

StableGarment: Garment-Centric Generation via Stable Diffusion

Rui Wang^{1†*}, Hailong Guo^{1†}, Jiaming Liu^{2†}, Huaxia Li², Haibo Zhao², Xu Tang², Yao Hu², Hao Tang³, and Peipei Li^{1‡}

¹ Beijing University of Posts and Telecommunications

² Xiaohongshu Inc.

³ Carnegie Mellon University

<https://raywang335.github.io/stablegarment.github.io/>



Fig. 1: The proposed StableGarment can perform various garment-centric generation tasks. Given a garment input, it could 1) utilize text prompts or control signals to generate a realistic model wearing the garment, 2) support switching stylized models to generate stylized models wearing the garment, and 3) conventional virtual try-on tasks. While performing those tasks, the details of the garment could be well preserved and the garments warping are visually natural.

Abstract. In this paper, we introduce StableGarment, a unified framework to tackle garment-centric(GC) generation tasks, including GC text-to-image, controllable GC text-to-image, stylized GC text-to-image, and

*Work done during internship at Xiaohongshu Inc.

†Equally contributed.

‡Corresponding author.

robust virtual try-on. The main challenge lies in retaining the intricate textures of the garment while maintaining the flexibility of pre-trained Stable Diffusion. Our solution involves the development of a garment encoder, a trainable copy of the denoising UNet equipped with additive self-attention (ASA) layers. These ASA layers are specifically devised to transfer detailed garment textures, also facilitating the integration of stylized base models for the creation of stylized images. Furthermore, the incorporation of a dedicated try-on ControlNet enables StableGarment to execute virtual try-on tasks with precision. We also build a novel data engine that produces high-quality synthesized data to preserve the model’s ability to follow prompts. Extensive experiments demonstrate that our approach delivers state-of-the-art (SOTA) results among existing virtual try-on methods and exhibits high flexibility with broad potential applications in various garment-centric image generation.

Keywords: Garment-centric generation, Diffusion models, Virtual try-on

1 Introduction

Recent advances in image generation have seen transformative developments, particularly with the emergence of text-to-image diffusion models trained on large datasets. Among these, Stable Diffusion, an open-source model referenced as [43], stands out for democratizing image generation from textual prompts for a wide user base. Such progress significantly impacts various application domains, notably within the fashion industry, which commands a considerable market presence.

In the realm of fashion, virtual try-on is a classic task that aims to superimpose given garments onto specific user images [13, 20, 29, 41, 56, 62, 62]. The development of diffusion models offers new levels of photorealism in generated images that were previously unattainable with Generative Adversarial Network (GAN)-based methods. Diffusion-based models not only achieve levels of realism previously deemed unattainable, but also excel in restoring intricate details and ensuring the images retain a natural appearance.

However, when extended beyond conventional virtual try-on tasks, existing methods face notable limitations. Garment merchants are in pursuit of creating varied product visuals, such as posters and display images, more cost-effectively. There is a dual demand: the ability for quick adjustments to models, poses, atmospheres, and backgrounds through textual prompts or reference conditions and the necessity for accurate depiction of textures and fabric dynamics. Stable diffusion’s adaptability for swift modifications presents a promising avenue. Recent advances utilizing stable diffusion models in virtual try-on [29] signal the potential for generating garment images via stable diffusion. However, prior works have not fully exploited its capabilities in text-to-image and stylized image creation and have failed to preserve the complete patterns, e.g., stripes and texts.

Therefore, merging the detailed representation of target garments with the adaptable nature of stable diffusion promises to benefit a broader spectrum of users, including merchants, consumers, and artists, by reducing the costs related to garment-related creativity and boosting commercial effectiveness. The question arises: How can we generate images from text prompts or control conditions while preserving the intricate details of specified garments? We address this question by introducing the concept of Garment-Centric (GC) Generation, which focuses on maintaining the fidelity of garment details while enabling flexibility in image creation.

To deal with this problem, we introduce StableGarment, a unified framework built upon Stable Diffusion. This framework is meticulously designed to release the full potential of Stable Diffusion. A garment encoder is devised to encode the details of the target garment. This encoder interfaces with the stable diffusion denoising UNet through an innovative additive self-attention(ASA) mechanism, enhancing the system’s versatility in text prompting and model switching. This approach to self-attention facilitates model adaptability for creative try-on purposes. To empower the model with virtual try-on ability, a try-on controlnet is trained. It takes the input of user poses and image contexts and superimposes the garments onto the input image. Moreover, our restructured training dataset, enriched with varied text prompts, enhances the prompt following of the generated images.

To summarize, our contributions are threefold:

1. We propose a unified framework to address garment-centric (GC) generation tasks, encompassing GC text-to-image, controllable GC text-to-image, stylized GC text-to-image, and virtual try-ons within a single model.
2. We introduce an additive self-attention layer that allows for seamless model switching to stylized base-models and propose a data engine to enhance the models’ ability to follow prompts.
3. Our model’s performance is benchmarked against existing standards, where it demonstrates state-of-the-art performance among all competitors, underscoring the superiority of our approach.

2 Related Work

Subject-driven Generation. Subject-driven generation aims to generate the target subject with text prompts from given reference images. It can be classified into two main categories: test-time finetuning methods [18,31,44] and finetuning-free methods [33,47,54,58,61]. Finetuning-free approaches offer greater flexibility and are more promising for real-world applications. Typically, finetuning-free methods encode reference images into embeddings or image prompts without requiring additional finetuning. ELITE [54] proposes global and local mapping schemes, but suffers from limited fidelity. Instantbooth [47] employs an adapter structure trained on domain-specific data for subject-driven generation without finetuning. IP-Adapter [58] encodes images into prompts, while BLIP-Diffusion [33] enables efficient zero-shot setups. However, these methods focus on general

subject learning and primarily learn subject representations through the cross-attention layer, which can result in a significant loss of fine-grained details in the generated images. To address this limitation, we propose a garment encoder that captures multi-scale garment features, enabling the generation of results with enhanced fine-grained details.

Stable Diffusion. Stable Diffusion is a powerful text-to-image model derived from Latent Diffusion Models (LDM) [43]. Its ability to generate high-quality images and perform flexible text editing has contributed significantly to various vision generation tasks. However, there are two significant challenges to achieving a satisfactory generation: (1) controllability and (2) personality. In the context of controllability, ControlNet [59] is one of the most effective methods for controlling Stable Diffusion. It consists of a trainable copy of the UNet’s encoder with zero initialization and provides structural guidance during image generation. In terms of personality, IP-Adapter [58] is the first method to utilize image prompts for personality learning and has shown promising results in maintaining semantic consistency. However, when dealing with subjects that contain complicated patterns or text, these methods fail to maintain fidelity. To address this issue, reference-only [39] methods have been proposed to enhance the details of the learning subject and demonstrate more convincing results. Reference-only nets often consist of a complete copy of UNet, which helps them capture multi-scale texture features across the UNet’s feature spaces. Based on these observations, we choose to fully utilize these techniques to solve garment-centric generation tasks.

Virtual Try-on. Virtual try-on approaches can be categorized into 3D-based [7, 21, 30, 38] and image-based methods [9, 17, 22, 23, 28, 29, 32, 40, 41, 51]. Image-based methods are more promising because of their lightweight nature and the ability to generate reasonable results using large-scale try-on datasets. Image-based virtual try-on methods mainly consist of two stages: warping and blending. Warping methods such as Thin Plate Spline (TPS) [16, 22] and flow-based approaches [9, 14, 20, 56] have limitations in handling complex deformations and lack structural information about the garment [12]. Segmentation maps and DensePose are often used as additional conditions to compensate for information loss caused by occlusion [14, 32, 34].

