# StyleDiffusion: Controllable Disentangled Style Transfer via Diffusion Models

Zhizhong Wang[*], Lei Zhao[†], Wei Xing

College of Computer Science and Technology, Zhejiang University

{endywon, cszhl, wxing}@zju.edu.cn

## Abstract

*Content and style (C-S) disentanglement is a fundamental problem and critical challenge of style transfer. Existing approaches based on explicit definitions (e.g., Gram matrix) or implicit learning (e.g., GANs) are neither interpretable nor easy to control, resulting in entangled representations and less satisfying results. In this paper, we propose a new C-S disentangled framework for style transfer without using previous assumptions. The key insight is to explicitly extract the content information and implicitly learn the complementary style information, yielding interpretable and controllable C-S disentanglement and style transfer. A simple yet effective CLIP-based style disentanglement loss coordinated with a style reconstruction prior is introduced to disentangle C-S in the CLIP image space. By further leveraging the powerful style removal and generative ability of diffusion models, our framework achieves superior results than state of the art and flexible C-S disentanglement and trade-off control. Our work provides new insights into the C-S disentanglement in style transfer and demonstrates the potential of diffusion models for learning well-disentangled C-S characteristics.*

## 1. Introduction

Given a reference style image, *e.g.*, *Starry Night* by Vincent Van Gogh, style transfer aims to transfer its artistic style, such as colors and brushstrokes, to an arbitrary content target. To achieve such a goal, it must first properly separate the style from the content and then transfer it to another content. This raises two fundamental challenges: (1) "how to disentangle content and style (C-S)" and (2) "how to transfer style to another content".

To resolve these challenges, valuable efforts have been devoted. Gatys *et al*. [19] proposed *A Neural Algorithm of Artistic Style* to achieve style transfer, which *explicitly* defines the high-level features extracted from a pre-trained Convolutional Neural Network (CNN) (*e.g.*, VGG [76]) as

content, and the feature correlations (*i.e.*, Gram matrix) as style. This approach acquires visually stunning results and inspires a large number of successors [35, 30, 53, 1, 95]. Despite the successes, by diving into the essence of style transfer, we observed three problems with these approaches: (1) The C-S are not completely disentangled. Theoretically, the C-S representations are intertwined. For example, matching the content representation of an image may also match its Gram matrix, and vice versa. (2) What CNN learned is a black box rugged to interpret [97], which makes the C-S definitions [19] uninterpretable and hard to control. (3) The transfer process is modeled as a separate optimization of content loss and style loss [19], so there lacks a deep understanding of the relationship between C-S. These problems usually lead to unbalanced stylizations and disharmonious artifacts [6], as will be shown in later Fig. 3.

On the other hand, disentangled representation learning [27] provides other ideas to *implicitly* disentangle C-S, either supervised [47, 37] or unsupervised [9, 98]. For style transfer, Kotovenko *et al*. [45] utilized fixpoint triplet style loss and disentanglement loss to enforce a GAN [21]-based framework to learn separate C-S representations in an unsupervised manner. Similarly, TPFR [79] learned to disentangle C-S in latent space via metric learning and two-stage peer-regularization, producing high-quality images even in the zero-shot setting. While these approaches successfully enforce properties "encouraged" by the corresponding losses, they still have three main problems: (1) Well-disentangled models seemingly cannot be identified without supervision [57, 70], which means the unsupervised learning [45, 79] may not achieve truly disentangled C-S, as will be shown in later Fig. 3. (2) These approaches are all based on GANs and thus often confined to the GAN predefined domains, *e.g.*, a specific artist's style domain [75]. (3) The implicitly learned C-S representations are still black boxes that are hard to interpret and control [57].

Facing the challenges above, in this paper, we propose a new C-S disentangled framework for style transfer *without using previous assumptions* such as Gram matrix [19] or GANs [45]. Our key insight stems from the fact that the definition of an image's style is much more complex

---

than its content, *e.g.*, we can easily identify the content of a painting by its structures, semantics, or shapes, but it is intractable to define the style [67, 22, 38, 87]. Therefore, we can bypass such a dilemma by *explicitly* extracting the content information and *implicitly* learning its *complementary* style information. Since we strictly constrain style as the *complement* of content, the C-S can be completely disentangled, and the control of disentanglement has been transformed into the control of content extraction. It achieves both controllability and interpretability.

However, achieving plausible and controllable content extraction is also non-trivial because the contents extracted from the content images and style images should share the same content domain, and the details of the extracted contents should be easy to control. To this end, we resort to recent developed diffusion models [28, 78] and introduce a *diffusion-based style removal module* to smoothly dispel the style information of the content and style images, extracting the domain-aligned content information. Moreover, owing to the strong generative capability of diffusion models, we also introduce a *diffusion-based style transfer module* to better learn the disentangled style information of the style image and transfer it to the content image. The style disentanglement and transfer are encouraged via a simple yet effective *CLIP [68]-based style disentanglement loss*, which induces the transfer mapping of the content image's content to its stylization (*i.e.*, the stylized result) to be aligned with that of the style image's content to its stylization (*i.e.*, the style image itself) in the CLIP image space. By further coordinating with a *style reconstruction prior*, it achieves both generalized and faithful style transfer. We conduct comprehensive comparisons and ablation study to demonstrate the effectiveness and superiority of our framework. With the well-disentangled C-S, it achieves very promising stylizations with fine style details, well-preserved contents, and a deep understanding of the relationship between C-S.

In summary, our contributions are threefold:

- We propose a novel C-S disentangled framework for style transfer, which achieves more interpretable and controllable C-S disentanglement and higher-quality stylized results.
- We introduce diffusion models to our framework and demonstrate their effectiveness and superiority in controllable style removal and learning well-disentangled C-S characteristics.
- A new CLIP-based style disentanglement loss coordinated with a style reconstruction prior is introduced to disentangle C-S in the CLIP image space.

## 2. Related Work

**Neural Style Transfer (NST).** The pioneering work of Gatys *et al.* [19] has opened the era of NST [34]. Since

then, this task has experienced tremendous progress, including efficiency [35, 52, 90], quality [23, 89, 55, 10, 7, 1, 46, 83, 56, 6, 92, 32, 99, 12, 96, 84], generality [5, 30, 53, 65, 13, 33, 29, 85, 95, 59, 93], and diversity [80, 86, 88]. Despite these successes, the essence of these approaches is mostly based on the *explicitly* defined C-S representations, such as Gram matrix [19], which have several limitations as discussed in Sec. 1. In our work, we propose new disentangled C-S representations *explicitly* extracted or *implicitly* learned by diffusion models, achieving more effective style transfer and higher-quality results.

**Disentangled Representation Learning (DRL).** The task of DRL [27] aims at modeling the factors of data variations [51]. Earlier works used labeled data to factorize representations in a supervised manner [37]. Recently, unsupervised settings have been largely explored [42], especially for disentangling style from content [98, 31, 51, 40, 91, 45, 66, 70, 8, 48]. However, due to the dependence on GANs [21], their C-S disentanglement is usually restricted in the GAN pre-defined domains (*e.g.*, Van Gogh's style domain). Besides, disentanglement cannot be effectively achieved without providing sufficient data [57]. In contrast, our framework learns the disentangled style from a single style image, and the disentanglement can be easily achieved by providing only a few ($\sim$50) content images for training.

