

# SGDiff: A Style Guided Diffusion Model for Fashion Synthesis

Zhengwentai Sun  
The Hong Kong Polytechnic University  
Hong Kong, Hong Kong  
zhengwt.sun@connect.polyu.hk

Yanghong Zhou  
The Hong Kong Polytechnic University  
Hong Kong, Hong Kong  
yanghong.zhou@connect.polyu.hk

Honghong He  
The Hong Kong Polytechnic University  
Hong Kong, Hong Kong  
21039747r@connect.polyu.hk

P. Y. Mok\*  
The Hong Kong Polytechnic University  
Hong Kong, Hong Kong  
Laboratory for Artificial Intelligence in Design  
Hong Kong, Hong Kong  
tracy.mok@polyu.edu.hk



Figure 1: Synthesizing clothing images with specific attributes and artistic styles: text-only synthesis results vs SGDiff results.

## ABSTRACT

This paper reports on the development of a **novel style guided diffusion model (SGDiff)** which overcomes certain weaknesses inherent in existing models for image synthesis. The proposed SGDiff combines image modality with a pretrained text-to-image diffusion model to facilitate creative fashion image synthesis. It addresses the limitations of text-to-image diffusion models by incorporating supplementary style guidance, substantially reducing training costs, and overcoming the difficulties of controlling synthesized styles with text-only inputs. This paper also introduces a new dataset – SG-Fashion, specifically designed for fashion image synthesis applications, offering high-resolution images and an extensive range of garment categories. By means of comprehensive ablation study, we examine the application of classifier-free guidance to a variety of conditions and validate the effectiveness of the proposed model for generating fashion images of the desired

categories, product attributes, and styles. The contributions of this paper include a novel classifier-free guidance method for multi-modal feature fusion, a comprehensive dataset for fashion image synthesis application, a thorough investigation on conditioned text-to-image synthesis, and valuable insights for future research in the text-to-image synthesis domain. The code and dataset are available at: <https://github.com/taited/SGDiff>.

## CCS CONCEPTS

• Computing methodologies → Computer vision tasks.

## KEYWORDS

fashion synthesis, style guidance, text-to-image, denoising diffusion probabilistic models

\*P. Y. Mok is the corresponding author (tracy.mok@polyu.edu.hk).

## ACM Reference Format:

Zhengwentai Sun, Yanghong Zhou, Honghong He, and P. Y. Mok. 2023. SGDiff: A Style Guided Diffusion Model for Fashion Synthesis. In *Proceedings of the 31st ACM International Conference on Multimedia (MM '23)*, October 29–November 3, 2023, Ottawa, ON, Canada. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3581783.3613806>

## 1 INTRODUCTION

Creative fashion synthesis, more specifically fashion image synthesis, holds enormous potential in the fashion industry, as it enables the generation of a diverse range of designs, which can inspire and foster innovation. By synthesizing fashion product images with specific attributes and styles, designers can rapidly explore a wide array of design concepts and ideas, reducing the time and resources needed for prototyping. For instance, Figure 1 (row (a)) shows that a single jumpsuit category may encompass numerous attributes, including sleeve length, neckline type, prints, and overall styles. The creative fashion synthesis may involve an arbitrary combination of these attributes and styles, as shown in rows (b) and (c) of Figure 1. Even though the existing fashion synthesis approaches [1, 6, 11, 17, 18, 21, 23, 43, 44, 49] can achieve promising results, they rely on different input modalities, such as category labels of textual modality, design sketches of visual modality, or other specialized inputs, which limit their ability to capture the large diversity of fashion product attributes and styles. Among the different input modalities, textual modality offers significant advantages in creative fashion synthesis since it provides a more accessible and flexible means to convey various garment categories and attribute information. Designers can easily communicate their design ideas by text descriptions, enabling a more streamlined design process and a broader exploration of garment variations.

The recent advancement in Denoising Diffusion Probabilistic Models (DDPM) [7, 15, 27, 38, 39] has radically improved text-to-image generation, achieving stunning results with reasonable semantics. DDPMs can generate images that align with the input text descriptions, whereas the process of image synthesis is formulated as recovering the target image from an initial noisy image. Nichol *et al.* [28] proposed GLIDE, a UNet-like structure for posterior probability estimation in the denoising process, to incorporate text conditions so as to control synthesis directions. Furthermore, Rombach *et al.* [33] investigated LDM model synthesizing high-resolution images with reasonable semantics using a Variational Autoencoder (VAE) to compress images into latent space and applying diffusion models to learn denoising in the latent space. Building on LDM [33], several studies [3, 13] explored detailed image editing using text descriptions as control. Even though these approaches can effectively synthesize images with reasonable and desired semantics, there are considerable training costs involved. For example, the GLIDE was trained on several hundred million text-image pairs while the LDM was trained on the LAION-5B [36] dataset. The huge model sizes limit the possible downstream applications. On the other hand, generating high-quality images that capture the essence of the desired design semantics, based only on text descriptions, is challenging due to the high dimensionality and variability in visual modality. The current methods could only control the synthesized results with simple descriptions, such as colors, but not visual styles like fabric textures, because describing different abstract styles with natural human languages is itself challenging.

In this paper, a novel approach called **Style-Guide Diffusion Model (SGDiff)** is developed to address the drawbacks of text-to-image diffusion models in fashion image synthesis. SGDiff is inspired by the old idiom *"A picture is worth a thousand words"* and that style information could be better described and conveyed by

images than by texts. Incorporating style guidance as control for image synthesis presents several challenges. Firstly, laborious data annotation is involved for selecting representative images as style guidance for the training purpose. Secondly, the alignment of features in different modalities remains a challenging problem when both style and text are simultaneously required to control a model, especially for domain specific semantics not being well covered by large paired image-text dataset. To circumvent the first challenge, we formulate the synthesis process as image reconstruction. By learning to reconstruct a garment from text descriptions and a randomly cropped image patch as style guidance, the proposed SGDiff can synthesize fashion garments that reflect both the text and style. To address the second issue, the proposed SGDiff model utilizes the image encoder of the Contrastive Language-Image Pre-Training (CLIP) model [31] to convert style images into semantic representations. Furthermore, a Skip Cross-Attention (SCA) module is specially designed and applied to integrate image modality with text modality. Such network design is very different from the existing CLIP guided methods [5, 25, 29], which align features in latent space using the distinct image and text encoders of CLIP and optimize the input latent variable as additional loss guidance. Existing methods suffer from the low extension ability for downstream applications and the optimization-based approach results in very slow image synthesis. Instead, with the help of the SCA module, we could fix the pretrained text-to-image diffusion and only fine tune the style (image) encoder and SCA module, significantly reducing the computational costs and addressing the multi-modal alignment problem. Moreover, most existing diffusion models only consider applying classifier-free guidance based on a single condition, whereas SGDiff explores the optimized way of applying multi-condition classifier-free guidance to the diffusion model. The key contributions of this paper are summarized as follows:

- i. A **new task** for creative fashion synthesis is addressed that both texts and style images are used to control the synthesis of fashion images under specific garment categories, attributes and styles.
- ii. SGDiff – a novel approach is developed that integrates image modality to a pretrained text-to-image diffusion model, enabling creative fashion synthesis with style guidance. To the best of our knowledge, this is **the first network proposal** integrating CLIP with the classifier-free guidance approach for modality fusion aiming toward conditioned image generation.
- iii. With the innovative network design, a **new network training strategy** is presented that significantly reduces training costs, only requiring fine-tuning the image encoder and modality fusion module rather than the entire network.
- iv. A **SG-Fashion dataset** is specifically constructed, which features high-resolution images and covers a wide range of garment categories. The proposed method has been validated both on this SG-Fashion and the Polyvore datasets.
- v. This is **the first of its kind of thorough investigation for extending classifier-free guidance to multiple conditions**, providing valuable insights for future research in the text-to-image synthesis domain.