Diffusion Probabilistic Models (DPMs) have recently shown promising results in various image synthesis and editing tasks, such as text-to-image synthesis [61], image-to-image translation [35], and image inpainting [37]. In the context of virtual try-on, researchers have explored the use of DPMs to address three key challenges: garment-agnostic preservation, garment detail recovery, and warping accuracy. To preserve garment-agnostic information, many inpainting-based methods [13, 57] have been proposed to learn the implicit warping process by inpainting the missing garment region. These methods can learn reasonable warping results but often fail to preserve fine details, such as logos and text, on the garment. To recover garment details, several approaches have adopted advanced techniques, including ControlNet [59] and image prompt learning, similar to IP-Adapter [58]. DCI-VTON [20] combines a diffusion model with a warping

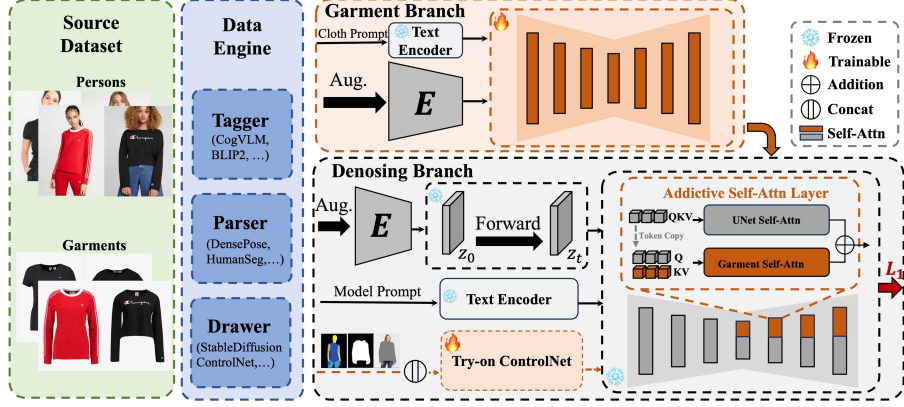


Fig. 2: Overview of our framework, consisting of a data engine, garment encoder, and try-on ControlNet. The data engine preserves the model’s capacity to follow prompts, while the garment encoder with additive self-attention layer captures garment details. Meanwhile, the try-on ControlNet is designed for virtual try-on tasks.

module, while LaDI-VTON [40] generates pseudo-words that represent garments, inspired by textual inversion [19]. TryOnDiffusion [62] employs cascaded diffusion models with conditions like human pose and garment pose, and StableVITON [29] leverages a reference garment ControlNet to enhance garment details. WarpDiffusion [36] incorporates warping garments into diffusion models using local cross attention. In this work, we utilize the garment encoder to recover garment textures and a try-on ControlNet to learn the implicit warping process, effectively addressing the challenges faced by previous methods.

3 Proposed Method

In this section, we begin by exploring the preliminaries on diffusion models as detailed in Section 3.1. Subsequently, we delve into the architecture of StableGarment in Section 3.2, outlining its overall structure, the details of the garment encoder, and the functionality of the try-on control network. The final part, described in Section 3.3, provides an in-depth explanation of our dataset along with the training and inference methodologies employed.

3.1 Preliminary

Diffusion Model (DM). DM [26, 49] belongs to the category of generative models that denoise from a Gaussian prior \mathbf{x}_T to target data distribution \mathbf{x}_0 by means of an iterative denoising procedure. The common loss used in DM is:

$$L_{simple}(\theta) := \mathbb{E}_{\mathbf{x}_0, t, \epsilon} \left[\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|_2^2 \right], \quad (1)$$

where \mathbf{x}_t is a noisy image constructed by adding noise $\epsilon \in \mathcal{N}(\mathbf{0}, \mathbf{1})$ to the natural image \mathbf{x}_0 and the network $\epsilon_\theta(\cdot)$ is trained to predict the added noise. At inference time, data samples can be generated from Gaussian noise $\epsilon \in \mathcal{N}(\mathbf{0}, \mathbf{1})$ using

the predicted noise $\epsilon_{\theta}(\mathbf{x}_t, t)$ at each timestep t with samplers like DDPM [26] or DDIM [48].

Latent Diffusion Model (LDM). LDM [43] is proposed to model image representations in autoencoder’s latent space. LDM significantly speeds up the sampling process and facilitates text-to-image generation by incorporating additional text conditions. The LDM loss is:

$$L_{LDM}(\theta) := \mathbb{E}_{z_0, t, \epsilon} \left[\|\epsilon - \epsilon_{\theta}(z_t, t, \tau_{\theta}(\mathbf{c}_t))\|_2^2 \right], \quad (2)$$

where z_0 represents image latents and $\tau_{\theta}(\cdot)$ refers to the BERT text encoder [15] used to encode text description \mathbf{c}_t .

Stable Diffusion (SD). SD is a widely adopted text-to-image diffusion model based on LDM. Compared to LDM, SD is trained on the large LAION [46] dataset and replaces BERT with the pre-trained CLIP [42] text encoder.

3.2 Network Architecture

Overview. The overall architecture of StableGarment is depicted in Fig. 2. The model is segmented into three primary components: a garment encoder, a try-on ControlNet, and a pre-trained stable diffusion model. The parameters of the pre-trained stable diffusion, including the VAE and UNet, are fixed. A trainable duplicate of the UNet is developed on top of the original, specifically designed to process garment input. Additionally, a try-on ControlNet accepts input from both the pose image and the image context. Both the garment encoder and the try-on ControlNet are integrated with the denoising UNet in an additive fashion, thereby offering the flexibility for seamless model switching.

Garment Encoder. The garment encoder, integral to our model, processes the garment input with a critical objective: to capture the intricate details of the garment. Drawing inspiration from recent studies [27, 39, 58, 59], we have devised our encoder composed of garment UNet and a pretrained VAE encoder. The garment UNet is a trainable copy of the denoising UNet and interacts with the denoising UNet via additive self-attention (ASA) layers. The underlying rationale for this configuration is to infuse reference information across multiple scales, thereby extracting the finest details from the reference image while simultaneously preserving its capacity for model switching.

The additive self-attention mechanism is specially designed to allow the encoder integration into various stylized base models. Although numerous studies have explored self-attention for connecting denoising UNet with a reference UNet, these implementations typically adopt a concatenation self-attention(CSA) mechanism, as suggested by the reference-only ControlNet [39]. However, our observations indicate that CSA often results in noticeable degradation during base-model switching under complex prompts. Inspired by the design principles of ControlNet [59] and IP-Adapter [58], which both retain the capacity for prompt following and model switching, we have redesigned the connection strategy in an additive format. This revision significantly enhances the system’s flexibility and effectiveness in model switching.

Formally, given a reference garment image I_G as input, it first undergoes the VAE encoder E to generate the garment latent, and then the garment UNet E_g generates multi-layered key-value pairs $\{(K_n, V_n)_{n=1}^N\}$, where n denotes the index of self-attention layer of the UNet structure. The K s and V s are then queried by the Q_n^u , the query tensor of the corresponding n -th self-attention layer of the denoising UNet. The output value S_n^u of the n -th self-attention layer of the denoising UNet is formulated as Eq. (3).

$$\begin{aligned} \{(K_n^g, V_n^g)_{n=1}^N\} &= E_g(E(I_G)), \\ S_n^u &= \text{Attention}(Q_n^u, K_n^g, V_n^g) + \text{Attention}(Q_n^u, K_n^u, V_n^g), \end{aligned} \quad (3)$$

where $\text{Attention}(Q, K, V)$ is the attention operator that calculate attention via query Q , key K , and value V .

Try-on ControlNet. To enable the StableGarment framework for virtual try-on tasks, we have developed a try-on ControlNet. This network is engineered to integrate the target body shape and image context into the denoising UNet workflow. The target body shape is precisely captured using DensePose [55], while the image context includes elements such as the background, skin tone, and attributes. A garment mask m is employed to highlight areas in need of inpainting. Acknowledging the tight correlation between image context, inpainting area, and the pose, these conditions are concatenated into a unified input to the try-on Controlnet. Starting with an image context I , a garment mask m is produced following the methodology outlined in [14]. Consequently, the inputs for the try-on ControlNet comprise the garment mask m , the masked image $I_m = I \cdot m$, alongside pose guidance p . The control signals of the ControlNet are directly infused into the denoising UNet as feature residuals. This setup enables the try-on ControlNet to deliver precise, pixel-aligned control signals for virtual try-ons, drawing from the provided image context to ensure accurate garment alignment.

Formally, the try-on ControlNet E_c processes the three aforementioned conditions as in Eq. (4), and the residuals of the features are added to the denoising UNet.

$$R^u = E_c(I_m \odot m \odot p), \quad (4)$$

where R^u represents feature residuals of try-on ControlNet, and \odot denotes the concatenation operation.

Text Prompt. To preserve the model’s ability to follow prompts, we dispatch distinct prompts to both the garment UNet and the denoising UNet. Specifically, the garment UNet receives prompts of the garment category, facilitating its understanding of the input garment’s dynamic nature. Concurrently, the denoising UNet is prompted with descriptions of the target image, maintaining its proficiency in generating images responsive to a diverse range of textual cues.