**Diffusion Models.** Diffusion models [77] such as denoising diffusion probabilistic models (DDPMs) [28, 63] have recently shown great success in image generation [78, 14, 17], image manipulation [62, 2, 41], and text-conditional synthesis [64, 74, 69, 71, 24, 4, 54]. These works have demonstrated the power of diffusion models to achieve higher-quality results than other generative models like VAEs [81], auto-regressive models [16], flows [44], and GANs [39]. Inspired by them, we introduce a diffusion-based style removal module and a style transfer module in our framework. These modules can smoothly remove the style information of images and better learn the recovery of it to achieve higher-quality style transfer results. *To the best of our knowledge, our work is the first to introduce diffusion models to the field of neural style transfer.*

## 3. Background

Denoising diffusion probabilistic models (DDPMs) [77, 28] are latent variable models that consist of two diffusion processes, *i.e.*, a forward diffusion process and a reverse diffusion process. The forward process is a fixed Markov Chain that sequentially produces a series of latents $x_1, ..., x_T$ by gradually adding Gaussian noises at each timestep $t \in [1, T]$:

$$q(x_t|x_{t-1}) := \mathcal{N}(\sqrt{1 - \beta_t}x_{t-1}, \beta_t\mathbf{I}), \quad (1)$$

where $\beta_t \in (0, 1)$ is a fixed variance schedule. An important property of the forward process is that given clean data
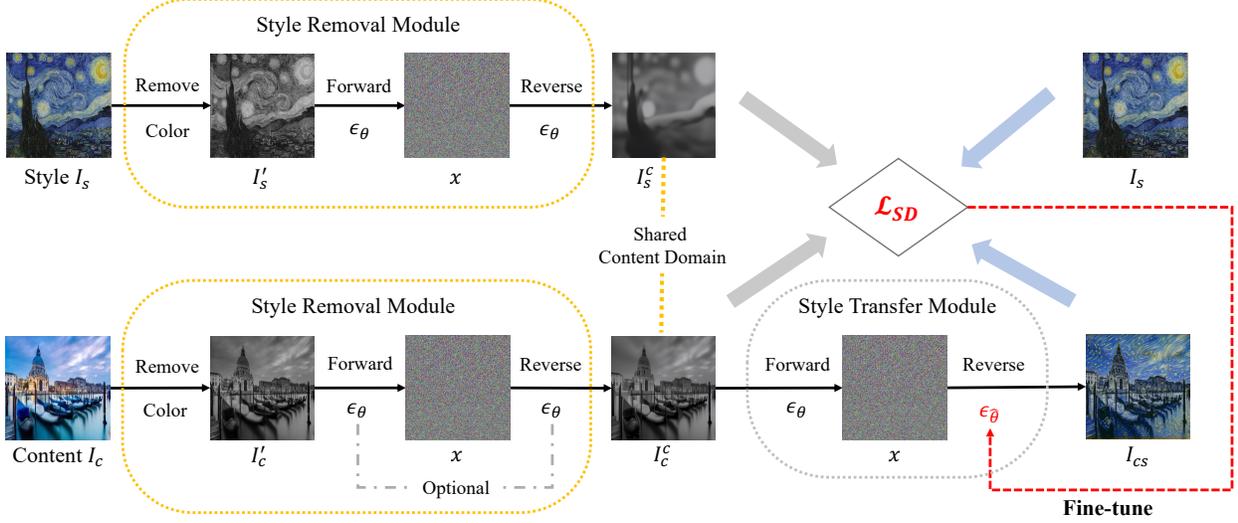
Figure 1. **Overview of our proposed StyleDiffusion.** The content image $I_c$ and style image $I_s$ are first fed into a diffusion-based style removal module to explicitly extract the domain-aligned content information. Then, the content of $I_c$ is fed into a diffusion-based style transfer module to obtain the stylized result $I_{cs}$. During training, we fine-tune the style transfer module via a CLIP-based style disentanglement loss $\mathcal{L}_{SD}$ coordinated with a style reconstruction prior (see details in Sec. 4.3, we omit it here for brevity) to implicitly learn the disentangled style information of $I_s$.

$x_0$, $x_t$ can be directly sampled as:

$$q(x_t|x_0) := \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1-\bar{\alpha}_t)\mathbf{I}),$$
$$x_t := \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1-\bar{\alpha}_t}\epsilon, \qquad (2)$$

where $\alpha_t := 1-\beta_t$ and $\bar{\alpha}_t := \prod_{s=0}^t \alpha_s$. Noise $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ has the same dimensionality as data $x_0$ and latent $x_t$.

The reverse process generates a reverse sequence by sampling the posteriors $q(x_{t-1}|x_t)$, starting from a Gaussian noise sample $x_T \sim \mathcal{N}(0, \mathbf{I})$. However, since $q(x_{t-1}|x_t)$ is intractable, DDPMs learn parameterized Gaussian transitions $p_\theta(x_{t-1}|x_t)$ with a learned mean $\mu_\theta(x_t, t)$ and a fixed variance $\sigma_t^2\mathbf{I}$ [28]:

$$p_\theta(x_{t-1}|x_t) := \mathcal{N}(\mu_\theta(x_t, t), \sigma_t^2\mathbf{I}), \qquad (3)$$

where $\mu_\theta(x_t, t)$ is the function of a noise approximator $\epsilon_\theta(x_t, t)$. Then, the reverse process can be expressed as:

$$x_{t-1} := \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon_\theta(x_t, t)\right) + \sigma_t\mathbf{z}, \qquad (4)$$

where $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$ is a standard Gaussian noise independent of $x_t$. $\epsilon_\theta(x_t, t)$ is learned by a deep neural network [72] through optimizing the following loss:

$$\min_\theta \| \epsilon_\theta(x_t, t) - \epsilon \|^2. \qquad (5)$$

Later, instead of using the fixed variances, Nichol and Dhariwal [63] presented a strategy for learning the variances. Song *et al.* [78] proposed DDIM, which formulates

an alternative non-Markovian noising process that has the same forward marginals as DDPM but allows a different reverse process:

$$x_{t-1} := \sqrt{\bar{\alpha}_{t-1}}f_\theta(x_t, t) + \sqrt{1-\bar{\alpha}_{t-1}-\sigma_t^2}\epsilon_\theta(x_t, t) + \sigma_t\mathbf{z}, \qquad (6)$$

where $f_\theta(x_t, t)$ is the predicted $x_0$ at timestep $t$ given $x_t$ and $\epsilon_\theta(x_t, t)$:

$$f_\theta(x_t, t) := \frac{x_t - \sqrt{1-\bar{\alpha}_t}\epsilon_\theta(x_t, t)}{\sqrt{\bar{\alpha}_t}}. \qquad (7)$$

Changing the choice of $\sigma_t$ values in Eq. (6) can achieve different reverse processes. Especially when $\sigma_t = 0$, which is called DDIM [78], the reverse process becomes a deterministic mapping from latents to images, which enables nearly perfect inversion [41]. Besides, it can also accelerate the reverse process with much fewer sampling steps [14, 41].

## 4. Method

Our task can be described as follows: given a style image $I_s$ and an arbitrary content image $I_c$, we want to first disentangle the content and style of them and then transfer the style of $I_s$ to the content of $I_c$. To achieve so, as stated in Sec. 1, our key idea is to explicitly extract the content information and then implicitly learn the *complementary* style information. Since our framework is built upon diffusion models [28, 78], we dub it *StyleDiffusion*.

Fig. 1 shows the overview of our StyleDiffusion, which consists of three key ingredients: I) a diffusion-based style

removal module, II) a diffusion-based style transfer module, and III) a CLIP-based style disentanglement loss coordinated with a style reconstruction prior. In the following subsections, we will introduce each of them in detail.

## 4.1. Style Removal Module

The style removal module aims at removing the style information of the content and style images, explicitly extracting the domain-aligned content information. Any reasonable content extraction operation can be used, depending on how the users define the content. For instance, users may want to use the structural outline as the content, so they can extract the outlines [36, 94] here. However, as discussed in Sec. 1, one challenge is *controllability* since the control of C-S disentanglement has been transformed into the control of content extraction. To this end, we introduce a diffusion-based style removal module to achieve both plausible and controllable content extraction.