## 2 RELATED WORK

### 2.1 Fashion Synthesis

Fashion synthesis, an emerging research area within the broader field of computer vision and generative models, concentrates on generating and manipulating fashion-related images, such as clothing and accessories as well as fashion models. Virtual try-on (VTON) has generated considerable attention in some recent studies [6, 11, 17, 21, 23, 43], which typically employ human parsing maps and pose estimation techniques to transfer textures from a desired garment onto a target person. Although these VTON approaches successfully synthesize consistent clothing attributes, they primarily focus on human-centric scenarios.

Several recent studies have investigated garment-centric fashion synthesis, with the aim to generate novel and diverse clothing items. For example, Jiang *et al.* [18] developed FashionG to transfer styles onto a garment without changing its original image content. Other researchers [8, 44, 49] explored the synthesis of compatible fashion based on a given garment image as a query. *These aforementioned studies are all using visual modality input as control for image synthesis, their ability to control the detail attributes of the generated fashion is rather limited.*

Text-to-image fashion synthesis remains relatively unexplored compared to other fashion synthesis approaches. Zhu *et al.* [51] proposed a method that uses textual descriptions to edit images of garments worn by humans. Zhang *et al.* [47] developed an ARMANI model for fashion synthesis based on multi-modal inputs including text descriptions and edge or regional detail in image modality. *Although the above approaches successfully enable control over the synthesized garments, they generally fail to achieve detailed control of the synthesized textures or styles.*

### 2.2 CLIP Model Guided Modality Fusion

The CLIP model, introduced by OpenAI [31], has revolutionized the field of computer vision by leveraging the power of large-scale transformers trained on both images and text. One of the main strengths of the CLIP model is its zero-shot learning capability, namely no learning is needed, which allows it to handle new tasks without requiring any task-specific fine-tuning. Its zero-shot capability has been exploited in various applications, such as image classification [9, 46], object detection [37, 41], and semantic segmentation [24, 48, 50].

CLIP models have been integrated with generative models like GANs [25, 35] and VQ-VAEs [5] to produce impressive results in various tasks, from text-to-image synthesis to image editing. For example, StyleCLIP [29] utilizes a pretrained StyleGAN [20] and the CLIP model to align image and text features within the style space. VQGAN-CLIP [5] uses CLIP as additional guidance to control the generation direction in pretrained generative model. FuseDream [25] is a training-free method integrating the latent generation space with CLIP embeddings. DALL-E [32] combines the CLIP model with a discrete VAE to generate high-quality images from textual descriptions. *All these models adopt a training-free pipeline and treat the CLIP model as a gradient guidance to interpret the generation of latent space. Although these methods could integrate pretrained generation models with CLIP for text-to-image synthesis, they synthesize every image as a separate optimization process, which*

*are computationally costly, and they fail to capture domain-specific text descriptions.*

### 2.3 Text-to-Image Diffusion Models

Diffusion models have recently emerged as a powerful branch of generative models, demonstrating their superior capabilities of handling image, text, audio as well as other modalities of data [22, 26, 28, 33, 40]. These models aim to learn the data distribution by performing a Markov chain, simulating the data generation process in reverse [7, 15, 27, 38, 39].

Despite the many research studies are focusing on synthesizing high-resolution images using diffusion models, there is a growing body of research that is interested in more controlled synthesis. Hertz *et al.* [13] investigated a Prompt-to-Prompt mechanism of text-to-image generation, where text features activate feature maps through cross-modal attention. InstructPix2Pix [3] combines the large pretrained language model GPT3 [4] and the state-of-the-art text-to-image LDM [33] model to synthesize a dataset for text-driven image editing. *Although these methods can synthesize images with corresponding semantics, they are trained on large open-domain datasets and have difficulty in capturing terms specific to the fashion domain.* Recently, Textual Inversion [10] and DreamBooth [34] can adapt pre-trained diffusion models with new styles. *Model retraining is, however, needed for every new style.*

## 3 METHOD

### 3.1 Preliminaries

Diffusion models utilize a Markov chain process, motivated by non-equilibrium thermodynamics, to simulate forward diffusion process. Given an image  $\mathbf{x}_0 \sim q(\mathbf{x})$ , the forward diffusion process adds small amount of Gaussian noises to  $\mathbf{x}_0$  in  $\mathcal{T}$  steps, producing a set of noisy images  $\mathbf{x}_1, \dots, \mathbf{x}_{\mathcal{T}}$ . The step sizes  $\beta_t$  are controlled by a variance schedule  $\{\beta_t \in (0, 1)\}_{t=1}^{\mathcal{T}}$ :

$$q(\mathbf{x}_{1:\mathcal{T}} | \mathbf{x}_0) = \prod_{t=1}^{\mathcal{T}} q(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad (1a)$$

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}). \quad (1b)$$

By reparameterization, let  $\alpha_t = 1 - \beta_t$  and  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ , the forward process can sample  $\mathbf{x}_t$  at arbitrary timestep  $t$  directly from  $\mathbf{x}_0$ :

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}). \quad (2)$$

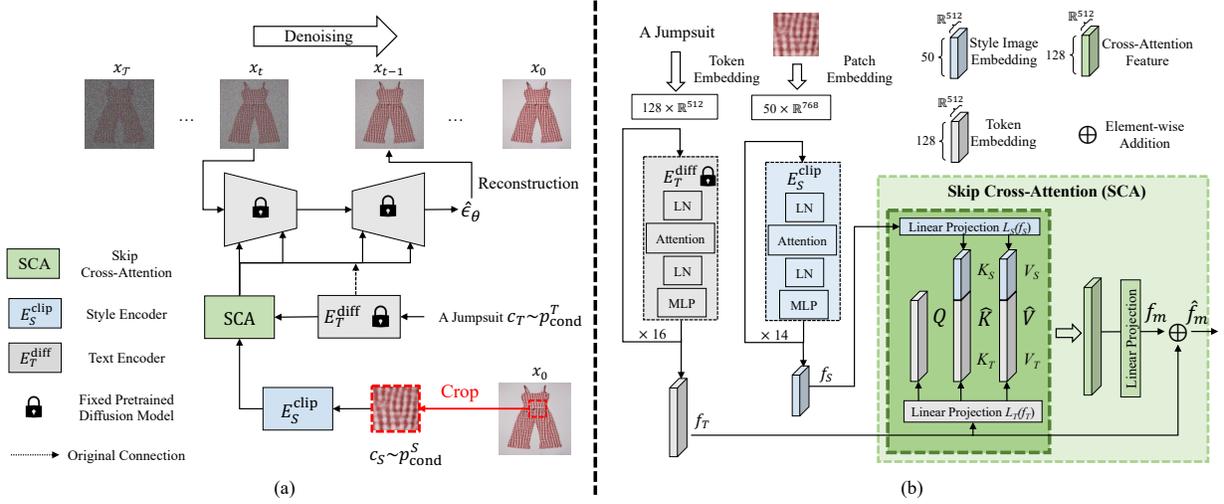
When  $\mathcal{T} \rightarrow \infty$ , the image  $\mathbf{x}_0$  will be diffused to a standard Gaussian noise  $\mathbf{x}_{\mathcal{T}} \sim \mathcal{N}(0, \mathbf{I})$ . Given a Gaussian noise, a neural network model is then learned to approximate the conditional probabilities to reverse the diffusion process  $p_{\theta}$  as follows:

$$p_{\theta}(\mathbf{x}_{0:\mathcal{T}}) = p(\mathbf{x}_{\mathcal{T}}) \prod_{t=1}^{\mathcal{T}} p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) \quad (3a)$$

$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t)), \quad (3b)$$

where  $\boldsymbol{\mu}_{\theta}$  and  $\boldsymbol{\Sigma}_{\theta}$  are the approximated mean and variance of the reversed Gaussian distribution. By simplifying  $\boldsymbol{\Sigma}_{\theta}$  as constant  $\beta_t$ ,  $\boldsymbol{\mu}_{\theta}$  is tractable [15]:

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_t \right). \quad (4)$$



**Figure 2: The proposed Style-Guided Diffusion Model (SGDiff) network (a) an overview and (b) detail model: SGDiff takes two inputs, a text condition ( $c_T$ ) for garment attributes and a style condition ( $c_S$ ) for style guidance, and leverages the Skip Cross-Attention (SCA) module and a pretrained CLIP image encoder for efficient training and resource utilization.**

With  $x_t$  known during training, the network  $\epsilon_\theta$  is reparameterized to predict noise  $\epsilon_t$  from input  $x_t$  at time step  $t$  with this simplified objective [16]:

$$\mathcal{L}_t^{\text{simple}} = \mathbb{E}_{t \sim [1, T], x_0, \epsilon_t} \left[ \left\| \epsilon_t - \epsilon_\theta \left( \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \epsilon_t, t \right) \right\|^2 \right]. \quad (5)$$

For brevity,  $\epsilon_\theta(x_t, t)$  is denoted as  $\epsilon_\theta(x_t)$  hereafter in this article.

### 3.2 SGDiff Overview

SGDiff aims to achieve detailed control over synthesized fashion images in terms of both correct garment attributes and garment textures (styles). Controlling detailed garment textures using natural language is challenging, therefore, the proposed SGDiff, as illustrated in Figure 2, takes two inputs: a text condition ( $c_T$ ) describing the garment attributes and a style condition ( $c_S$ ) guiding the synthesized garment texture. The text encoder  $E_T^{\text{diff}}$  of the diffusion model encodes the semantic representation  $f_T$ , and the style encoder  $E_S^{\text{clip}}$  of a pretrained CLIP model encodes the style representation  $f_S$ . The diffusion network  $\epsilon_\theta$  estimates the noise  $\hat{\epsilon}_t$  as follows:

$$\hat{\epsilon}_t = \epsilon_\theta \left( x_t, E_T^{\text{diff}}(c_T), E_S^{\text{clip}}(c_S) \right). \quad (6)$$

To avoid labor-intensive data annotation, the conditioned image synthesis is formulated as an image reconstruction task, as shown in Figure 2(a), in which a randomly image patch cropped from the garment image is taken as style condition  $c_S$ , the model is then trained to reconstruct garment according to the style guidance  $c_S$ .

To achieve efficient training, we have the pre-trained text-to-image diffusion model fine-tuned on a domain-specific dataset using text as input condition, according to a classifier-free guidance approach [16]. Next, by fixing the diffusion network parameters, we optimize the specially designed SCA module and fine-tune a pretrained image encoder  $E_S^{\text{clip}}$  with multiple conditions of text

description and style guidance, which will be discussed in detail in Section 3.5.

### 3.3 Skip Cross-Attention Module

Figure 2(b) illustrates the process of integrating two different modalities, namely text description of garment attributes  $c_T$  and image of style guidance  $c_S$ , in the proposed SGDiff model. The integration of the two input modalities is achieved through the specially designed Skip Cross-Attention (SCA) module.

Both encoders,  $E_T^{\text{diff}}$  and  $E_S^{\text{clip}}$ , employ transformer-based structures and the output features  $f_T \in \mathbb{R}^{128 \times 512}$  and  $f_S \in \mathbb{R}^{50 \times 512}$  represent two modalities of input. Such aligned features of  $f_T$  and  $f_S$  enable easy integration of the two representations by attention mechanism [42]. To do so, the semantic representation  $f_T$  is linearly projected into query and key-value pairs:

$$Q, K_T, V_T = L_T(f_T), \quad (7)$$

where  $L_T$  represents linear projection, and query  $Q$  and key-value pairs  $K_T, V_T$  all have size  $\mathbb{R}^{128 \times 512}$ . The style representation  $f_S$  is projected into key-value pairs only:

$$K_S, V_S = L_S(f_S). \quad (8)$$

The style key-value pairs are concatenated with text key-value pairs:

$$\hat{K} = K_S (+) K_T \quad \text{and} \quad \hat{V} = V_S (+) V_T, \quad (9)$$

where (+) denotes length-wise concatenation.

Specifically, the semantic representation  $f_T$  is chosen as query  $Q$  because it provides key attribute information for garment synthesis. With  $f_T$  as query, style representation  $f_S$  is aligned with the garment attributes in order to improve the quality of the synthesized images. The cross-attention is implemented by integrating the key-value pairs from both modalities as follows:

$$f_m = \text{Attention}(Q, \hat{K}, \hat{V}) = \text{softmax} \left( \frac{Q\hat{K}^T}{\sqrt{d_k}} \right) \hat{V}. \quad (10)$$