3.3 Training Strategy and Inference

Dataset Preparation. For the preparation of our dataset, we employ a pre-trained model to generate high-quality synthesized garment data, which is crucial

for a precise and controllable garment-centric learning process. The primary objective of the data engine is to generate garment-centric images across a wide range of text prompts, thereby preserving the model’s ability to follow text prompts throughout the learning phase. A detailed demonstration can be found in the supplementary material. Our data engine is divided into three modules:

- **Parser:** The parser module is tasked with creating a parse map and dense map from an input image. It incorporates a pre-trained parsing segmentation network alongside Detectron2 [55], facilitating accurate segmentation of the garment.
- **Tagger:** The tagger module generates detailed text descriptions for target garment images. It utilizes CogVLM-chat [52] for generating text descriptions, and GPT4 to create inpainting templates, which are subsequently used by the drawer module.
- **Drawer:** The drawer module functions as the synthesized data generator. It uses ControlNet [59] in conjunction with a pre-trained inpainting model, epiCRealism [4], to produce the final synthetic garment images.

Learning Objectives. The training for the StableGarment model incorporates a two-stage strategy. Initially, in the first phase, we utilize synthetic data to train our garment encoder. This training focuses on accurately capturing the nuances of clothing details while maintaining the ability for textual modifications. In the second stage, we freeze our garment encoder and only train our try-on ControlNet. We persist with the same learning objective, incorporating tryon-specific conditions. The loss functions for the two stages are expressed as Eq. (5):

$$\begin{aligned} L &= \mathbb{E}_{z_0, t, I_g, \epsilon \in \mathcal{N}(0,1)} [\|\epsilon - \epsilon_\theta\|_2^2], \\ L &= \mathbb{E}_{z_0, t, I_g, p, m, I, \epsilon \in \mathcal{N}(0,1)} [\|\epsilon - \epsilon_\theta\|_2^2], \end{aligned} \quad (5)$$

where I_g and I are the garment reference image and the image context respectively, p is the pose guidance from the image context, and m denotes the corresponding garment mask.

Inference. As illustrated in Table 1, StableGarment is capable of dealing with various tasks. In addition to using different components for different tasks, the inference pipelines are different. In the context of garment-centric generation tasks, our method is flexible and generalizable to other third-party plugins, e.g., stylized base models and ControlNets. Specifically, for virtual try-on tasks, we employ an inpainting pipeline to achieve both virtual try-on consistency and cloth-agnostic region preservation. Concretely, each updating denoising step can be rewritten as:

$$\epsilon'_{t-1} = m \cdot \epsilon_{t-1} + (1 - m) \cdot z_{t-1}^{src}, \quad (6)$$

where z_{t-1}^{src} represents the noised image latents from the source model input and ϵ_{t-1} represents predicted noise. This adjustment ensures the seamless preservation of background elements and model details.

Table 1: The StableGarment is capable of performing various garment-centric(GC) tasks. The configurations for performing different tasks are listed in the table.

Task	Pipeline	Garment encoder	Base model	ControlNet
GC t2i	text2image	✓	stable diffusion v1.5	×
stylized GC t2i	text2image	✓	stylized base-model	×
controllable GC t2i	text2image	✓	any base-model	any ControlNet
virtual try-on	inpainting	✓	stable diffusion v1.5	try-on ControlNet

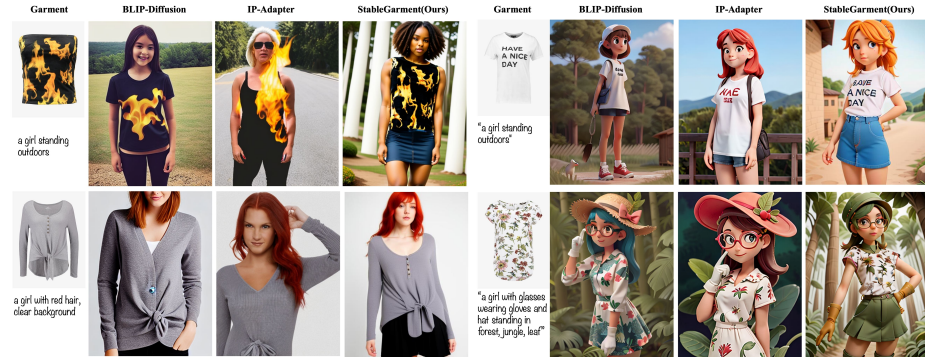


Fig. 3: Comparison with subject-driven generation methods.

4 Experiments

Baselines. For subject-driven generation, we compare our method with three finetuning-free method, including ELITE [54], IP-Adapter [58] and BLIP-Diffusion [33]. In the context of virtual try-on task, we compare our method with three GAN-based virtual try-on methods, VITON-HD [14], HR-VITON [32] and GP-VTON [56], and three Diffusion-based virtual try-on methods, LADI-VTON [40], DCI-VTON [20] and StableVITON [29]. We also compare our method with Diffusion-based inpainting methods, Paint-by-Example [57] and AnyDoor [13]. To ensure a fair comparison, we utilize the official implementations for all the aforementioned methods at 512×384 resolution.

Datasets. For subject-driven generation, we introduce a garment-centric generation benchmark derived from a subset of the VITON-HD test set. This benchmark, which includes 14 garments and 6 text prompts, aims to cover a broad spectrum of garment attributes such as color, shape, pattern, size, and type. The prompts, crafted through GPT-4, integrate a diverse range of skin tones, hairstyles, backgrounds, and additional attributes, enabling a robust evaluation against subject-driven generation approaches. The details of proposed benchmark will be provided in the supplementary material.

In the context of the virtual try-on task, we employ two widely-used, publicly available high-resolution datasets for virtual try-on: VITON-HD [14] and Dress Code [41]. We keep the same train-test split as the previous methods [29, 40] on both of these datasets. In our study, we only use the official datasets for training the virtual try-on tasks to ensure fairness.



Fig. 4: Qualitative comparison with baselines on VITON-HD dataset.

Evaluation Metrics. For virtual try-on task, we keep consistent with previous methodologies and apply SSIM [53] and LPIPS [60] metrics for evaluation under paired settings, and FID [25] and KID [10] metrics for unpaired settings. Following previous studies [8, 45, 50, 61], we utilize the DINO-M metric to assess the fidelity of the generated image to the reference image. The DINO-M metric compares a masked version of the try-on image with the garment image and computes the similarity using image embeddings extracted by the DINO [11] model, reflecting the identity of the garment. As the objectives of subject-driven generation differ from try-on, we employ three distinct metrics CLIP-T, CLIP-I [24] and Aesthetic score [46] to compare with subject-driven generation methods.

We adopt two metrics, human preference and human scores, to integrate human judgment into our evaluation process. In our experimental setup, we compile composite images generated by various methods for 100 randomly selected pairs from the test set. Following [29], totally 200 evaluators were asked to evaluate the three aspects of the try-on results, garment identity, try-on quality and garment-agnostic preservation. In each aspect, human preference showcases the frequency with which each method was selected as the top choice in these baselines, and human scores reflect the weighted evaluation scores based on their relative ranking order. The formula computation and our implementation details will be provided in the supplementary material.

4.1 Qualitative Results

Subject-driven Generation Results. Fig. 3 displays our StableGarment can generate reasonable outputs following the different prompts while preserving



Fig. 5: Comparison on Dress Code dataset.

clothing texture under different diffusion base models. BLIP-Diffusion and IP-Adapter can only preserve partial texture features of the garment, making them prone to losing details such as shape and text. In contrast, our method not only achieves better fidelity in texture, shape, and detail preservation but also offers adaptability in style transfer when switching base models.

Virtual Try-on Results. Fig. 4 shows a visual comparison of our outputs to other baselines on VITON-HD. Although previous methods synthesize reasonable warping garments and preserve the background well, they tend to generate blurry try-on results with manifest artifacts. Diffusion-based methods greatly improve the garment generation quality, but most of these methods fail to capture accurate clothing texture, e.g. text and tiny patterns. Due to the meaningful feature spaces of our proposed garment encoder, the try-on cloth successfully preserves the slight details. For example, our method accurately captures the ‘ADIDAS’ on the cloths in the last row. The quantitative results on VTION-HD convincingly verify that our method can achieve a more realistic try-on effect for these garments. More examples of our virtual try-on results will be reported in the supplementary materials.

We also provide a visual comparison of our outputs to other baselines on the Dress Code in Fig. 5. Compared to GP-VTON and LADI-VTON, our try-on results exhibit a more natural appearance.