Given an input image, *e.g.*, the style image $I_s$, since the color is an integral part of style [50], our style removal module first removes its color by a commonly used ITU-R 601-2 luma transform [20]. The obtained grayscale image is denoted as $I_s'$. Then, we leverage a pre-trained diffusion model [14] $\epsilon_\theta$ to remove the style details such as brushstrokes and textures of $I_s'$, extracting the content $I_s^c$. The insight is that the pre-trained diffusion model can help eliminate the domain-specific characteristics of input images and align them to the pre-trained domain [11, 41]. We assume that images with different styles belong to different domains, but their contents should share the same domain. Therefore, we can pre-train the diffusion model on a surrogate domain, *e.g.*, the photograph domain, and then use this domain to construct the contents of images. After pre-training, the diffusion model can convert the input images from diverse domains to the latents $x$ via the forward process and then inverse them to the photograph domain via the reverse process. In this way, the style characteristics can be ideally dispelled, leaving only the contents of the images.

Specifically, in order to obtain the results with fewer sampling steps and ensure that the content structures of the input images can be well preserved, we adopt the deterministic DDIM [78] sampling as the reverse process (Eq. (8)), and the ODE approximation of its reversal [41] as the forward process (Eq. (9)):

$$x_{t-1} = \sqrt{\bar{\alpha}_{t-1}} f_\theta(x_t, t) + \sqrt{1 - \bar{\alpha}_{t-1}} \epsilon_\theta(x_t, t), \quad (8)$$

$$x_{t+1} = \sqrt{\bar{\alpha}_{t+1}} f_\theta(x_t, t) + \sqrt{1 - \bar{\alpha}_{t+1}} \epsilon_\theta(x_t, t), \quad (9)$$

where $f_\theta(x_t, t)$ is defined in Eq. (7). The forward and reverse diffusion processes enable us to easily control the intensity of style removal by adjusting the number of return step $T_{remov}$ (see details in later Sec. 5.1). With the increase of $T_{remov}$, more style characteristics will be removed, and

the main content structures are retained, as will be shown in later Sec. 5.3. Note that for content images that are photographs, the diffusion processes are optional[1] since they are already within the pre-trained domain, and there is almost no style except the colors to be dispelled. The superiority of diffusion-based style removal against other operations, such as Auto-Encoder (AE) [53]-based style removal, can be found in *supplementary material (SM)*.

## 4.2. Style Transfer Module

The style transfer module aims to learn the disentangled style information of the style image and transfer it to the content image. A common generative model like AEs [30] can be used here. However, inspired by the recent great success of diffusion models [14, 41], we introduce a diffusion-based style transfer module, which can better learn the disentangled style information in our framework and achieve higher-quality and more flexible stylizations (see Sec. 5.3).

Given a content image $I_c$, denote $I_c^c$ is the content of $I_c$ extracted by the style removal module (Sec. 4.1). We first convert it to the latent $x$ using a pre-trained diffusion model $\epsilon_\theta$. Then, guided by a CLIP-based style disentanglement loss coordinated with a style reconstruction prior (Sec. 4.3), the *reverse process* of the diffusion model is fine-tuned ($\epsilon_\theta \rightarrow \epsilon_{\hat{\theta}}$) to generate the stylized result $I_{cs}$ referenced by the style image $I_s$. Once the fine-tuning is completed, *any content image can be manipulated into the stylized result with the disentangled style of the style image $I_s$*. To make the training easier and more stable, we adopt the deterministic DDIM forward and reverse processes in Eq. (8) and Eq. (9) during the fine-tuning. However, at inference, the stochastic DDPM [28] forward process (Eq. (2)) can also be used directly to help obtain diverse results [86] (Sec. 5.3).

## 4.3. Loss Functions and Fine-tuning

Enforcing the style transfer module (Sec. 4.2) to learn and transfer the disentangled style information should address two key questions: (1) "how to regularize the learned style is disentangled" and (2) "how to aptly transfer it to other contents". To answer these questions, we introduce a novel CLIP-based style disentanglement loss coordinated with a style reconstruction prior to train the networks.

**CLIP-based Style Disentanglement Loss.** Denote $I_c^c$ and $I_s^c$ are the respective contents of the content image $I_c$ and the style image $I_s$ extracted by the style removal module (Sec. 4.1). We aim to learn the disentangled style information of the style image $I_s$ *complementary* to its content $I_s^c$. Therefore, a straightforward way to obtain the disentangled style information is a direct subtraction:

$$D_s^{px} = I_s - I_s^c. \quad (10)$$

---

[1] Unless otherwise specified, we do not use the diffusion processes for content images in order to better maintain the content structures.
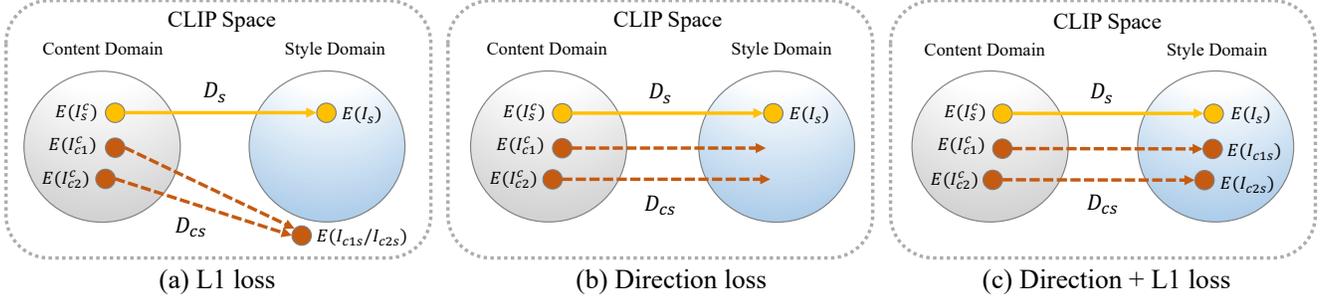
Figure 2. **Illustration of different loss functions** to transfer the disentangled style information. (a) L1 loss cannot guarantee the stylized results are within the style domain and may suffer from a collapse problem. (b) Direction loss aligns the disentangled directions but cannot realize accurate mappings. (c) Combining L1 loss and direction loss is able to achieve accurate one-to-one mappings from the content domain to the style domain.

However, the simple pixel differences do not contain meaningful semantic information, thus cannot achieve plausible results [19, 45]. To address this problem, we can formulate the disentanglement in a latent semantic space:

$$D_s = E(I_s) - E(I_s^c), \tag{11}$$

where $E$ is a well-pre-trained projector. Specifically, since $I_s$ and $I_s^c$ have similar contents but with different styles, the projector $E$ must have the ability to distinguish them in terms of the style characteristics. In other words, as we define that images with different styles belong to different domains, the projector $E$ should be able to distinguish the domains of $I_s$ and $I_s^c$. Fortunately, inspired by the recent vision-language model CLIP [68] that encapsulates knowledgeable semantic information of not only the photograph domain but also the artistic domain [18, 69, 49], we can use its image encoder as our projector $E$ off the shelf. The open-domain CLIP space here serves as a good metric space to measure the "style distance" between content and its stylized result. This "style distance" thus can be interpreted as the disentangled style information. Note that here the style is implicitly defined as the *complement* of content, which is fundamentally different from the Gram matrix [19] that is an explicit style definition independent of content (see comparisons in Sec. 5.3). The comparisons between CLIP space and other possible spaces can be found in *SM*.