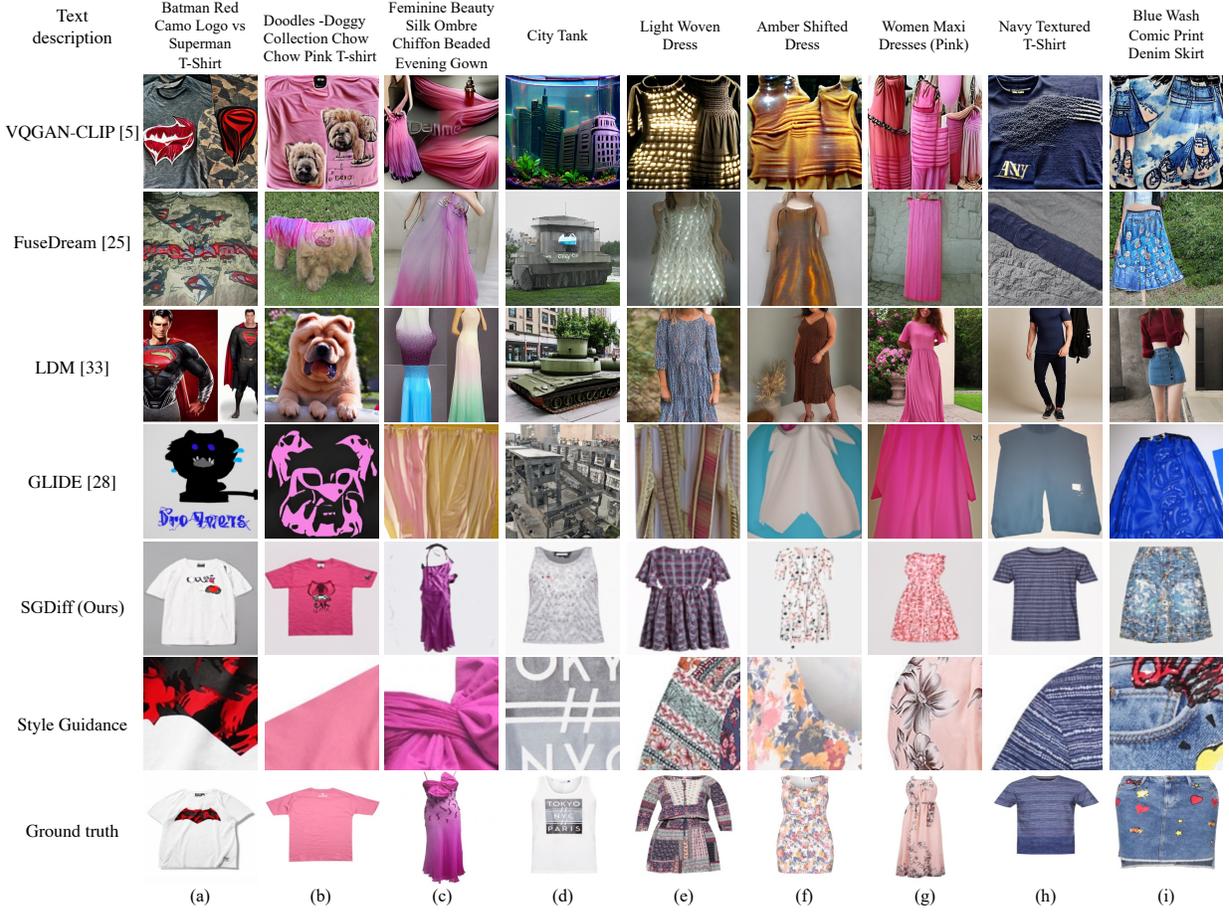


Figure 3: Qualitative comparison of SGDiff with state-of-the-art (SOTA) approaches. The 2nd and 3rd rows illustrate the results of CLIP-based methods of VQGAN-CLIP [5] and FuseDream [25], while the 4th and 5th rows illustrate the results of diffusion-based methods of LDM [33] and GLIDE [28]. The 6th row illustrates SGDiff’s ability to incorporate style images (the 7th row) into text conditions (the 1st row), successfully synthesizing garments with the desired textures.

Finally, the skip connection is applied, as shown in Figure 2:

$$\hat{f}_m = f_m + f_t. \quad (11)$$

The SCA module enables effective integration of text and image modalities, allowing the SGDiff model to control the synthesized texture without any reduction in semantic control.

### 3.4 Training Objectives

As discussed in Section 3.1, diffusion models implicitly learn to reconstruct an image from Gaussian noise. The network  $\epsilon_\theta$  estimates the noise in the current input noisy image  $x_t$ . The training objective of DDPM (Eq. (5)), however, does not address condition constraints explicitly. Therefore, SGDiff employs perceptual loss, in addition to Eq. (5), to govern image synthesis. To this end, the reconstructed image  $\hat{x}_0$  is obtained at every time step  $t$ , according to the estimated noise  $\hat{\epsilon}_t$  by Eq. (6):

$$\hat{x}_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} \left( x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}_t \right). \quad (12)$$

The Perceptual Loss [19] is then calculated by:

$$\mathcal{L}_t^{\text{perc}} = \mathbb{E}_m \|\psi_m(\hat{x}_0) - \psi_m(x_0)\|_2, \quad (13)$$

where  $\psi_m$  denotes the  $m$ -th layer of VGG. Following [19], the layers of *relu1\_2*, *relu2\_2*, *relu3\_2*, *relu4\_2*, and *relu5\_2* are used in Eq. (13). The overall training objective with Perceptual Loss, adapted from [27], is as follows:

$$\mathcal{L} = \lambda_s \mathcal{L}_t^{\text{simple}} + \mathcal{L}_t^{\text{vlb}} + \lambda_p \mathcal{L}_t^{\text{perc}}, \quad (14)$$

where  $\lambda_s$  and  $\lambda_p$  are balancing weights for the corresponding losses.

### 3.5 Multi-Modal Conditions

Classifier-free guidance [16] has obvious advantages over classifier guidance [7] for conditioned generation with DDPMs. For more flexible control, the proposed SGDiff also adopts classifier-free guidance approach [16], in which the model  $\epsilon_\theta$  is trained with conditional state  $c$  and unconditional state  $\emptyset$  according to a certain probability  $c \sim p_{\text{cond}}$ :

$$\hat{\epsilon}_\theta(x_t, c) = \epsilon_\theta(x_t, \emptyset) + s [\epsilon_\theta(x_t, c) - \epsilon_\theta(x_t, \emptyset)]. \quad (15)$$

Nevertheless, the above approach (Eq. (15)) does not address more complex situation where conditions are multiple, happen in different combinations at varied probabilities. Until recently, InstructPix2Pix [3] suggested different weights for two conditions:

$$\begin{aligned} \hat{\epsilon}_\theta(x_t, c_1, c_2) = & \epsilon_\theta(x_t, \emptyset, \emptyset) \\ & + s_1 [\epsilon_\theta(x_t, c_1, \emptyset) - \epsilon_\theta(x_t, \emptyset, \emptyset)] \\ & + s_2 [\epsilon_\theta(x_t, c_1, c_2) - \epsilon_\theta(x_t, c_1, \emptyset)], \end{aligned} \quad (16)$$

where  $s_1$  and  $s_2$  indicate the weight scale of condition  $c_1 \sim p_{\text{cond}}^1$  and  $c_2 \sim p_{\text{cond}}^2$ , respectively. In [3], however, it was not discussed either the order of  $c_1$  and  $c_2$  or the weight scales  $s_1$  and  $s_2$ .

In the current task, Eq. (16) is applied by setting the two conditions as  $c_T$  and  $c_S$ . The SGDiff is subjected to two conditions with independent conditional probability  $p_{\text{cond}}^S = 0.8$  and  $p_{\text{cond}}^T = 0.8$ . In model training, like all text-to-image diffusion models, the unconditional state  $\emptyset$  of textual condition  $c_T$  is set to padding token. The unconditional state  $\emptyset$  of style guidance  $c_S$  is done by inputting a blank (background only) patch image.