4.2 Quantitative Results

Comparison with Subject-driven Generation Methods. We propose a garment-centric benchmark as described above to evaluate our method and other baselines, which includes 14 garments and 6 different prompts. Table 2 provides a detailed quantitative evaluation of these methods. Overall, both IP-Adapter and our method demonstrate the ability to generate accurate outputs that adhere to the given text descriptions. Moreover, our method clearly outperforms previous state-of-the-art fine-tuning-free methods on the CLIP-I and Aesthetic score metrics, showcasing the capability of our garment encoder to capture comprehensive and fine-grained features from reference garments.

Table 2: Quantitative comparison of subject-driven generation methods. The best and the second best results are denoted as **Bold** and underline, respectively.

Model	CLIP-T \uparrow	CLIP-I \uparrow	AS \uparrow
BLIP-Diffusion [33]	0.187	<u>0.692</u>	4.973
IP-Adapter [58]	0.274	0.628	<u>5.165</u>
ELITE [54]	0.244	0.617	4.468
Ours	<u>0.262</u>	0.719	5.255

Table 3: Quantitative comparison on VITON-HD dataset. We multiply KID by 100 for better comparison. The best and the second best results are denoted as **Bold** and underline, respectively.

Model	LPIPS \downarrow	SSIM \uparrow	FID \downarrow	KID \downarrow	DINO-M \uparrow
VITON-HD [14]	0.117	0.862	12.12	0.323	0.615
HR-VTON [32]	0.105	0.868	11.27	0.273	0.634
LADI-VTON [40]	0.092	0.875	9.37	0.158	0.633
DCI-VTON* [20]	0.092	0.876	9.44	0.164	0.648
StableVITON [29]	0.083	0.866	<u>8.19</u>	<u>0.118</u>	0.647
GP-VTON* [56]	<u>0.083</u>	0.887	9.83	0.141	0.675
Paint-by-Example [57]	0.217	0.776	14.17	0.825	0.590
AnyDoor [13]	0.139	0.828	12.73	0.544	0.641
Ours	0.077	<u>0.877</u>	7.98	0.104	<u>0.661</u>

Table 4: User Study on virtual try-on task. We evaluate our methods with other four baselines on three metrics: garment identity, try-on quality and garment-agnostic preservation. The best and the second best results are denoted as **Bold** and underline, respectively.

Model	Human Preference(%) \uparrow			Human Scores \uparrow		
	Identity	Quality	Preservation	Identity	Quality	Preservation
GP-VTON [56]	<u>15.35</u>	13.82	15.85	<u>2.98</u>	2.72	2.68
DCI-VTON [20]	13.68	12.51	15.13	2.82	2.64	2.78
LADI-VTON [40]	10.88	12.56	14.06	2.12	2.22	2.41
StableVITON [29]	13.40	<u>14.32</u>	<u>22.77</u>	2.69	<u>2.99</u>	3.02
Ours	46.70	46.80	27.30	3.15	3.12	<u>2.99</u>

Comparison with Visual Try-on Methods. We compare our method with existing baselines on the VITON-HD dataset and report the results in Table 3. Our method demonstrates competitive performance across all metrics compared to the baselines, particularly in the unpaired metrics(i.e., FID and KID). The success of our approach can be attributed to the garment encoder and DensePose-based try-on ControlNet, which effectively capture fine-grained garment features and generate high-quality try-on images. Furthermore, we utilize DINO-M to evaluate the recovery of cloth texture in the unpaired setting. Under this metric, our method displays competitive performance, further demonstrating the superiority of our method in the virtual try-on task.

Table 5: Ablation study for garment-centric generation. For User result “a / b”, a is frequency that each method is chosen as the best result for restoring the clothes, and b represents the best generated result following the text description.

Ablation Setups	SD1.5			Anythingv5 [1]			Disney Pixar [3]		
	CLIP-T \uparrow	User(%) \uparrow	AS \uparrow	CLIP-T \uparrow	User(%) \uparrow	AS \uparrow	CLIP-T \uparrow	User(%) \uparrow	AS \uparrow
Ours(w/o ASA)	0.328	32.57/38.24	5.430	0.332	12.29/25.14	5.860	0.316	10.78/22.30	5.631
Ours(w/o Synthesized data)	0.254	35.04/28.71	5.107	0.275	21.83/14.86	5.625	0.280	24.71/12.99	6.02
Ours(full)	0.327	32.38/33.04	5.598	0.321	65.87/60.00	5.919	0.314	64.50/64.71	5.865

4.3 User Study

To further evaluate the generation quality of our model, we conduct a user study to measure both the realism of the generated images and their coherence with the inputs given to the virtual try-on model. As shown in Table 4, our method outperforms other baselines in garment texture and try-on quality, demonstrating the effectiveness of our approach in preserving garment texture and maintaining try-on quality. Although our method does not achieve the best performance in the traditional quantitative metrics, it still showcases satisfactory performance in garment-agnostic preservation.

4.4 Ablation Study

We take the 512×384 resolution on the VITON-HD dataset as the basic setting and perform ablation studies to validate the effectiveness of each component of our method.

Additive Self-Attention. To study the effectiveness of our ASA operator, we built a model with concatenated self-attention and trained with the same setup. The qualitative and quantitative comparison can be found in Fig. 6 and Table 5. When testing under standard setup, the quantitative performance of both models is comparable. However, when switching the base models into stylized ones, it is visually obvious that the images generated via the CSA model often suffer from unavoidable artifacts. Especially when the text prompt becomes complex, such as a detailed background description, the CSA model often misinterprets the attributes of the garment. Instead, the ASA counterpart generates stylized images with better quality. Both details of the garment and style of the base model can be well preserved by the ASA model.

Effects of Synthesized Data. We also ablate the benefits of using synthesized data. The comparing setups are: 1) training with VITON-HD and 2) training with our synthesized data. From Table 5, we found that on the GC generation benchmark, their performances are close under the standard settings. However, when testing the prompt following capacity on other stylized base models, the model trained with synthesized data significantly outperforms its counterpart. From Fig. 6, the images generated by the model trained without synthesized data would display monotonous backgrounds, which reflects a characteristic of the VITON-HD dataset.

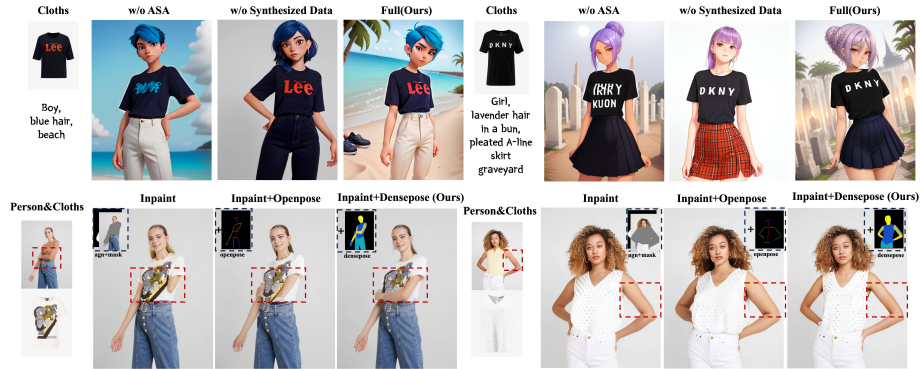


Fig. 6: Effects of our proposed components.

Table 6: Ablation study for try-on ControlNet.

Ablation Setups	LPIPS↓	SSIM↑	FID↓	KID↓	DINO-M↑
Inpainting	0.097	0.852	9.11	0.303	0.661
+OpenPose	0.080	0.870	8.60	0.251	0.663
+DensePose(Ours)	0.077	0.877	8.02	0.110	0.661

Comparison with Different Try-on ControlNet. We evaluate three types of try-on ControlNet, with three different control inputs: (1) garment-agnostic input and mask, (2) garment-agnostic input, mask and OpenPose, and (3) garment-agnostic input, mask, and DensePose. Fig. 6 shows a visual comparison of these three different setups of try-on ControlNets. All of these ControlNets enable to preserve detailed clothing texture due to our proposed garment encoder. However, without the aid of geometric information, it is difficult for our model to generate accurate human body, including arms and hands. In comparison to OpenPose control, DensePose control offers more geometric information. As shown in Fig. 6, replacing OpenPose with DensePose as a condition results in the disappearance of artifacts caused by occlusion in the generated images, with the thickness of the arms also better matching source model. We also conduct a comprehensive quantitative experiment to verify the effectiveness of our proposed DensePose-based try-on ControlNet. As presented in Table 6, when adding the OpenPose condition on top of the base inpainting ControlNet, all metrics improve, as the model gains more information about the human body structure. Replacing OpenPose with DensePose leads to improvements in the detailed morphology of the human body, and all metrics also show enhancements.