After obtaining the disentangled style information $D_s$, the next question is how to properly transfer it to other contents. A possible solution is directly optimizing the L1 loss:

$$D_{cs} = E(I_{cs}) - E(I_c^c),$$
$$\mathcal{L}_{SD}^{L1} = \| D_{cs} - D_s \|, \tag{12}$$

where $I_{cs}$ is the stylized result, $D_{cs}$ is the disentangled style information of $I_{cs}$. However, as illustrated in Fig. 2 (a) and further validated in later Sec. 5.3, minimizing the L1 loss cannot guarantee the stylized result $I_{cs}$ is within the style domain of the style image $I_s$. It is because L1 loss only

minimizes the absolute pixel difference (*i.e.*, Manhattan distance); thus, it may produce stylized images that satisfy the Manhattan distance but deviate from the target style domain in the transfer direction. Besides, it may also lead to a collapse problem where a stylized output meets the same Manhattan distance with different contents in the latent space.

To address these problems, we can further constrain the disentangled directions as follows:

$$\mathcal{L}_{SD}^{dir} = 1 - \frac{D_{cs} \cdot D_s}{\| D_{cs} \| \| D_s \|}. \tag{13}$$

This direction loss aligns the transfer direction of the content image's content to its stylization (*i.e.*, the stylized result) with the direction of the style image's content to its stylization (*i.e.*, the style image itself), as illustrated in Fig. 2 (b). Collaborated with this loss, the L1 loss $\mathcal{L}_{SD}^{L1}$ thus can achieve accurate one-to-one mappings from contents in the content domain to their stylizations in the style domain, as illustrated in Fig. 2 (c).

Finally, our style disentanglement loss is defined as a compound of $\mathcal{L}_{SD}^{L1}$ and $\mathcal{L}_{SD}^{dir}$:

$$\mathcal{L}_{SD} = \lambda_{L1} \mathcal{L}_{SD}^{L1} + \lambda_{dir} \mathcal{L}_{SD}^{dir}, \tag{14}$$

where $\lambda_{L1}$ and $\lambda_{dir}$ are hyper-parameters set to 10 and 1 in our experiments. Since our style information is induced by the difference between content and its stylized result, we can deeply understand the relationship between C-S through learning. As a result, the style can be naturally and harmoniously transferred to the content, leading to better stylized images, as will be shown in later Fig. 3.

**Style Reconstruction Prior.** To fully use the prior information provided by the style image and further elevate the stylization effects, we integrate a style reconstruction prior into the fine-tuning of the style transfer module. Intuitively, given the content $I_s^c$ of the style image $I_s$, the style transfer module should be capable of recovering it to the original style image as much as possible. Therefore, we can define

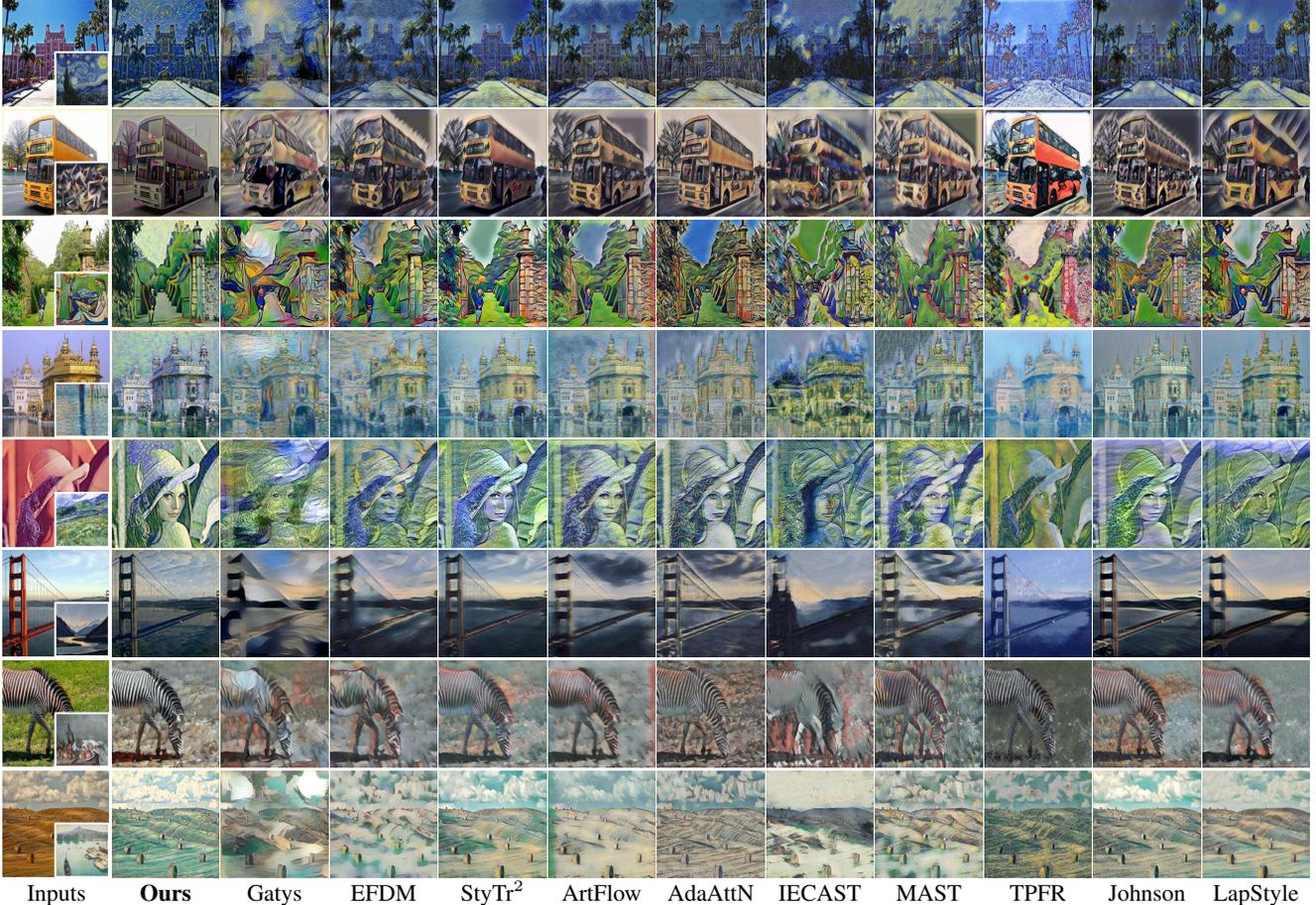| Inputs | **Ours** | Gatys | EFDM | StyTr$^2$ | ArtFlow | AdaAttN | IECAST | MAST | TPFR | Johnson | LapStyle |

Figure 3. **Qualitative comparisons** with state of the art. Zoom-in for better comparison. Please see more in *SM*.

a style reconstruction loss as follows:

$$\mathcal{L}_{SR} = \parallel I_{ss} - I_s \parallel, \tag{15}$$

where $I_{ss}$ is the stylized result given $I_s^c$ as content. We optimize it separately before optimizing the style disentanglement loss $\mathcal{L}_{SD}$. The detailed fine-tuning procedure can be found in *SM*. The style reconstruction prior helps our model recover the style information more sufficiently. It also provides a good initialization for the optimization of $\mathcal{L}_{SD}$, which helps the latter give full play to its ability, thus producing higher-quality results (see later Sec. 5.3).

