, 7, 11

[34] 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 *CVPR*, pages 7850–7859, 2020. 2, 3

[35] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 2, 3, 4, 8

[36] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 6

[37] Yuxin Zhang, Nisha Huang, Fan Tang, Haibin Huang, Chongyang Ma, Weiming Dong, and Changsheng Xu. Inversion-based style transfer with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10146–10156, 2023. 2

[38] Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Malik, and Alexei A Efros. View synthesis by appearance flow. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 286–301. Springer, 2016. 3

[39] Luyang Zhu, Dawei Yang, Tyler Zhu, Fitsum Reda, William Chan, Chitwan Saharia, Mohammad Norouzi, and Ira Kemelmacher-Shlizerman. Tryondiffusion: A tale of two unets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4606–4615, 2023. 3

# StableVITON: Learning Semantic Correspondence with Latent Diffusion Model for Virtual Try-On

## Supplementary Material

## A. Implementation details

**Architecture Details.** We adopt the autoencoder and the denoising U-Net of the Stable Diffusion v1.4 [25]. We initialize our denoising U-Net with the weights of the U-Net from the Paint-by-Example [33]. The U-Net's encoder and decoder both consist of 12 blocks, involving three downsampling and upsampling steps each. As a result, when StableVITON receives the input image of $64 \times 48$ resolution (after going through the autoencoder), three intermediate feature maps are generated for each of the following resolutions: $8 \times 6$, $16 \times 12$, $32 \times 24$, $64 \times 48$. We use the feature maps at resolutions other than $8 \times 6$ as inputs for each of the nine zero cross-attention blocks. Similarly, for the spatial encoder following the U-Net's encoder structure, we utilized the nine feature maps at resolutions other than $8 \times 6$ as the key and value inputs for the cross-attention layers.

**Training & Inference Details.** We train the model using an AdamW optimizer with a fixed learning rate of 1e-4 for 360k iterations, employing a batch size of 32. Then, we finetune the model with the attention total variation weight hyper-parameter $\lambda_{ATV} = 0.001$, using the same learning rate and batch size for 36K iterations. We train for about 100 hours using four NVIDIA A100 GPUs. For augmentation, we simultaneously applied horizontal flip (p=0.5) to both the clothing and the UNet's input condition, and independently applied Random Shift (limit=0.2, p=0.5) and Random Scale (limit=0.2, p=0.5) to both the clothing and UNet's input. We simultaneously applied HSV adjustments (limit=5, p=0.5) and contrast adjustments (limit=0.3, p=0.5) to both the clothing and $\mathbf{x}_0$. To prevent the issue of facial distortion, we finetune the decoder of the autoencoder separately on the training datasets VITON-HD [3] and Dress-Code [19]. In training the decoder, we use the AdamW optimizer with a learning rate of 5e-5 and a batch size of 32 for 10k iterations on each dataset. For inference, we employ the pseudo linear multi-step (PLMS) [17] sampler, with the number of sampling steps set to 50.

## B. User Study Details.

In the user study, participants were asked to evaluate which of the two images, one generated by the baseline and the other by StableVITON, was superior in terms of 1) fidelity, 2) person attributes, 3) clothing identity, and 4) background quality (with respect to the cross-dataset setting). The questions for each criterion were as follows:

- Fidelity: Choose which image better exhibits resemblance to reality in aspects such as the human body and color harmony.
- Person Attributes: Choose which image better maintains features like skin tone, pose, and appearance from the input image.
- Clothing Identity: Choose which image better preserves characteristics such as the design, logo, and shape of the input clothing.
- Background Quality: Choose which image better maintains the background of the input image.

## C. Additional Qualitative Results

In Fig. 11, we present the generation results on VITON-HD dataset, using the models trained on DressCode upper body dataset. GAN-based models (*i.e.*, HR-VITON and GP-VTON) show significant artifacts around the target person, as evidenced in quantitative results (see Table 2 in the main paper), leading to high FID and KID scores. In addition, diffusion-based models, while providing a plausible appearance, fail to preserve the details of the clothing.

## D. StableVITON at High Resolution

To synthesize high-fidelity images, we have further trained StableVITON at the higher resolution of $1024 \times 768$. Instead of starting from scratch, we experimentally observed that progressively training StableVITON with $1024 \times 768$ resolution images leads to faster convergence. We conducted additional training for 85k iterations using the same training settings as StableVITON.

**Qualitative Results.** We present the results generated by StableVITON trained on VITON-HD dataset at $1024 \times 768$ resolution for images in VITON-HD, DressCode, SHHQ-1.0, and web-crawled datasets, in Fig. 13, Fig. 14, Fig. 15, and Fig. 16, respectively. We observe that there is a clearer preservation of facial or clothing details at $1024 \times 768$ resolution.

## E. Limitations & Discussion

Even when fine-tuning the decoder of the autoencoder on the virtual try-on dataset, it remains challenging to preserve very fine details of the face or clothing. As a result, as shown in Fig. 5 of the main paper, subtle variations in facial features, such as the eyes, can be observed. However, we effectively address these issues related to fine details by increasing the model's resolution, as demonstrated in Fig. 13.

In our experiments, we observed that our model fails to preserve objects occluding the person or accessories such as bracelets and watches attached to the target person. This issue arises from the model's inability to incorporate additional information, apart from clothing, during the sampling process to fill the masked regions of the agnostic map. We leave such preservation issues as future work.