5 Conclusion

In this paper, we propose our novel unified framework, StableGarment, to solve garment-driven generation tasks, including garment-centric text-to-image, garment-centric controllable text-to-image, garment-centric stylized image generation,

and virtual try-on with flexible inputs. We introduce a garment encoder to capture detailed clothing features, which consists of a trainable copy of UNet with designed additive self-attention (ASA) layers. These designs help our encoder capture detailed clothing features and switch between stylized base models. Furthermore, to preserve the model’s ability to follow prompts, we propose a novel data engine to generate high-quality synthesized data. We also present a carefully designed try-on ControlNet to address the classic virtual try-on task. Extensive experiments have demonstrated that our approach delivers state-of-the-art (SOTA) results among existing virtual try-on methods and exhibits high flexibility with broad potential applications in various related fields.

References

1. Anything v5. <https://civitai.com/models/9409?modelVersionId=30163> (2023) 13
2. Densepose controlnet. <https://civitai.com/models/120149/controlnet-for-densepose> (2023) 20
3. Disney pixar cartoon type a. <https://civitai.com/models/65203/disney-pixar-cartoon-type-a> (2023) 13
4. epicrealism v5-inapinting. <https://civitai.com/models/25694?modelVersionId=134361> (2023) 8, 20
5. Ip-adapter-faceid. <https://huggingface.co/h94/IP-Adapter-FaceID> (2023) 22
6. stabilityai/sd-vae-ft-mse-original. <https://huggingface.co/stabilityai/sd-vae-ft-mse-original> (2023) 23
7. Aggarwal, A., Wang, J., Hogue, S., Ni, S., Budagavi, M., Guo, X.: Layered-garment net: Generating multiple implicit garment layers from a single image. In: Proceedings of the Asian Conference on Computer Vision. pp. 3000–3017 (2022) 4
8. Avrahami, O., Aberman, K., Fried, O., Cohen-Or, D., Lischinski, D.: Break-a-scene: Extracting multiple concepts from a single image. arXiv preprint arXiv:2305.16311 (2023) 10
9. Bai, S., Zhou, H., Li, Z., Zhou, C., Yang, H.: Single stage virtual try-on via deformable attention flows. In: European Conference on Computer Vision. pp. 409–425 (2022) 4
10. Bińkowski, M., Sutherland, D.J., Arbel, M., Gretton, A.: Demystifying mmd gans. arXiv preprint arXiv:1801.01401 (2018) 10
11. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 9650–9660 (2021) 10
12. Chen, C.Y., Chen, Y.C., Shuai, H.H., Cheng, W.H.: Size does matter: Size-aware virtual try-on via clothing-oriented transformation try-on network. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 7513–7522 (2023) 4
13. Chen, X., Huang, L., Liu, Y., Shen, Y., Zhao, D., Zhao, H.: Anydoor: Zero-shot object-level image customization. arXiv preprint arXiv:2307.09481 (2023) 2, 4, 9, 12
14. Choi, S., Park, S., Lee, M., Choo, J.: Viton-hd: High-resolution virtual try-on via misalignment-aware normalization. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 14131–14140 (2021) 4, 7, 9, 12, 21
15. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018) 6
16. Duchon, J.: Splines minimizing rotation-invariant semi-norms in sobolev spaces. In: Constructive Theory of Functions of Several Variables: Proceedings of a Conference Held at Oberwolfach April 25–May 1, 1976. pp. 85–100. Springer (1977) 4
17. Fenocchi, E., Morelli, D., Cornia, M., Baraldi, L., Cesari, F., Cucchiara, R.: Dual-branch collaborative transformer for virtual try-on. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2247–2251 (2022) 4

18. Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A.H., Chechik, G., Cohen-Or, D.: An image is worth one word: Personalizing text-to-image generation using textual inversion. arXiv preprint arXiv:2208.01618 (2022) [3](#)
19. Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A.H., Chechik, G., Cohen-Or, D.: An image is worth one word: Personalizing text-to-image generation using textual inversion. arXiv preprint arXiv:2208.01618 (2022) [5](#)
20. Gou, J., Sun, S., Zhang, J., Si, J., Qian, C., Zhang, L.: Taming the power of diffusion models for high-quality virtual try-on with appearance flow. In: Proceedings of the 31st ACM International Conference on Multimedia. pp. 7599–7607 (2023) [2](#), [4](#), [9](#), [12](#), [21](#)
21. Halimi, O., Stuyck, T., Xiang, D., Bagautdinov, T., Wen, H., Kimmel, R., Shiratori, T., Wu, C., Sheikh, Y., Prada, F.: Pattern-based cloth registration and sparse-view animation. ACM Transactions on Graphics (TOG) pp. 1–17 (2022) [4](#)
22. Han, X., Wu, Z., Wu, Z., Yu, R., Davis, L.S.: Viton: An image-based virtual try-on network. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7543–7552 (2018) [4](#)
23. Han, X., Wu, Z., Wu, Z., Yu, R., Davis, L.S.: Viton: An image-based virtual try-on network. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7543–7552 (2018) [4](#)
24. Hessel, J., Holtzman, A., Forbes, M., Bras, R.L., Choi, Y.: Clipscore: A reference-free evaluation metric for image captioning. arXiv preprint arXiv:2104.08718 (2021) [10](#)
25. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems (2017) [10](#)
26. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. Advances in neural information processing systems pp. 6840–6851 (2020) [5](#), [6](#)
27. Hu, L., Gao, X., Zhang, P., Sun, K., Zhang, B., Bo, L.: Animate anyone: Consistent and controllable image-to-video synthesis for character animation. arXiv preprint arXiv:2311.17117 (2023) [6](#)
28. Issenhuth, T., Mary, J., Calauzenes, C.: 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. pp. 619–635 (2020) [4](#)
29. Kim, J., Gu, G., Park, M., Park, S., Choo, J.: Stableviton: Learning semantic correspondence with latent diffusion model for virtual try-on. arXiv preprint arXiv:2312.01725 (2023) [2](#), [4](#), [5](#), [9](#), [10](#), [12](#), [21](#)
30. Korosteleva, M., Lee, S.H.: Neuraltailor: Reconstructing sewing pattern structures from 3d point clouds of garments. ACM Transactions on Graphics (TOG) pp. 1–16 (2022) [4](#)
31. Kumari, N., Zhang, B., Zhang, R., Shechtman, E., Zhu, J.Y.: Multi-concept customization of text-to-image diffusion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1931–1941 (2023) [3](#)
32. Lee, S., Gu, G., Park, S., Choi, S., Choo, J.: High-resolution virtual try-on with misalignment and occlusion-handled conditions. In: European Conference on Computer Vision. pp. 204–219 (2022) [4](#), [9](#), [12](#), [21](#)
33. Li, D., Li, J., Hoi, S.: Blip-diffusion: Pre-trained subject representation for controllable text-to-image generation and editing. Advances in Neural Information Processing Systems **36** (2024) [3](#), [9](#), [12](#)