## 5. Experimental Results

### 5.1. Implementation Details

We use ADM diffusion model [14] pre-trained on ImageNet [73] and adopt a fast sampling strategy [41]. Specifically, instead of sequentially conducting the diffusion processes until the last timestep $T$ (*e.g.*, 1000), we accelerate them by performing up to $T_{\{\cdot\}} < T$ (which is called return step), *i.e.*, $T_{remov} = 601$ for style removal and $T_{trans} =$

301 for style transfer. Moreover, as suggested by [41], we further accelerate the forward and reverse processes with fewer discretization steps, *i.e.*, $(S_{for}, S_{rev}) = (40, 40)$ ($S_{for}$ for forward process and $S_{rev}$ for reverse process) for style removal, and $(S_{for}, S_{rev}) = (40, 6)$ for style transfer. When fine-tuning or inference, we can adjust $T_{remov}$ or $T_{trans}$ to flexibly control the degree of style removal and C-S disentanglement, as will be shown in Sec. 5.3. To fine-tune the model for a target style image, we randomly sample 50 images from ImageNet as the content images. We use Adam optimizer [43] with an initial learning rate of 4e-6 and increase it linearly by 1.2 per epoch. All models are fine-tuned with 5 epochs. See more details in *SM*.

### 5.2. Comparisons with Prior Arts

We compare our StyleDiffusion against ten state-of-the-art (SOTA) methods [19, 95, 12, 1, 56, 6, 13, 79, 35, 55]. For fair comparisons, all these methods are fine-tuned or trained on the target styles similar to our approach.

**Qualitative Comparisons.** As can be observed in Fig. 3, due to the entangling of C-S representations, Gatys [19]

| | **Ours** | Gatys | EFDM | StyTr$^2$ | ArtFlow | AdaAttN | IECAST | MAST | TPFR | Johnson | LapStyle |
|---|---|---|---|---|---|---|---|---|---|---|---|
| SSIM ↑ | **0.672** | 0.311 | 0.316 | 0.537 | 0.501 | 0.542 | 0.365 | 0.392 | 0.536 | 0.634 | 0.657 |
| CLIP Score ↑ | **0.741** | 0.677 | 0.607 | 0.531 | 0.546 | 0.577 | 0.646 | 0.590 | 0.644 | 0.537 | 0.595 |
| Style Loss ↓ | 0.837 | **0.111** | 0.178 | 0.216 | 0.258 | 0.310 | 0.284 | 0.229 | 0.989 | 0.364 | 0.274 |
| User Study — Style | - | 43.1% | 41.2% | 39.3% | 36.4% | 37.2% | 33.8% | 39.1% | 14.5% | 42.8% | 47.3% |
| User Study — Overall | - | 26.0% | 38.1% | 44.0% | 34.2% | 43.9% | 32.7% | 32.2% | 22.6% | 43.4% | 46.2% |
| Training Time/h | ∼0.4 | - | ∼3 | ∼4 | ∼3 | ∼3 | ∼3 | ∼3 | ∼10 | ∼1 | ∼3 |
| Testing Time/s | 5.612 | 10.165 | 0.028 | 0.168 | 0.204 | 0.076 | 0.034 | 0.066 | 0.302 | 0.015 | 0.008 |

Table 1. **Quantitative comparisons** with state of the art. The training/testing time is measured with an Nvidia Tesla A100 GPU, and the testing time is averaged on images of size 512×512 pixels. ↑: Higher is better. ↓: Lower is better.

and EFDM [95] often produce unsatisfying results with distorted contents (*e.g.*, rows 1-3) and messy textures (*e.g.*, rows 4-8). StyTr$^2$ [12] and ArtFlow [1] improve the results by adopting more advanced networks [82, 44], but they may still produce inferior results with halo boundaries (*e.g.*, rows 2-3) or dirty artifacts (*e.g.*, rows 4-6). AdaAttN [56] performs per-point attentive normalization to preserve the content structures better, but the stylization effects may be degraded in some cases (*e.g.*, rows 1, 2, 4, and 5). IECAST [6] utilizes contrastive learning and external learning for style transfer, so fine-tuning it on a single style image would result in degraded results. MAST [13] uses multi-adaptation networks to disentangle C-S. However, since it still relies on the C-S representations of [19], the results usually exhibit messy textures and conspicuous artifacts. TPFR [79] is a GAN-based framework that learns to disentangle C-S in latent space. As the results show, it cannot recover correct style details and often generates deviated stylizations, which signifies that it may not learn truly disentangled C-S representations [57]. Like our method, Johnson [35] and LapStyle [55] also train separate models for each style. However, due to the trade-off between C-S losses of [19], they may produce less-stylized results or introduce unnatural patterns (*e.g.*, rows 1-6).

By contrast, our StyleDiffusion completely disentangles C-S based on diffusion models. Therefore, it can generate high-quality results with sufficient style details (*e.g.*, rows 1-4) and well-preserved contents (*e.g.*, rows 5-8). Compared with the previous methods that tend to produce mixed results of content and style, our approach can better consider the relationship between them. Thus, the stylizations are more natural and harmonious, especially for challenging styles such as cubism (*e.g.*, row 2) and oil painting (*e.g.*, rows 1, 3, 4, and 5).

**Quantitative Comparisons.** We also resort to quantitative metrics to better evaluate our method, as shown in Tab. 1. We collect 32 content and 12 style images to synthesize 384 stylized results and compute the average Structural Similarity Index (SSIM) [1] to assess the content similarity. To evaluate the style similarity, we calculate the CLIP image similarity score [68] and Style Loss [19, 30] between the style images and the corresponding stylized results. As

shown in Tab. 1, our method obtains the highest SSIM and CLIP Score while the Style Loss is relatively higher than other methods. It is because these methods are directly trained to optimize Style Loss. Nevertheless, the Style Loss achieved by our method is still comparable and lower than the GAN-based TPFR [79]. Furthermore, it is noteworthy that our method can also incorporate Style Loss to enhance the performance in this regard (see later Sec. 5.3).

**User Study.** As style transfer is highly subjective and CLIP Score and Style Loss are biased to the training objective, we additionally resort to user study to evaluate the style similarity and overall stylization quality. We randomly select 50 C-S pairs for each user. Given each C-S pair, we show the stylized results generated by our method and a randomly selected SOTA method side by side in random order. The users are asked to choose (1) which result transfers the style patterns better and (2) which result has overall better stylization effects. We obtain 1000 votes for each question from 20 users and show the percentage of votes that existing methods are preferred to ours in Tab. 1. The lower numbers indicate our method is more preferred than the competitors. As the results show, our method is superior to others in both style consistency and overall quality.

**Efficiency.** As shown in the bottom two rows of Tab. 1, our approach requires less training time than others as it is fine-tuned on only a few (∼50) content images. When testing, our approach is faster than the optimization-based method Gatys [19], albeit slower than the remaining feedforward methods due to the utilization of diffusion models. We discuss it in later Sec. 6, and more timing and resource details can be found in *SM*.