**Background masking:** Apart from inputting a blank image patch as unconditional state, the background color in RGB space may also appear in the foreground. To avoid confusion, we mask the background pixel values to -255 to distinguish them from the normal RGB values. Such masking technique allows the model to focus more on the foreground texture. The effectiveness of such background masking setting will be evaluated in Section 4.4.

**Condition order and weight scales:** In order to explore the effect of the condition order, by setting  $c_1 = c_S$  and  $c_2 = c_T$ , alternatively  $c_1 = c_T$  and  $c_2 = c_S$ , in Eq. (16), and  $s_T = 1$ , this will result in

$$\begin{aligned} \hat{\epsilon}_\theta(x_t, c_S, c_T) = & (s_S - 1) [\epsilon_\theta(x_t, c_S, \emptyset) - \epsilon_\theta(x_t, \emptyset, \emptyset)] \\ & + \epsilon_\theta(x_t, c_S, c_T) \end{aligned} \quad (17a)$$

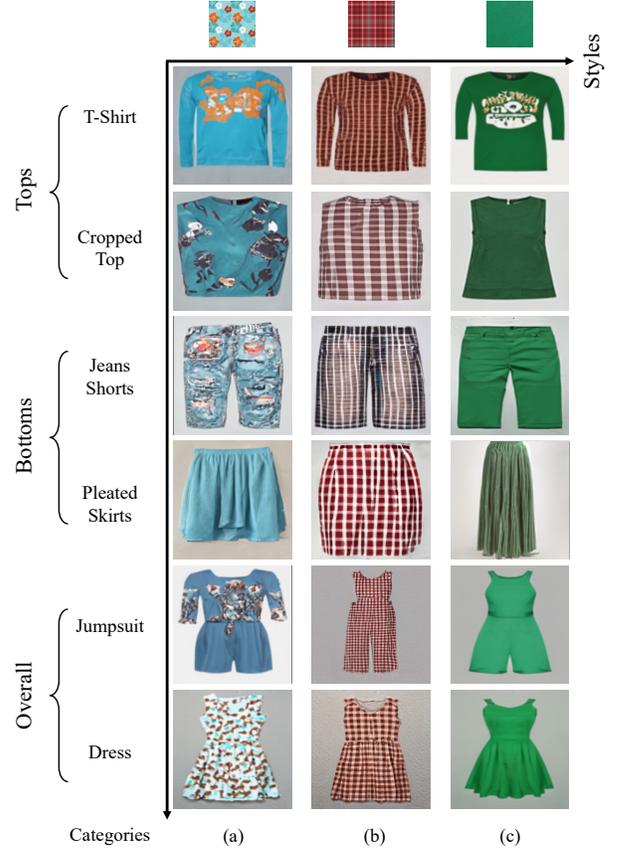
$$\begin{aligned} \hat{\epsilon}_\theta(x_t, c_T, c_S) = & (s_S - 1) [\epsilon_\theta(x_t, c_T, c_S) - \epsilon_\theta(x_t, c_T, \emptyset)] \\ & + \epsilon_\theta(x_t, c_T, c_S). \end{aligned} \quad (17b)$$

In our implementation, the model  $\epsilon_\theta$  takes  $c_S$  and  $c_T$  simultaneously, the two terms  $\epsilon_\theta(x_t, c_S, c_T)$  and  $\epsilon_\theta(x_t, c_T, c_S)$  are therefore equivalent. Comparing Eq. (17a) with (17b), thus  $[\epsilon_\theta(x_t, c_S, \emptyset) - \epsilon_\theta(x_t, \emptyset, \emptyset)] = [\epsilon_\theta(x_t, c_T, c_S) - \epsilon_\theta(x_t, c_T, \emptyset)]$ . It implies that if the style condition and text condition are independent, the condition order will not have a significant impact on the image generation. Moreover, the weight scale serves to adjust the influence of style guidance. When  $s_S > s_T$  (i.e.  $s_S > 1$  when  $s_T = 1$ ), it introduces a positive conditioned direction to the denoising processing, emphasizing the influence of condition is guiding the synthesis. The multi-condition synthesis will be further evaluated in Section 4.4.

## 4 EXPERIMENTS

### 4.1 Datasets and Implementation Details

In this study, we prepared a SG-Fashion dataset with 17,000 fashion product images downloaded from e-commerce websites including ASOS, Uniqlo and H&M. We set aside 1,700 images as the test set. The dataset covers 72 product categories, encompassing most types of garment items. Since our SGDiff does not rely on textural descriptions, we use the original product titles as text descriptions. Apart from the SG-Fashion dataset, we also experimented on the publicly available dataset of Polyvore [12] with the same settings.



**Figure 4: Illustration of SGDiff’s capability to synthesize garments across various categories and styles, using style guidance of different colors.**

GLIDE [28] was adopted as the backbone text-to-image diffusion model, which uses a low-resolution generation model for size  $64 \times 64$  and a super-resolution model to upsample the generated low-resolution image to the size of  $256 \times 256$ . We fine tuned the generation model and directly employed the super-resolution model as the pretrained text-to-image model. For the pretrained CLIP image encoder, we chose vision transformer of ViT/32. To speed up the synthesis process, we adopted DDIM [38] scheduler with 100 sampling steps for all diffusion-based models.

The backbone model (GLIDE) was fine tuned on the domain-specific dataset that the AdamW optimizer was used with a learning rate of  $1e^{-4}$ , and the model was optimized for 235,000 iterations. Due to GPU limitations, we set the batch size to 8 and trained the GLIDE on a single RTX 3090 GPU. We also used AdamW but with a learning rate of  $1e^{-5}$  for training the SGDiff with 50,000 iterations for all experiments on a single RTX 3090 GPU. In terms of the SCA module, we adopted multi-head attention with 4 heads. In all experiments, we set  $\lambda_s = 1$  and  $\lambda_p = 0.001$  in Eq. (14). Since the training of SGDiff fixes the parameters of the pretrained backbone, we can use a larger batch size of 16. For SGDiff training, we cropped a single texture patch from the foreground. To ensure this cropped patch provides sufficient style information, we applied BASNet [30] for background masking.

## 4.2 Qualitative Evaluation

The qualitative evaluation compares the SGDiff results with several SOTA text-to-image generation methods, including LDM [33] and GLIDE [28] for diffusion-based methods, and FuseDream [25] and VQGAN-CLIP [5] for CLIP-guided GAN-based methods. All selected SOTA methods have zero-shot capability. Figure 3 presents a comprehensive qualitative comparison of these methods. Generally speaking, FuseDream and LDM could synthesize garments in most cases, while VQGAN-CLIP and GLIDE could only synthesize fabrics. The proposed SGDiff could successfully implement the fashion synthesis with desired clothing category and style. Specifically, when synthesizing a garment with complex text descriptions (see examples in columns (a), (b), and (c)), the other methods tend to ignore the key message but capture part of the semantics like Batman logo, pink doggy, or silk, while SGDiff tends to synthesize clothing and consider the style guidance to control the synthesized textures. Moreover, semantic confusion is one of main challenges in text-to-image synthesis. For instance, ‘Tank’ refers to a specific type of upper clothing in the fashion domain. Column (d) of Figure 3 shows that both the diffusion-based and CLIP-based approaches have difficulty in capturing domain-specific semantics, while only SGDiff could synthesize a tank garment with specified styles. The other columns present cases when offering textual descriptions like amber, light and pink, although the other SOTA methods could synthesize clothing with textures that are similar to the descriptions, they show greater differences to the ground truth images comparing to SGDiff. In conclusion, SGDiff is suitable for fashion synthesis since it could capture the garment category and desired styles.