34. Li, K., Chong, M.J., Zhang, J., Liu, J.: Toward accurate and realistic outfits visualization with attention to details. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 15546–15555 (2021) [4](#)
35. Li, P., Wang, R., Huang, H., He, R., He, Z.: Pluralistic aging diffusion autoencoder. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 22613–22623 (2023) [4](#)
36. Li, X., Kampffmeyer, M., Dong, X., Xie, Z., Zhu, F., Dong, H., Liang, X., et al.: Warpdiffusion: Efficient diffusion model for high-fidelity virtual try-on. arXiv preprint arXiv:2312.03667 (2023) [5](#)
37. Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R., Van Gool, L.: Repaint: Inpainting using denoising diffusion probabilistic models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11461–11471 (2022) [4](#)
38. Majithia, S., Parameswaran, S.N., Babar, S., Garg, V., Srivastava, A., Sharma, A.: Robust 3d garment digitization from monocular 2d images for 3d virtual try-on systems. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 3428–3438 (2022) [4](#)
39. Mikubill: sd-webui-controlnet (2023), <https://github.com/Mikubill/sd-webui-controlnet>, gitHub repository [4](#), [6](#)
40. Morelli, D., Baldrati, A., Cartella, G., Cornia, M., Bertini, M., Cucchiara, R.: Ladvton: Latent diffusion textual-inversion enhanced virtual try-on. arXiv preprint arXiv:2305.13501 (2023) [4](#), [5](#), [9](#), [12](#), [21](#)
41. Morelli, D., Fincato, M., Cornia, M., Landi, F., Cesari, F., Cucchiara, R.: Dress code: High-resolution multi-category virtual try-on. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2231–2235 (2022) [2](#), [4](#), [9](#), [21](#)
42. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763 (2021) [6](#)
43. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10684–10695 (2022) [2](#), [4](#), [6](#)
44. Ruiz, N., Li, Y., Jampani, V., Pritch, Y., Rubinstein, M., Aberman, K.: 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. pp. 22500–22510 (2023) [3](#)
45. Ruiz, N., Li, Y., Jampani, V., Pritch, Y., Rubinstein, M., Aberman, K.: 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. pp. 22500–22510 (2023) [10](#)
46. Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al.: Laion-5b: An open large-scale dataset for training next generation image-text models. Advances in Neural Information Processing Systems pp. 25278–25294 (2022) [6](#), [10](#)
47. Shi, J., Xiong, W., Lin, Z., Jung, H.J.: Instantbooth: Personalized text-to-image generation without test-time finetuning. arXiv preprint arXiv:2304.03411 (2023) [3](#)
48. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. arXiv preprint arXiv:2010.02502 (2020) [6](#)

49. Song, Y., Sohl-Dickstein, J., Kingma, D.P., Kumar, A., Ermon, S., Poole, B.: Score-based generative modeling through stochastic differential equations. arXiv preprint arXiv:2011.13456 (2020) [5](#)
50. Voynov, A., Chu, Q., Cohen-Or, D., Aberman, K.: $p+$: Extended textual conditioning in text-to-image generation. arXiv preprint arXiv:2303.09522 (2023) [10](#)
51. Wang, B., Zheng, H., Liang, X., Chen, Y., Lin, L., Yang, M.: Toward characteristic-preserving image-based virtual try-on network. In: Proceedings of the European conference on computer vision (ECCV). pp. 589–604 (2018) [4](#)
52. Wang, W., Lv, Q., Yu, W., Hong, W., Qi, J., Wang, Y., Ji, J., Yang, Z., Zhao, L., Song, X., Xu, J., Xu, B., Li, J., Dong, Y., Ding, M., Tang, J.: Cogvlm: Visual expert for pretrained language models (2023) [8](#), [20](#)
53. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. IEEE transactions on image processing pp. 600–612 (2004) [10](#)
54. Wei, Y., Zhang, Y., Ji, Z., Bai, J., Zhang, L., Zuo, W.: Elite: Encoding visual concepts into textual embeddings for customized text-to-image generation. arXiv preprint arXiv:2302.13848 (2023) [3](#), [9](#), [12](#)
55. Wu, Y., Kirillov, A., Massa, F., Lo, W.Y., Girshick, R.: Detectron2. <https://github.com/facebookresearch/detectron2> (2019) [7](#), [8](#), [20](#)
56. Xie, Z., Huang, Z., Dong, X., Zhao, F., Dong, H., Zhang, X., Zhu, F., Liang, X.: Gp-vton: Towards general purpose virtual try-on via collaborative local-flow global-parsing learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 23550–23559 (2023) [2](#), [4](#), [9](#), [12](#), [21](#)
57. Yang, B., Gu, S., Zhang, B., Zhang, T., Chen, X., Sun, X., Chen, D., Wen, F.: Paint by example: Exemplar-based image editing with diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18381–18391 (2023) [4](#), [9](#), [12](#), [21](#)
58. Ye, H., Zhang, J., Liu, S., Han, X., Yang, W.: Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. arXiv preprint arXiv:2308.06721 (2023) [3](#), [4](#), [6](#), [9](#), [12](#), [22](#)
59. Zhang, L., Rao, A., Agrawala, M.: Adding conditional control to text-to-image diffusion models. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 3836–3847 (2023) [4](#), [6](#), [8](#), [22](#)
60. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 586–595 (2018) [10](#)
61. Zhang, Y., Liu, J., Song, Y., Wang, R., Tang, H., Yu, J., Li, H., Tang, X., Hu, Y., Pan, H., et al.: Ssr-encoder: Encoding selective subject representation for subject-driven generation. arXiv preprint arXiv:2312.16272 (2023) [3](#), [4](#), [10](#)
62. Zhu, L., Yang, D., Zhu, T., Reda, F., Chan, W., Saharia, C., Norouzi, M., Kemelmacher-Shlizerman, I.: Tryondiffusion: A tale of two unets. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4606–4615 (2023) [2](#), [5](#)

Supplementary

In this supplementary material, we have provided additional details and experiments to support our main paper. We have introduced our data engine and implementation details in Section A and Section B, respectively. In Section C, we have presented additional quantitative experiments, including the evaluation of the cross-dataset analysis between VITON-HD and Dress Code training datasets. We have also introduced our garment benchmark in Section D. In Section E, we have provided more information about our user study, including metric computation and query format. We have discussed more applications in Section F and the limitations of our method in Section G.

A Data Engine

The overall framework for our data engine are as shown in Fig. 7. The complete workflow is as follows:

1. Adopt human segmentation model and DensePose prediction model [55] to get parsing map and dense map.
2. Combine predicted maps with morphological operations to estimate the garment-agnostic mask and corresponding image content.
3. Utilize GPT4 to generate inpaint prompts and employ the DensePose ControlNet [2] and Stable Diffusion inpainting model [4] to inpaint the synthesized outputs.
4. Use CogVLM [52] as a tagger to label each image with a detailed description.
5. Repeat steps 1-4 to generate satisfactory synthesized outputs.

B Implementation Details

For subject-driven generation, we employed Stable Diffusion V1-5 as the pre-trained diffusion model and trained it on our synthesized dataset. Our synthesized dataset consists of 11,436 images with corresponding DensePose, text description and parsing images. We adopt the flip operation for data augmentation. The model underwent 50k iterations of training on 8 H800 GPUs, with a batch size of 8 per GPU and a learning rate of $1e-4$. Inference was performed using DDIM as the sampler, with a step size of 25 and a guidance scale set to 3. In the context of virtual try-on training, we follow the previous setup at the first stage and only train our garment encoder on official dataset. At the second stage, we frozen our garment encoder and only train our DensePose try-on ControlNet. We adopt flip, random shift and random scale for data augmentation. We initialize our ControlNet with the DensePose ControlNet [2]. We train our second stage for 20k iterations on 8 H800 GPUs, with a batch size of 8 per GPU and a learning rate of $1e-5$. Inference was performed using DDIM as the sampler with inpainting pipeline, with a step size of 25 and a guidance scale set to 3. The detail of our garment encoder is shown in Fig. 8.

C Additional Quantitative Experiments

Table 7: Quantitative comparison on cross data setting. The best and the second best results are denoted as **Bold** and underline, respectively.

Model	LPIPS↓	SSIM↑	FID↓	KID↓
VITON-HD [14]	0.187	0.853	44.26	2.882
HR-VTON [32]	0.108	0.909	19.97	0.735
LADI-VTON [40]	0.101	0.901	16.34	0.536
DCI-VTON* [20]	0.124	0.898	18.81	0.802
StableVITON* [29]	<u>0.060</u>	<u>0.911</u>	<u>12.58</u>	<u>0.170</u>
GP-VTON* [56]	0.385	0.887	65.71	6.601
Paint-by-Example [57]	0.087	0.889	14.17	0.478
Ours	0.046	0.944	11.15	0.063

In this section, we conduct a cross-dataset analysis between the VITON-HD [14] and Dress Code [41] training datasets. We report the results in Table 7, presenting only the metrics measured by previous methods. The table demonstrates the effectiveness of our proposed method, as we outperform other methods across all metrics. Furthermore, the qualitative results presented below verify the remarkable capacity of our approach in adapting to different garment domains while preserving intricate garment details.

D Garment Benchmark

We have developed a garment-centric generation benchmark specifically tailored for evaluating subject-driven generation methods, serving as a comparative platform against our own methodology. Primarily honing our model for the try-on task, we meticulously curated this benchmark from a subset of the VITON-HD test set. Comprising 14 distinct garments and 6 meticulously crafted text prompts, our benchmark offers a comprehensive evaluation framework. To ensure a diverse representation, we selected garments based on various attributes including color, shape, pattern, size, and more, although not exhaustive in its coverage. These prompts, generated using GPT-4 with a consistent format, mirror the essence of try-on tasks by focusing on the individual’s appearance while incorporating additional attributes such as skin tone and hairstyle.