### 5.3. Ablation Study

**Control of C-S Disentanglement.** A prominent advantage of our StyleDiffusion is that we can flexibly control the C-S disentanglement by adjusting the content extraction of the style removal module (Sec. 4.1). Fig. 4 demonstrates the continuous control achieved by adjusting the return step $T_{remov}$ of the style removal module. As shown in the top row, with the increase of $T_{remov}$, more style characteristics are dispelled, and the main content structures are retained. Correspondingly, when more style is removed in the
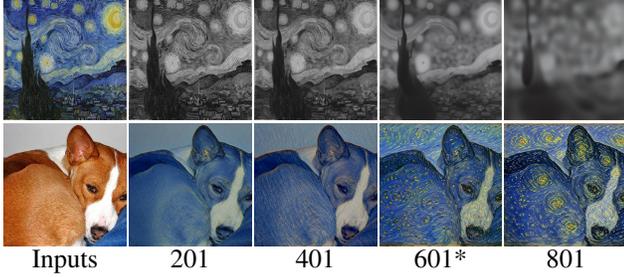
Figure 4. **Control of C-S disentanglement** by adjusting the return step $T_{remov}$ of the *style removal module*. The top row shows the extracted contents of the style image. The bottom row shows the corresponding stylized results. * denotes our default setting. Zoom-in for better comparison. *See SM for quantitative analyses.*



Figure 5. **Control of C-S trade-off** by adjusting the return step $T_{trans}$ of the *style transfer module*. The top row shows adjusting $T_{trans}$ at the **training** stage while fixing $T_{trans} = 301$ at the testing stage. The bottom row shows adjusting $T_{trans}$ at the **testing** stage while fixing $T_{trans} = 301$ at the training stage. * denotes our default setting. Zoom-in for better comparison. *See SM for quantitative analyses.*

top row, it will be aptly transferred to the stylized results in the bottom row, *e.g.*, the twisted brushstrokes and the star patterns. It validates that our method successfully separates style from content in a controllable manner and properly transfers it to other contents. Moreover, the flexible C-S disentanglement also makes our StyleDiffusion versatile for other tasks, such as photo-realistic style transfer (see *SM*).

**Superiority of Diffusion-based Style Transfer.** Although our style transfer module is not limited to the diffusion model, using it offers three main advantages: **(1)** *Flexible C-S trade-off control.* As shown in Fig. 5, we can flexibly control the C-S trade-off at both the training stage (top row) and the testing stage (bottom row) by adjusting the return step $T_{trans}$ of the diffusion model. With the increase of $T_{trans}$, more style characteristics are transferred, yet the content structures may be ruined (*e.g.*, the last column). When proper $T_{trans}$ is adopted, *e.g.*, $T_{trans} = 301$, the sweet spot can be well achieved. Interestingly, as shown in the last two columns of the bottom row, though the model is trained on $T_{trans} = 301$, we can extrapolate the style by using larger $T_{trans}$ (*e.g.*, 401) at the testing stage (but the results may be degraded when using too large $T_{trans}$, *e.g.*, 601). It provides a very flexible way for users to adjust the results according to their preferences. This property, however, cannot be simply achieved by using other models, *e.g.*, the widely used AEs [30, 53], since our framework does not involve any feature transforms [30, 53] or C-S losses trade-off [3]. **(2)** *Higher-quality stylizations.* Owing to the strong generative ability of the diffusion model, it can achieve higher-quality stylizations than other models. For comparison, we use the pre-trained VGG-AE [30, 49] as the style transfer module and fine-tune its decoder network for each style. As shown in column (b) of Fig. 6, though the results are still acceptable, they may produce distorted contents and inferior textures, clearly worse than the results generated by the diffusion model in column (a). This is also validated by the bottom quantitative scores. It signifies that the diffusion model can better learn the disentangled
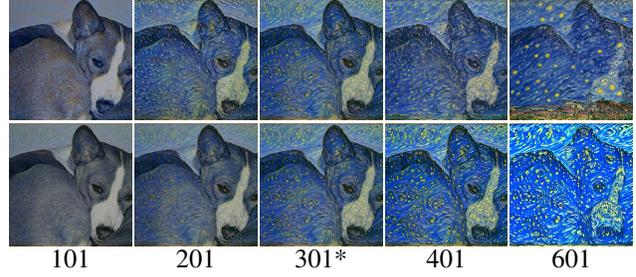


| | | (a) Diffusion | (b) AE |
|---|---|---|---|
| Style | Content | | |
| SSIM / CLIP Score: | | 0.672 / 0.741 | 0.526 / 0.702 |

Figure 6. **Diffusion-based vs. AE-based style transfer.**

content and style characteristics in our framework, helping produce better style transfer results. **(3)** *Diversified style transfer.* As mentioned in Sec. 4.2, during inference, we can directly adopt the stochastic DDPM [28] forward process (Eq. (2)) to obtain diverse results (see *SM*). The diverse results can give users endless choices to obtain more satisfactory results. However, using other models like AEs in our framework cannot easily achieve it [86].

**Loss Analyses.** To verify the effectiveness of each loss term used for fine-tuning our StyleDiffusion, we present ablation study results in Fig. 7 (a-d). **(1)** Using L1 loss $\mathcal{L}_{SD}^{L1}$ successfully transfers the cubism style like the blocky patterns in the top row, but the colors stray from the style images, especially in the bottom row. It is consistent with our earlier analyses in Sec. 4.3 that the L1 loss is prone to produce implausible results outside the style domain. **(2)** Adding direction loss $\mathcal{L}_{SD}^{dir}$ helps pull the results closer to the style domain. The textures are enhanced in the top row, and the colors are more plausible in the top and bottom rows. **(3)** By further coordinating with the style reconstruction prior $\mathcal{L}_{SR}$, the stylization effects are significantly elevated where the style information is recovered more suf-

| | | | | | | |
|---|---|---|---|---|---|---|
| Style | Content | (a) $\mathcal{L}_{SD}^{L1}$ | (b) + $\mathcal{L}_{SD}^{dir}$ | (c) + $\mathcal{L}_{SR}$* | (d) $\mathcal{L}_{SR}$ | (e) $\mathcal{L}_{Gram}$ | (f) + $\mathcal{L}_{SR}$ |
| SSIM / CLIP Score: | | 0.660 / 0.652 | 0.693 / 0.705 | 0.672 / 0.741 | 0.793 / 0.488 | 0.429 / 0.712 | 0.367 / 0.763 |

Figure 7. **Ablation study on loss functions.** * denotes our full model. Zoom-in for better comparison.

ficiently. It may be because it provides a good initialization for the optimization of $\mathcal{L}_{SD}^{L1}$ and $\mathcal{L}_{SD}^{dir}$, which helps them give full play to their abilities. As verified in Fig. 7 (d), using the style reconstruction alone cannot learn meaningful style patterns except for basic colors. All the above analyses are also supported by the bottom quantitative scores.

**Comparison with Gram Loss.** To further verify the superiority of our proposed losses, we replace them with the widely used Gram Loss [19, 30] in Fig. 7 (e-f). As can be observed, Gram Loss destroys the content structures severely, *e.g.*, the zebra head in the top row and the enlarged area in the bottom row. This is because it does not disentangle C-S and only matches the global statistics without considering the relationship between C-S. In contrast, our losses focus on learning the disentangled style information apart from the content, which is induced by the difference between content and its stylized result. Therefore, they can better understand the relationship between C-S, achieving more satisfactory results with fine style details and better-preserved contents, as validated by Fig. 7 (c) and the bottom quantitative scores. Furthermore, we also conduct comparisons between our proposed losses and Gram Loss [19, 30] on the AE baseline [30, 49] to eliminate the impact of diffusion models. As shown in Fig. 8 (a-b), our losses can achieve more satisfactory results than Gram Loss, which is consistent with the results in Fig. 7. Moreover, as shown in Fig. 8 (c), they can also be combined with Gram Loss to improve the performance on the Style Loss metric. However, it may affect the full disentanglement of C-S in our framework, which strays from our target and decreases the content preservation (see SSIM score in Fig. 8 (c)). Therefore, we do not incorporate Gram Loss in our framework by default.