In addition to the comparative analysis, Figure 4 illustrates the innovative capability of SGDiff in synthesizing garments across various categories and styles. With style guidance images under different color schemes, SGDiff effectively transfers styles from the guidance images to the synthesized garments, meeting the condition of garment attributes. Figure 4 shows a range of synthesized fashion under specific color scheme in each column, offering valuable inspiration for innovative fashion design. When conditioned generation are out of the training set, SGDiff can still exhibit a remarkable generative capability by successfully blending different condition combinations, e.g., the jeans shorts with red check and green patterns showed in columns (b) and (c) are not existed in the training data. Moreover, the style guidance appears in interesting variations in the generated fashion. These results highlight the versatility and robustness of the SGDiff model in the realm of fashion design.

## 4.3 Metrics and Quantitative Evaluation

Table 1 shows the quantitative evaluation, in which three metrics, including FID [14], LPIPS [45] and CLIP-Score (CS) [31], are used to assess and compare the performance of SGDiff with other SOTA methods. FID and LPIPS measure the distance in feature space, with FID focusing on the overall distribution statistics of the generated/synthesized images and the ground truths, while LPIPS computes the distance between each pair of synthesized image and the corresponding ground truth, lower the FID and LPIPS values higher the image quality. In contrast, the CLIP-score measures the

**Table 1: Quantitative evaluation and comparison of various SOTA methods**

Datasets	SG-Fashion			Polyvore		
	LPIPS ↓	FID ↓	CS ↑	LPIPS ↓	FID ↓	CS ↑
VQGAN-CLIP [5]	0.7364	95.84	22.20	0.7122	68.01	<b>39.65</b>
FuseDream [25]	0.7067	60.44	<b>38.03</b>	0.7032	<b>41.94</b>	38.53*
LDM [33]	0.7158	85.73	31.66*	0.7214	59.79	31.89
GLIDE [28]	0.6921	78.7	23.72	0.7164	63.85	23.28
Ground Truth	-	-	29.13	-	-	29.88
Baseline	0.5772*	36.13*	27.31	0.6637*	43.5	26.24
SGDiff (Ours)	<b>0.4474</b>	<b>32.06</b>	27.53	<b>0.6369</b>	41.98*	27.33

the best results are in **bold**, and the second best results are indicated with \*.

**Table 2: Consumption of synthesizing an image with resolution of  $256 \times 256$  on a RTX 3090 GPU**

	VQGAN-CLIP	FuseDream	LDM	GLIDE	Ours
Time	62s	171s	5.9s	9s	9.8s
Memory	5686M	9296M	6570M	5550M	5986M

semantic correspondence, namely the cosine similarity between synthesized images and their corresponding text descriptions, with higher scores indicating better alignment.

As shown in Table 1, SGDiff model performs the best in terms of LPIPS, comparing to other SOTA methods on both SG-Fashion and Polyvore datasets. SGDiff’s FID value is also the lowest for SG-Fashion dataset and only slightly lower than FuseDream for Polyvore dataset by 0.04%. This demonstrates that the SGDiff model can generate better images fulfilling the conditions without sacrificing the image quality. The CS of the SGDiff is higher than GLIDE and the baseline (i.e. GLIDE being fine tuned on the datasets), but lower than FuseDream and LDM, because FuseDream optimizes the BigGAN-256 [2] latent space using CLIP guidance and LDM leverages a vast text-to-image dataset consisting of billions of examples. Nevertheless, these methods did not consider the integration of the text feature and image feature for image generation, they indeed did not perform well in LPIPS and FID.

Table 2 compares the model memory and average time cost for synthesizing an image of size  $256 \times 256$  on a RTX 3090 GPU. As shown, the running time of the SGDiff model is much shorter than that of VQGAN-CLIP and FuseDream. Although the running time of the SGDiff model is slightly longer than LDM, the memory consumption is lower. Compared to the baseline, the increases in time and memory are relatively insignificant because we only fine tune the image encoder and modality fusion module. In summary, the SGDiff can be trained without much memory and can generate an image with good quality based on text and style conditions within 10 seconds on RTX 3090.

## 4.4 Ablation Study

Ablation study was conducted to evaluate the effect of each component of the proposed SGDiff on SG-Fashion dataset.

**Effectiveness of the SCA:** As demonstrated in Table 3, the comparison between the element-wise addition of features and the cross-attention (CA) method shows that CA is significantly

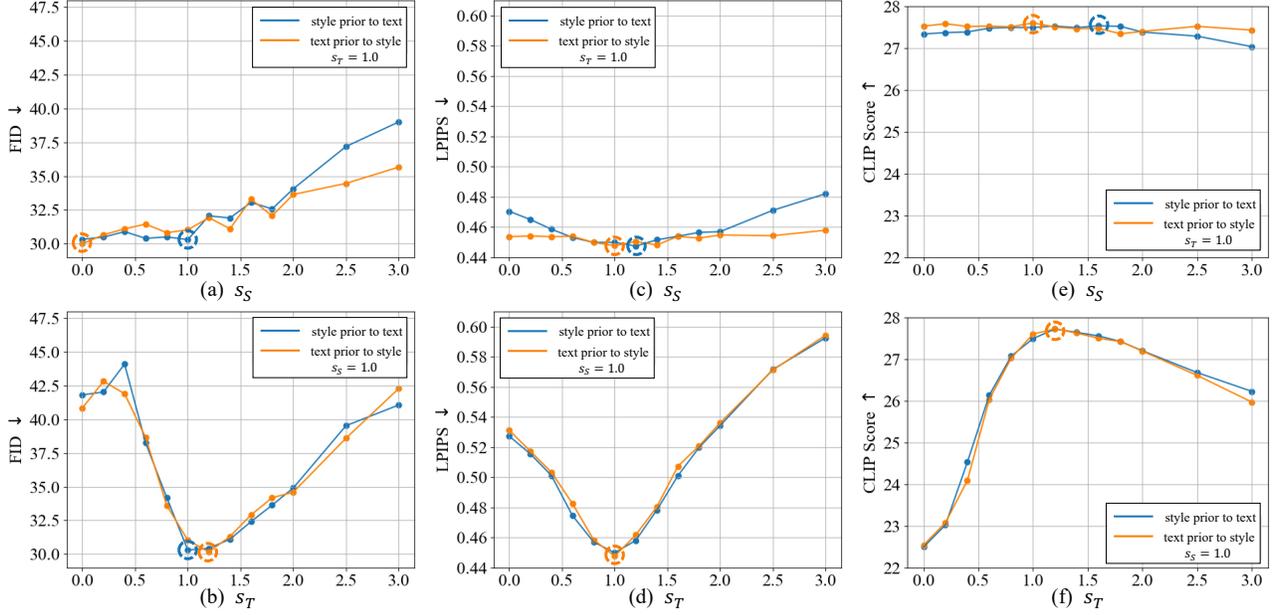


Figure 5: Ablation study on the impact of style and text guidance on the performance of SGDiff in terms of (a) and (b) for FID, (c) and (d) for LPIPS and (e) and (f) for CLIP-score. We set one conditional weight varies in range of  $[0, 0.2, 0.4, 0.6, 0.8, 1.0, 1.2, 1.4, 1.6, 1.8, 2.0, 2.5, 3.0]$  while the other conditional weight is fixed at 1.0.