The generated prompts and selected garments with corresponding ids are shown in the Fig.9 and Fig.10, respectively.

We compare our methods with three subject-driven methods, BLIP-Diffusion, IP-Adapter and ELITE. We take official implementation of IP-Adapter and ELITE and use BLIP-Diffusion implementation by diffusers. We set all parameters following examples provided by implementation, respectively.

E User Study

Metric Computation. In our user study to assess the effectiveness of virtual try-on tasks, we utilize two principal metrics: human preference and human scores. Participants are requested to rank the results from our model alongside those from comparative baseline models, focusing on three key aspects: garment identity, try-on quality, and the garment-agnostic preservation. Human preference is measured by the frequency at which each method is ranked as the preferred choice relative to the baselines. Human scores, on the other hand, are derived from a weighted scheme reflecting the participants’ rankings. The formula for computing the human scores is as follows:

$$S = \frac{\sum(f \times W)}{N}, \quad (7)$$

where S refers to the average comprehensive score of an option, f refers to the preference frequency, N refers to the total number of responses for the item, and W refers to the inverse rank assigned by a participant, multiplied by the total count of methods under comparison. This weighting system is designed to proportionally recognize the preferences indicated by participants. The format of our used query is as shown in Fig. 11.

F More Applications

Virtual try-on task. We present the results generated by our StableGarment model in Figures 12 and 13. Our model can learn accurate warping transformations while preserving the intricate details of the garments.

Combined with IP-Adapter [58]. Our model, when combined with the IP-Adapter [58], enables the generation of target individuals wearing target garments. We leverage the ID preservation capability of IP-Adapter FaceID-PlusV2 [5] to deliver an authentic try-on experience. The visual results, shown in Figure 16, demonstrate the remarkable compatibility of our methods with current plugins.

E-commerce model generation. Leveraging the capabilities of ControlNet [59], our model can generate e-commerce models guided by specific conditions, such as OpenPose and DensePose. We present the OpenPose-guided generation results in Figure 14.

Stylized garment-centric generation. Furthermore, by replacing the standard Stable Diffusion 1.5 model with other diverse base models, we can generate creative and stylized outputs while preserving the intricate details of the garments. We present these results in Figure 15.

G Limitations and Discussion

Our model faces two main challenges: the VAE reconstruction problem and the generation of incorrect accessories. Fig.17 shows that the standard VAE used

in Stable Diffusion 1.5 fails to preserve all the details of the garment through a simple reconstruction. Although we tried to employ a more advanced version of VAE [6], the problem cannot be fully solved. Therefore, research on how to better preserve detailed information may need further improvement. Meanwhile, in the context of incorrect accessory generation, we found that the trained model may tend to generate some incorrect accessories during inference, especially for the high-resolution version. This is often caused by inaccurate parsing conditions, e.g., garment-agnostic masks or DensePose. We leave this problem for future work.

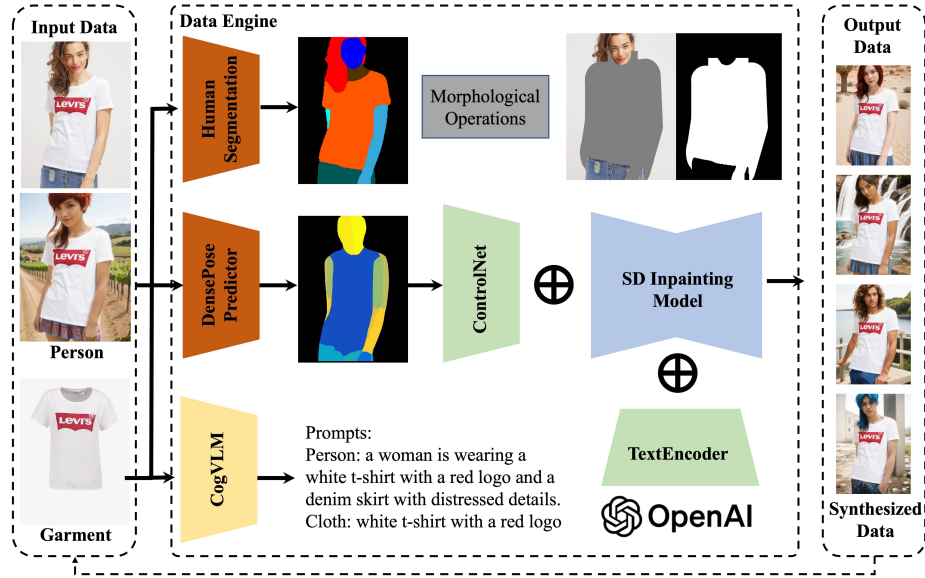


Fig. 7: Data engine framework.

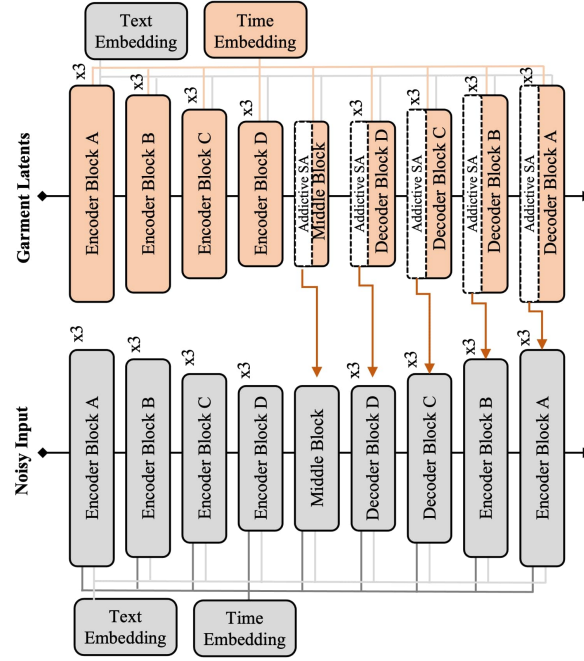


Fig. 8: The framework of our proposed garment encoder.

```
prompt_list = [
    "A woman with shoulder-length wavy blonde hair and fair skin.",
    "An elderly woman with elegant gray hair and deep-set eyes.",
    "A girl with long, straight jet-black hair and a serene expression.",
    "A female with tightly coiled afro-textured hair and rich brown skin.",
    "A teen with a short spiky platinum blond hair and a subtle tan.",
    "A person with long, flowing red hair and light green eyes."
]
```

Fig. 9: Evaluation prompts for subject-driven generation methods.

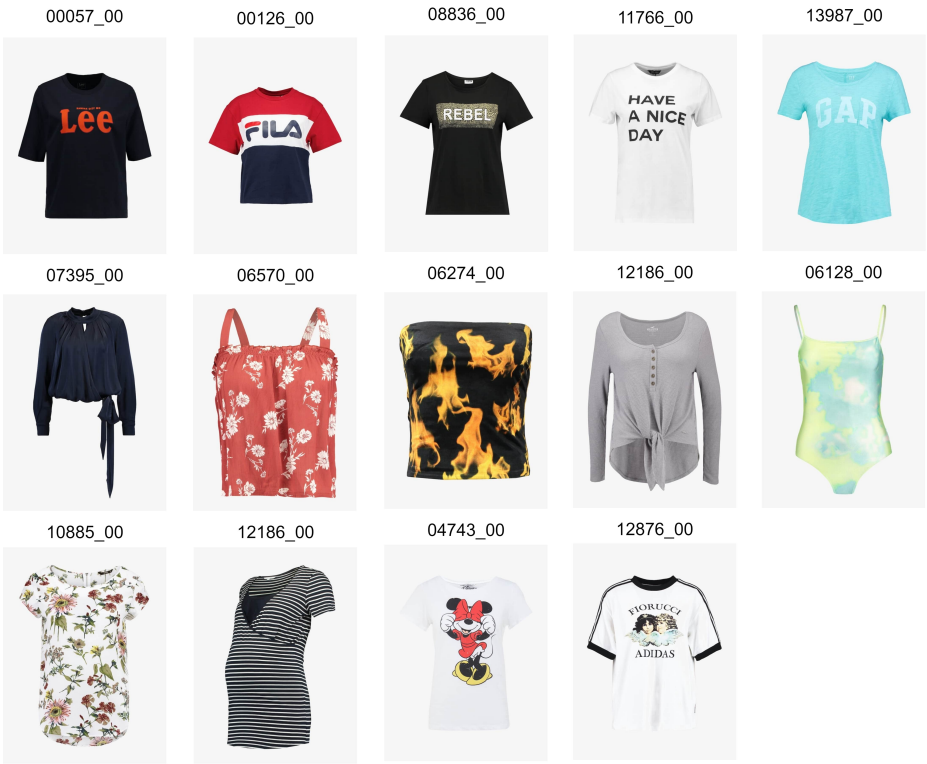


Fig. 10: The selected garments for evaluating subject-driven generation methods.