Figure 11. Qualitative comparison with baselines in a cross dataset setting (DressCode / VITON-HD). Best viewed when zoomed in.

| Clothing | Person | HR-VITON | Paint-by-Example | LADI-VTON | DCI-VTON | GP-VTON | Ours |
|----------|--------|----------|------------------|-----------|----------|---------|------|



| Clothing & Person | StableVITON | Clothing & Person | StableVITON |
|-------------------|-------------|-------------------|-------------|

Figure 12. Limitations of StableVITON. StableVITON fails to preserve objects occluding the person or accessories such as bracelets.

Figure 13. The generation results for the VITON-HD test dataset by StableVITON trained on VITON-HD dataset at $1024 \times 768$ resolution. Best viewed when zoomed in.

Figure 14. The generation results for the DressCode dataset by StableVITON trained on VITON-HD dataset at $1024 \times 768$ resolution. Best viewed when zoomed in.
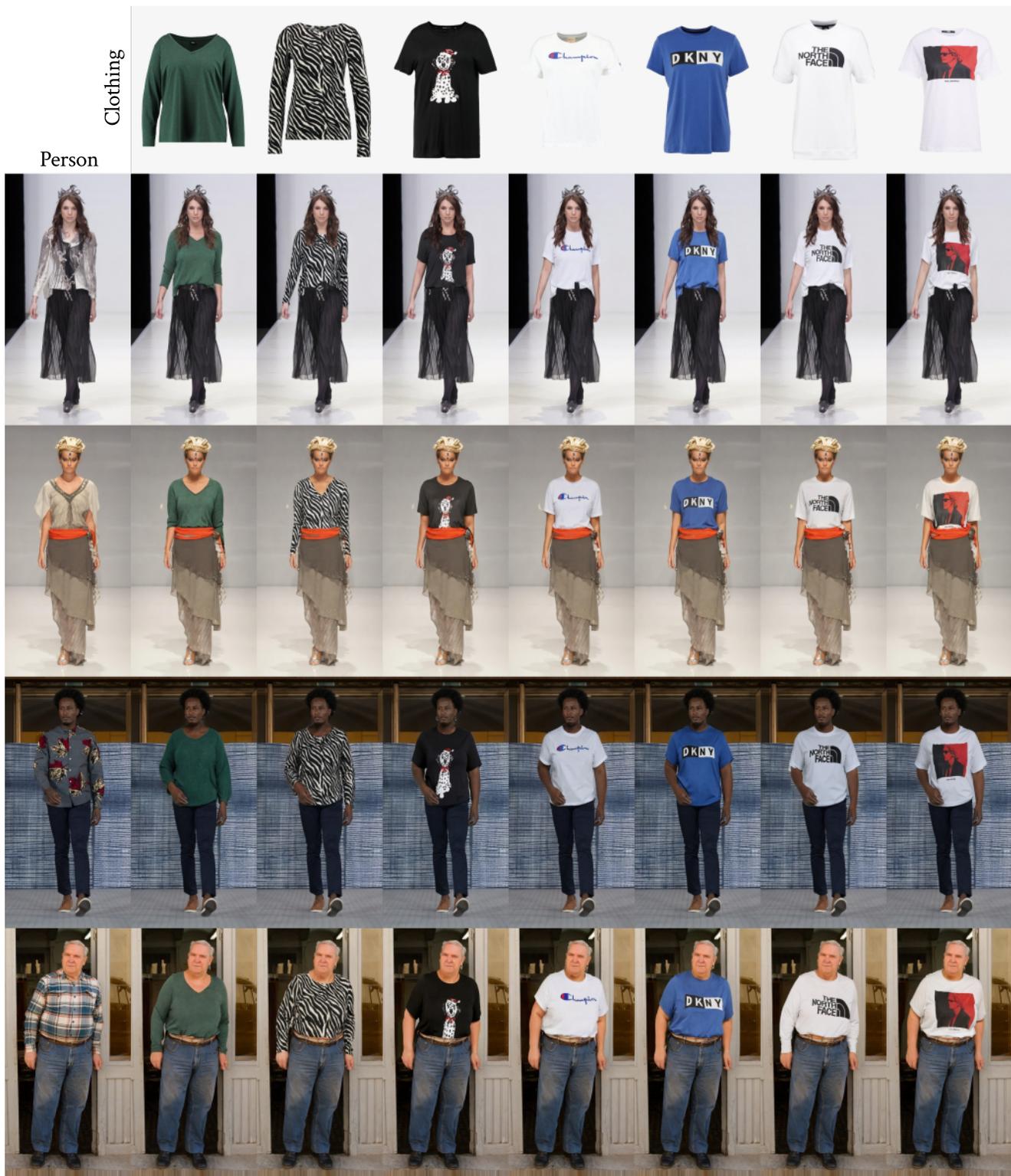
Figure 15. The generation results for the SHHQ-1.0 dataset by StableVITON trained on VITON-HD dataset at $1024 \times 768$ resolution. Best viewed when zoomed in.

Figure 16. The generation results for the web-crawled images by StableVITON trained on VITON-HD dataset at $1024 \times 768$ resolution. Best viewed when zoomed in.