Table 3: Ablation experiments on modality fusion methods and classifier-free approaches.

Classifier-free	Mask	Modality fusion	LPIPS ↓	FID ↓	CS ↑
Eq. (15)		$\oplus^1$	0.6833	42.63	25.63
Eq. (15)		CA <sup>2</sup>	0.5650	38.88	25.39
Eq. (15)		SCA	0.5607	39.21	25.98
Eq. (15)	✓	SCA	0.5695	37.22	26.06
Eq. (16)	✓	SCA	<b>0.4474</b>	<b>32.06</b>	<b>27.53</b>

<sup>1</sup>  $\oplus$  refers to an element-wise addition operation, where the features  $f_T$  and  $f_S$  are projected onto the same dimension before operation;

<sup>2</sup> CA indicates SCA module without skip connection, w.r.t. Eq. (10) without Eq. (11).

more effective in improving LPIPS and FID scores. However, it has the downside of causing a decline in semantic information, as CS decreases. To address this issue, the SCA module with skip connections was used. As shown in the third row of the table, SCA leads to improvements in both LPIPS and CS scores, demonstrating its ability to improve the similarity between synthesized images and ground truth images.

**Effectiveness of background masking:** As shown in Table 3, after applying background masking, the FID value decreases by 1.99 and the CS remains almost the same. This demonstrates that background masking is beneficial to improve image quality. The reason for slightly increased LPIPS is that LPIPS is sensitive to perceptual information, the lack of background may degrade LPIPS metric. However, the fashion synthesis task only focuses on the synthesized foreground, and the background could be easily removed.

**The orders and weights for different conditions:** Figure 5 displays the relationship between FID, LPIPS and CS with different conditional weights and order settings. The trend of setting text

prior to style is similar to setting style prior to text, indicating little impact on results with fixed  $s_S = 1$  and varying  $s_T$ . In addition, it can be seen from Figure 5 that the optimal values (see the circled dots of Figure 5) of  $s_S$  and  $s_T$  are almost in the range of 1.0 to 1.6. More specifically, we choose the setting of  $s_S = 1.2$ ,  $s_T = 1.0$ , and *style prior to text* as optimal. This setting achieves the best LPIPS which is important in controlling synthesized styles. The numerical results are shown in the last row of Table 3

## 5 CONCLUSIONS AND FUTURE WORK

This paper has reported on the development of a novel style guided diffusion model (SGDiff), overcoming inherent weaknesses in existing diffusion models for image synthesis. The proposed SGDiff has demonstrated its effectiveness in incorporating style guidance into pretrained text-to-image diffusion models. Without relying on large amounts of labelled data or computing resources, SGDiff is capable of achieving promising control over the synthesized textures, making it a valuable contribution to the field. As a future work, we plan to expand upon the capabilities of SGDiff by focusing on more detailed control over various aspects of the synthesized textures, such as color themes, patterns, and materials. By refining these controls, we aim to further improve the utility and applicability of the proposed model in diverse applications and domains.

## ACKNOWLEDGMENTS

The work described in this paper is supported in part by the Innovation and Technology Commission of Hong Kong under grant ITP/028/21TP and by the Laboratory for Artificial Intelligence in Design (Project Code: RP1-1) under InnoHK Research Clusters, Hong Kong Special Administrative Region.