## 6. Conclusion and Limitation

In this work, we present a new framework for more interpretable and controllable C-S disentanglement and style transfer. Our framework, termed *StyleDiffusion*, leverages



| | | | | |
|---|---|---|---|---|
| Style | Content | (a) $\mathcal{L}_{Gram}$ | (b) Ours | (c) Both |
| SSIM: | | 0.306 | **0.526** | 0.464 |
| CLIP Score: | | 0.586 | 0.702 | **0.728** |
| Style Loss: | | 0.263 | 0.732 | **0.231** |

Figure 8. **More loss function ablation study** on the AE baseline.

diffusion models to explicitly extract the content information and implicitly learn the complementary style information. A novel CLIP-based style disentanglement loss coordinated with a style reconstruction prior is also introduced to encourage the disentanglement and style transfer. Our method yields very encouraging stylizations, especially for challenging styles, and the experimental results verify its effectiveness and superiority against state of the art.

Currently, the framework still suffers from several limitations: (1) The model needs to be fine-tuned for each style, and arbitrary style transfer is left to our future work. (2) The efficiency is not fast enough due to the use of diffusion models. Further research in accelerating diffusion sampling would be helpful. (3) There are some failure cases analyzed in *SM*, which may help inspire future improvements. Moreover, our framework may also be applied to other image translation [31] or manipulation [66] tasks, and we would like to explore them in our future work.

# References

[1] Jie An, Siyu Huang, Yibing Song, Dejing Dou, Wei Liu, and Jiebo Luo. Artflow: Unbiased image style transfer via reversible neural flows. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 862–871, 2021. 1, 2, 6, 7, 14, 21, 22

[2] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18208–18218, 2022. 2

[3] Mohammad Babaeizadeh and Golnaz Ghiasi. Adjustable real-time style transfer. In *International Conference on Learning Representations (ICLR)*, 2019. 8

[4] Andreas Blattmann, Robin Rombach, Kaan Oktay, Jonas Müller, and Björn Ommer. Retrieval-augmented diffusion models. *Advances in Neural Information Processing Systems*, 35:15309–15324, 2022. 2

[5] Dongdong Chen, Lu Yuan, Jing Liao, Nenghai Yu, and Gang Hua. Stylebank: An explicit representation for neural image style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1897–1906, 2017. 2

[6] Haibo Chen, Lei Zhao, Zhizhong Wang, Huiming Zhang, Zhiwen Zuo, Ailin Li, Wei Xing, and Dongming Lu. Artistic style transfer with internal-external learning and contrastive learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 34:26561–26573, 2021. 1, 2, 6, 7, 14, 21, 22

[7] Haibo Chen, Lei Zhao, Zhizhong Wang, Huiming Zhang, Zhiwen Zuo, Ailin Li, Wei Xing, and Dongming Lu. Dualast: Dual style-learning networks for artistic style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 872–881, 2021. 2

[8] Haibo Chen, Lei Zhao, Huiming Zhang, Zhizhong Wang, Zhiwen Zuo, Ailin Li, Wei Xing, and Dongming Lu. Diverse image style transfer via invertible cross-space mapping. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14860–14869. IEEE Computer Society, 2021. 2

[9] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. *Advances in Neural Information Processing Systems (NeurIPS)*, 29, 2016. 1

[10] Jiaxin Cheng, Ayush Jaiswal, Yue Wu, Pradeep Natarajan, and Prem Natarajan. Style-aware normalized loss for improving arbitrary style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 134–143, 2021. 2

[11] Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, and Sungroh Yoon. Ilvr: Conditioning method for denoising diffusion probabilistic models. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14347–14356. IEEE, 2021. 4

[12] Yingying Deng, Fan Tang, Weiming Dong, Chongyang Ma, Xingjia Pan, Lei Wang, and Changsheng Xu. Stytr2: Image style transfer with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11326–11336, 2022. 2, 6, 7, 14, 21, 22

[13] Yingying Deng, Fan Tang, Weiming Dong, Wen Sun, Feiyue Huang, and Changsheng Xu. Arbitrary style transfer via multi-adaptation network. In *Proceedings of the 28th ACM International Conference on Multimedia (ACM MM)*, pages 2719–2727, 2020. 2, 6, 7, 14, 21, 22

[14] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems (NeurIPS)*, 34:8780–8794, 2021. 2, 3, 4, 6, 14

[15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2020. 17

[16] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12873–12883, 2021. 2

[17] Wan-Cyuan Fan, Yen-Chun Chen, DongDong Chen, Yu Cheng, Lu Yuan, and Yu-Chiang Frank Wang. Frido: Feature pyramid diffusion for complex scene image synthesis. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 579–587, 2023. 2

[18] Rinon Gal, Or Patashnik, Haggai Maron, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. Stylegan-nada: Clip-guided domain adaptation of image generators. *ACM Transactions on Graphics (TOG)*, 41(4):1–13, 2022. 5, 19

[19] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2414–2423, 2016. 1, 2, 5, 6, 7, 9, 14, 17, 21, 22

[20] Rafael C Gonzalez. *Digital image processing*. Pearson education india, 2009. 4, 14, 20

[21] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2014. 1, 2

[22] Daniel J Graham, James M Hughes, Helmut Leder, and Daniel N Rockmore. Statistics, vision, and the analysis of artistic style. *Wiley Interdisciplinary Reviews: Computational Statistics*, 4(2):115–123, 2012. 2

[23] Shuyang Gu, Congliang Chen, Jing Liao, and Lu Yuan. Arbitrary style transfer with deep feature reshuffle. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8222–8231, 2018. 2

[24] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In

*Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10696–10706, 2022. 2

[25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 17

[26] Mingming He, Dongdong Chen, Jing Liao, Pedro V Sander, and Lu Yuan. Deep exemplar-based colorization. *ACM Transactions on Graphics (TOG)*, 37(4):1–16, 2018. 20

[27] Irina Higgins, David Amos, David Pfau, Sebastien Racaniere, Loic Matthey, Danilo Rezende, and Alexander Lerchner. Towards a definition of disentangled representations. *arXiv preprint arXiv:1812.02230*, 2018. 1, 2

[28] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems (NeurIPS)*, 33:6840–6851, 2020. 2, 3, 4, 8, 19

[29] Kibeom Hong, Seogkyu Jeon, Huan Yang, Jianlong Fu, and Hyeran Byun. Domain-aware universal style transfer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 2

[30] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1501–1510, 2017. 1, 2, 4, 7, 8, 9, 17

[31] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 172–189, 2018. 2, 9

[32] Jing Huo, Shiyin Jin, Wenbin Li, Jing Wu, Yu-Kun Lai, Yinghuan Shi, and Yang Gao. Manifold alignment for semantically aligned style transfer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14861–14869, 2021. 2

[33] Yongcheng Jing, Xiao Liu, Yukang Ding, Xinchao Wang, Errui Ding, Mingli Song, and Shilei Wen. Dynamic instance normalization for arbitrary style transfer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 4369–4376, 2020. 2

[34] Yongcheng Jing, Yezhou Yang, Zunlei Feng, Jingwen Ye, Yizhou Yu, and Mingli Song. Neural style transfer: A review. *IEEE Transactions on Visualization and Computer Graphics (TVCG)*, 26(11):3365–3385, 2019. 2

[35] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 694–711. Springer, 2016. 1, 2, 6, 7, 14, 21, 22

[36] Henry Kang, Seungyong Lee, and Charles K Chui. Coherent line drawing. In *Proceedings of the 5th international symposium on Non-photorealistic animation and rendering*, pages 43–50, 2007. 4

[37] Theofanis Karaletsos, Serge Belongie, and Gunnar Rätsch. Bayesian representation learning with oracle constraints. *arXiv preprint arXiv:1506.05011*, 2015. 1, 2

[38] Sergey Karayev, Matthew Trentacoste, Helen Han, Aseem Agarwala, Trevor Darrell, Aaron Hertzmann, and Holger Winnemoeller. Recognizing image style. *arXiv preprint arXiv:1311.3715*, 2013. 2