*1. Which of the following methods do you think best preserves the details of clothing:



Sort in ascending order

- ☐ Method A
- ☐ Method B
- ☐ Method C
- ☐ Method D
- ☐ Method E

Fig. 11: Query format.



Fig. 12: Our try-on results on VITON-HD dataset at 512x384 resolution.



Fig. 13: Our try-on results on VITON-HD dataset at 512x384 resolution.

Cloths&OpenPose

Diverse E-ommerce Outputs

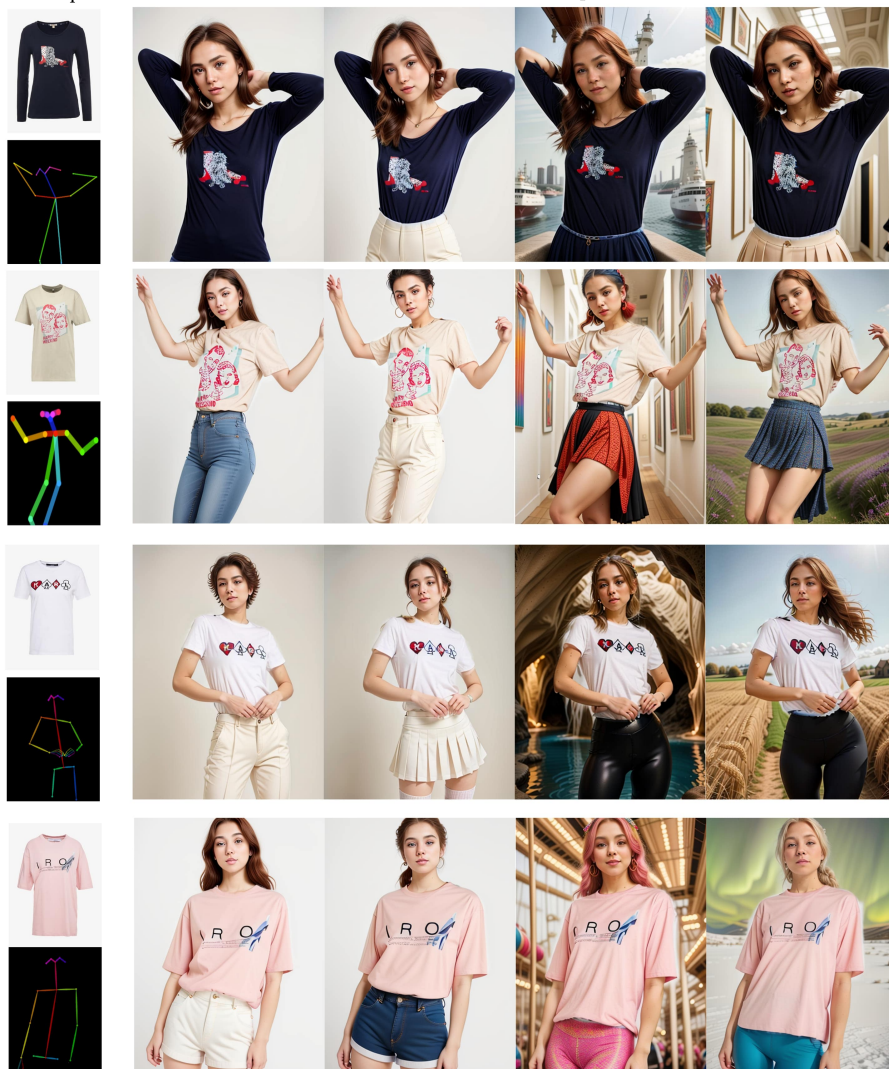


Fig. 14: E-commerce model generation.



Fig. 15: Stylized garment-centric generation.

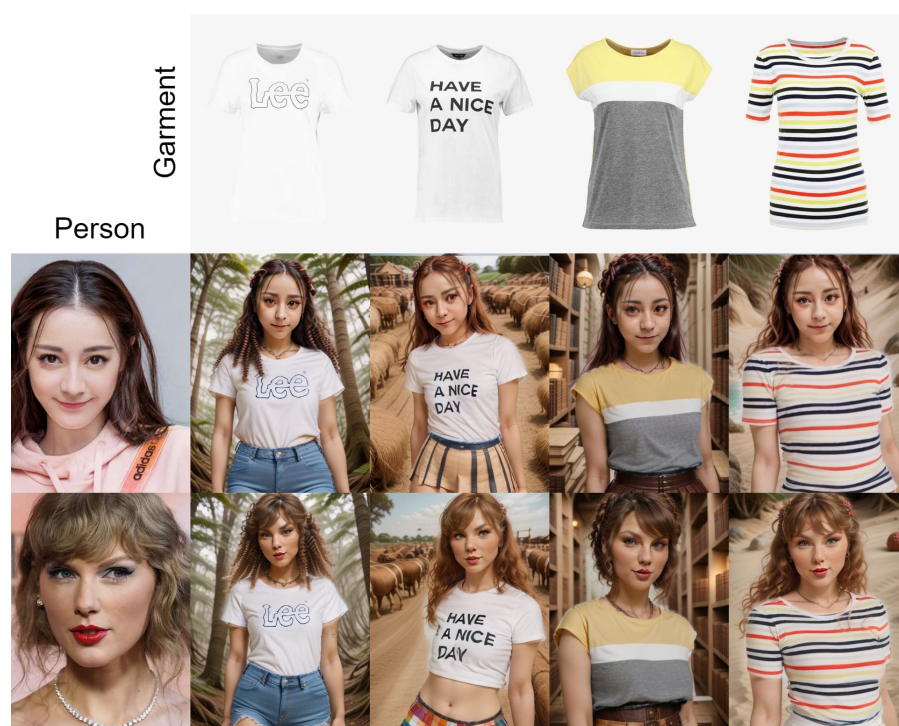


Fig. 16: Our text-to-image try-on results integrated with IP-Adapter.

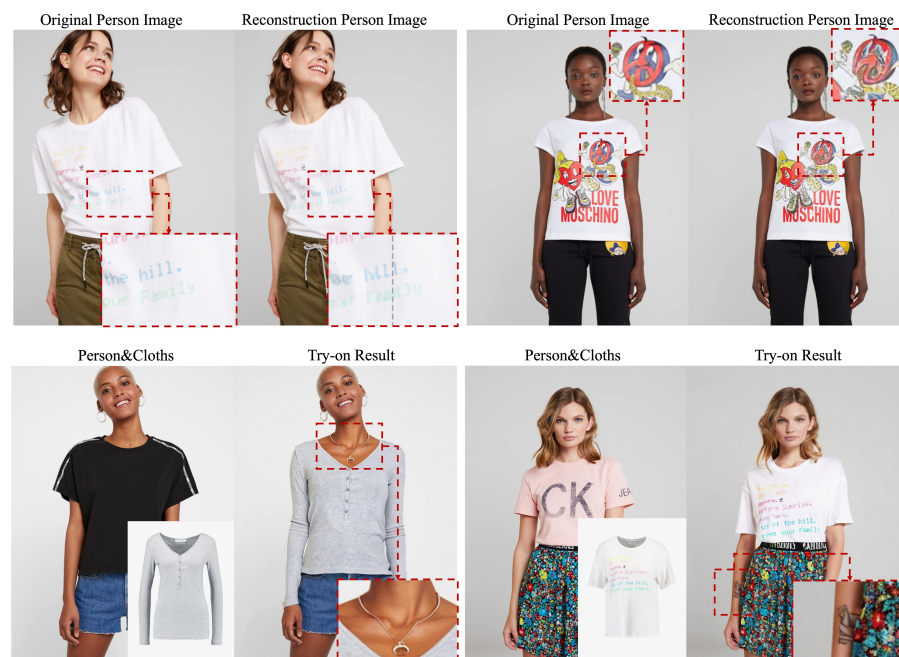


Fig. 17: Bad cases.