## REFERENCES

- [1] Kenan E Ak, Joo Hwee Lim, Jo Yew Tham, and Ashraf A Kassim. 2019. Attribute manipulation generative adversarial networks for fashion images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 10541–10550.
- [2] Andrew Brock, Jeff Donahue, and Karen Simonyan. 2019. Large Scale GAN Training for High Fidelity Natural Image Synthesis. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=B1xsqj09Fm>
- [3] Tim Brooks, Aleksander Holynski, and Alexei A. Efros. 2023. InstructPix2Pix: Learning to Follow Image Editing Instructions. In *CVPR*.
- [4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [5] Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hallahan, Louis Castricato, and Edward Raff. 2022. VQGAN-CLIP: Open Domain Image Generation and Editing with Natural Language Guidance. In *Computer Vision – ECCV 2022*, Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner (Eds.). Springer Nature Switzerland, Cham, 88–105.
- [6] Aiyu Cui, Daniel McKee, and Svetlana Lazebnik. 2021. Dressing in order: Recurrent person image generation for pose transfer, virtual try-on and outfit editing. In *Proceedings of the IEEE/CVF international conference on computer vision*. 14638–14647.
- [7] Prafulla Dhariwal and Alexander Quinn Nichol. 2021. Diffusion Models Beat GANs on Image Synthesis. In *Advances in Neural Information Processing Systems*, A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (Eds.). <https://openreview.net/forum?id=AAWuCVzaVt>
- [8] Yajuan Ding, PY Mok, Yunshan Ma, and Yi Bin. 2023. Personalized fashion outfit generation with user coordination preference learning. *Information Processing & Management* 60, 5 (2023), 103434.
- [9] Sepideh Esmailpour, Bing Liu, Eric Robertson, and Lei Shu. 2022. Zero-shot out-of-distribution detection based on the pre-trained model clip. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 36. 6568–6576.
- [10] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and Daniel Cohen-or. 2023. An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion. In *The Eleventh International Conference on Learning Representations*. <https://openreview.net/forum?id=NAQvF08TcyG>
- [11] Yuying Ge, Yibing Song, Ruimao Zhang, Chongjian Ge, Wei Liu, and Ping Luo. 2021. Parser-free virtual try-on via distilling appearance flows. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 8485–8493.
- [12] Xintong Han, Zuxuan Wu, Yu-Gang Jiang, and Larry S Davis. 2017. Learning fashion compatibility with bidirectional lstms. In *Proceedings of the 25th ACM international conference on Multimedia*. 1078–1086.
- [13] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or. 2023. Prompt-to-Prompt Image Editing with Cross-Attention Control. In *The Eleventh International Conference on Learning Representations*. <https://openreview.net/forum?id=CDixzkzeyb>
- [14] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc. [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/8a1d694707eb0fe65871369074926d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/8a1d694707eb0fe65871369074926d-Paper.pdf)
- [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851.
- [16] Jonathan Ho and Tim Salimans. 2021. Classifier-Free Diffusion Guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*. <https://openreview.net/forum?id=qw8AKxfYbI>
- [17] Bingwen Hu, Ping Liu, Zhedong Zheng, and Mingwu Ren. 2022. SPG-VTON: Semantic Prediction Guidance for Multi-Pose Virtual Try-on. *IEEE Transactions on Multimedia* 24 (2022), 1233–1246. <https://doi.org/10.1109/TMM.2022.3143712>
- [18] Shuhui Jiang, Jun Li, and Yun Fu. 2022. Deep Learning for Fashion Style Generation. *IEEE Transactions on Neural Networks and Learning Systems* 33, 9 (2022), 4538–4550. <https://doi.org/10.1109/TNNLS.2021.3057892>
- [19] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. 2016. Perceptual losses for real-time style transfer and super-resolution. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*. Springer, 694–711.
- [20] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2020. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 8110–8119.
- [21] Bo-Kyeong Kim, Geonmin Kim, and Soo-Young Lee. 2020. Style-Controlled Synthesis of Clothing Segments for Fashion Image Manipulation. *IEEE Transactions on Multimedia* 22, 2 (2020), 298–310. <https://doi.org/10.1109/TMM.2019.2929000>
- [22] Yichong Leng, Zehua Chen, Junliang Guo, Haohe Liu, Jiawei Chen, Xu Tan, Danilo Mandic, Lei He, Xiangyang Li, Tao Qin, et al. 2022. Binauralgrad: A two-stage conditional diffusion probabilistic model for binaural audio synthesis. *Advances in Neural Information Processing Systems* 35 (2022), 23689–23700.
- [23] Kathleen M Lewis, Srivatsan Varadharajan, and Ira Kemelmacher-Shlizerman. 2021. Tryongan: Body-aware try-on via layered interpolation. *ACM Transactions on Graphics (TOG)* 40, 4 (2021), 1–10.
- [24] Feng Liang, Bichen Wu, Xiaoliang Dai, Kungpeng Li, Yanan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. 2022. Open-vocabulary semantic segmentation with mask-adapted clip. *arXiv preprint arXiv:2210.04150* (2022).
- [25] Xingchao Liu, Chengyue Gong, Lemeng Wu, Shujian Zhang, Hao Su, and Qiang Liu. 2021. Fusedream: Training-free text-to-image generation with improved clip+ gan space optimization. *arXiv preprint arXiv:2112.01573* (2021).
- [26] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. 2022. SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations. In *International Conference on Learning Representations*. [https://openreview.net/forum?id=aBsCjcPu\\_tE](https://openreview.net/forum?id=aBsCjcPu_tE)
- [27] Alexander Quinn Nichol and Prafulla Dhariwal. 2021. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*. PMLR, 8162–8171.
- [28] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2021. GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models. In *International Conference on Machine Learning*. PMLR, 16784–16804.
- [29] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. 2021. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2085–2094.
- [30] Xuebin Qin, Zichen Zhang, Chenyang Huang, Chao Gao, Masood Dehghan, and Martin Jagersand. 2019. Basnet: Boundary-aware salient object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 7479–7489.
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- [32] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. In *International Conference on Machine Learning*. PMLR, 8821–8831.
- [33] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10684–10695.
- [34] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2023. 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*. 22500–22510.
- [35] Peter Schaldenbrand, Zhixuan Liu, and Jean Oh. 2022. StyleCLIPDraw: Coupling Content and Style in Text-to-Drawing Translation. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, Lud De Raedt (Ed.). International Joint Conferences on Artificial Intelligence Organization, 4966–4972. <https://doi.org/10.24963/ijcai.2022/688> AI and Arts.
- [36] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. 2022. LAION-5B: An open large-scale dataset for training next generation image-text models. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*. <https://openreview.net/forum?id=M3Y74vmsMcY>
- [37] Hengcan Shi, Munawar Hayat, Yicheng Wu, and Jianfei Cai. 2022. Proposal-CLIP: unsupervised open-category object proposal generation via exploiting clip cues. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9611–9620.
- [38] Jiaming Song, Chenlin Meng, and Stefano Ermon. 2021. Denoising Diffusion Implicit Models. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=StlgarCHLP>
- [39] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. 2021. Maximum likelihood training of score-based diffusion models. *Advances in Neural Information Processing Systems* 34 (2021), 1415–1428.
- [40] Xuan Su, Jiaming Song, Chenlin Meng, and Stefano Ermon. 2023. Dual Diffusion Implicit Bridges for Image-to-Image Translation. In *The Eleventh International Conference on Learning Representations*. <https://openreview.net/forum?id=5HLotVVGDe>
- [41] Zhu Teng, Yani Duan, Yan Liu, Baopeng Zhang, and Jianping Fan. 2021. Global to local: Clip-LSTM-based object detection from remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing* 60 (2021), 1–13.
- [42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc. [https://proceedings.neurips.cc/paper\\_files/paper/](https://proceedings.neurips.cc/paper_files/paper/)

- 2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf
- [43] Jun Xu, Yuanyuan Pu, Renzan Nie, Dan Xu, Zhengpeng Zhao, and Wenhua Qian. 2021. Virtual Try-on Network With Attribute Transformation and Local Rendering. *IEEE Transactions on Multimedia* 23 (2021), 2222–2234. <https://doi.org/10.1109/TMM.2021.3070972>
- [44] Cong Yu, Yang Hu, Yan Chen, and Bing Zeng. 2019. Personalized fashion design. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9046–9055.
- [45] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [46] Renrui Zhang, Wei Zhang, Rongyao Fang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. 2022. Tip-adapter: Training-free adaption of clip for few-shot classification. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXV*. Springer, 493–510.
- [47] Xujie Zhang, Yu Sha, Michael C. Kampffmeyer, Zhenyu Xie, Zequn Jie, Chengwen Huang, Jianqing Peng, and Xiaodan Liang. 2022. ARMANI: Part-Level Garment-Text Alignment for Unified Cross-Modal Fashion Design. In *Proceedings of the 30th ACM International Conference on Multimedia (Lisboa, Portugal) (MM '22)*. Association for Computing Machinery, New York, NY, USA, 4525–4535. <https://doi.org/10.1145/3503161.3548230>
- [48] Chong Zhou, Chen Change Loy, and Bo Dai. 2022. Extract free dense labels from clip. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVIII*. Springer, 696–712.
- [49] Dongliang Zhou, Haijun Zhang, Qun Li, Jianghong Ma, and Xiaofei Xu. 2022. COutfitGAN: Learning to Synthesize Compatible Outfits Supervised by Silhouette Masks and Fashion Styles. *IEEE Transactions on Multimedia* (2022), 1–15. <https://doi.org/10.1109/TMM.2022.3185894>
- [50] Ziqin Zhou, Bowen Zhang, Yinjie Lei, Lingqiao Liu, and Yifan Liu. 2022. ZegCLIP: Towards Adapting CLIP for Zero-shot Semantic Segmentation. *arXiv preprint arXiv:2212.03588* (2022).
- [51] Shizhan Zhu, Raquel Urtasun, Sanja Fidler, Dahua Lin, and Chen Change Loy. 2017. Be your own prada: Fashion synthesis with structural coherence. In *Proceedings of the IEEE international conference on computer vision*. 1680–1688.