[39] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4401–4410, 2019. 2

[40] Hadi Kazemi, Seyed Mehdi Iranmanesh, and Nasser Nasrabadi. Style and content disentanglement in generative adversarial networks. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 848–856. IEEE, 2019. 2

[41] Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. Diffusionclip: Text-guided diffusion models for robust image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2426–2435, 2022. 2, 3, 4, 6, 14, 19

[42] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *International Conference on Machine Learning (ICML)*, pages 2649–2658. PMLR, 2018. 2

[43] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 6

[44] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *Advances in Neural Information Processing Systems (NeurIPS)*, 31, 2018. 2, 7

[45] Dmytro Kotovenko, Artsiom Sanakoyeu, Sabine Lang, and Bjorn Ommer. Content and style disentanglement for artistic style transfer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4422–4431, 2019. 1, 2, 5

[46] Dmytro Kotovenko, Matthias Wright, Arthur Heimbrecht, and Bjorn Ommer. Rethinking style transfer: From pixels to parameterized brushstrokes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12196–12205, 2021. 2

[47] Tejas D Kulkarni, William F Whitney, Pushmeet Kohli, and Josh Tenenbaum. Deep convolutional inverse graphics network. *Advances in Neural Information Processing Systems (NeurIPS)*, 28, 2015. 1

[48] Gihyun Kwon and Jong Chul Ye. Diagonal attention and style-based gan for content-style disentanglement in image generation and translation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13980–13989, 2021. 2

[49] Gihyun Kwon and Jong Chul Ye. Clipstyler: Image style transfer with a single text condition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18062–18071, 2022. 5, 8, 9

[50] Berel Lang. *The concept of style*. Cornell University Press, 1987. 4

[51] Hsin-Ying Lee, Hung-Yu Tseng, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. Diverse image-to-image translation via disentangled representations. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 35–51, 2018. 2

[52] Xueting Li, Sifei Liu, Jan Kautz, and Ming-Hsuan Yang. Learning linear transformations for fast image and video

style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

[53] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Universal style transfer via feature transforms. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 386–396, 2017. 1, 2, 4, 8, 17

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

[55] Tianwei Lin, Zhuoqi Ma, Fu Li, Dongliang He, Xin Li, Errui Ding, Nannan Wang, Jie Li, and Xinbo Gao. Drafting and revision: Laplacian pyramid network for fast high-quality artistic style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5141–5150, 2021. 2, 6, 7, 14, 21, 22

[56] Songhua Liu, Tianwei Lin, Dongliang He, Fu Li, Meiling Wang, Xin Li, Zhengxing Sun, Qian Li, and Errui Ding. Adaattn: Revisit attention mechanism in arbitrary neural style transfer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6649–6658, 2021. 2, 6, 7, 14, 21, 22

[57] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In *International Conference on Machine Learning (ICML)*, pages 4114–4124. PMLR, 2019. 1, 2, 7

[58] Cewu Lu, Li Xu, and Jiaya Jia. Contrast preserving decolorization with perception-based quality metrics. *International Journal of Computer Vision (IJCV)*, 110(2):222–239, 2014. 20

[59] Haofei Lu and Zhizhong Wang. Universal video style transfer via crystallization, separation, and blending. In *Proceedings of the International Joint Conferences on Artificial Intelligence Organization (IJCAI), Vienna, Austria*, pages 23–29, 2022. 2

[60] Fujun Luan, Sylvain Paris, Eli Shechtman, and Kavita Bala. Deep photo style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4990–4998, 2017. 19

[61] Calvin Luo. Understanding diffusion models: A unified perspective. *arXiv preprint arXiv:2208.11970*, 2022. 17

[62] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *International Conference on Learning Representations (ICLR)*, 2021. 2

[63] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning (ICML)*, pages 8162–8171. PMLR, 2021. 2, 3

[64] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob Mcgrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning (ICML)*, pages 16784–16804. PMLR, 2022. 2

[65] Dae Young Park and Kwang Hee Lee. Arbitrary style transfer with style-attentional networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5880–5888, 2019. 2

[66] Taesung Park, Jun-Yan Zhu, Oliver Wang, Jingwan Lu, Eli Shechtman, Alexei Efros, and Richard Zhang. Swapping autoencoder for deep image manipulation. *Advances in Neural Information Processing Systems (NeurIPS)*, 33:7198–7211, 2020. 2, 9

[67] DM Parker and Jan B Deregowski. *Perception and artistic style*. Elsevier, 1991. 2

[68] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, pages 8748–8763. PMLR, 2021. 2, 5, 7, 17, 19

[69] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 2, 5

[70] Xuanchi Ren, Tao Yang, Yuwang Wang, and Wenjun Zeng. Rethinking content and style: Exploring bias for unsupervised disentanglement. In *2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, pages 1823–1832. IEEE, 2021. 1, 2

[71] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, 2022. 2

[72] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-assisted Intervention*, pages 234–241. Springer, 2015. 3

[73] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 6, 14, 17

[74] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Raphael Gontijo-Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2

[75] Artsiom Sanakoyeu, Dmytro Kotovenko, Sabine Lang, and Bjorn Ommer. A style-aware content loss for real-time hd style transfer. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 698–714, 2018. 1

[76] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 1, 17

[77] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using

nonequilibrium thermodynamics. In *International Conference on Machine Learning (ICML)*, pages 2256–2265. PMLR, 2015. 2

[78] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations (ICLR)*, 2020. 2, 3, 4, 19

[79] Jan Svoboda, Asha Anoosheh, Christian Osendorfer, and Jonathan Masci. Two-stage peer-regularized feature recombination for arbitrary image style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13816–13825, 2020. 1, 6, 7, 14, 21, 22

[80] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6924–6932, 2017. 2

[81] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017. 2

[82] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017. 7

[83] Pei Wang, Yijun Li, and Nuno Vasconcelos. Rethinking and improving the robustness of image style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 124–133, 2021. 2

[84] Zhizhong Wang, Zhanjie Zhang, Lei Zhao, Zhiwen Zuo, Ailin Li, Wei Xing, and Dongming Lu. Aesust: Towards aesthetic-enhanced universal style transfer. In *Proceedings of the 30th ACM International Conference on Multimedia (ACM MM)*, pages 1095–1106, 2022. 2

[85] Zhizhong Wang, Lei Zhao, Haibo Chen, Ailin Li, Zhiwen Zuo, Wei Xing, and Dongming Lu. Texture reformer: Towards fast and universal interactive texture transfer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 2624–2632, 2022. 2

[86] Zhizhong Wang, Lei Zhao, Haibo Chen, Lihong Qiu, Qihang Mo, Sihuan Lin, Wei Xing, and Dongming Lu. Diversified arbitrary style transfer via deep feature perturbation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7789–7798, 2020. 2, 4, 8, 19

[87] Zhizhong Wang, Lei Zhao, Haibo Chen, Zhiwen Zuo, Ailin Li, Wei Xing, and Dongming Lu. Evaluate and improve the quality of neural style transfer. *Computer Vision and Image Understanding (CVIU)*, 207:103203, 2021. 2

[88] Zhizhong Wang, Lei Zhao, Haibo Chen, Zhiwen Zuo, Ailin Li, Wei Xing, and Dongming Lu. Divswapper: Towards diversified patch-based arbitrary style transfer. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence (IJCAI)*, pages 4980–4987, 2022. 2, 19

[89] Zhizhong Wang, Lei Zhao, Sihuan Lin, Qihang Mo, Huiming Zhang, Wei Xing, and Dongming Lu. Glstylenet: exquisite style transfer combining global and local pyramid features. *IET Computer Vision*, 14(8):575–586, 2020. 2

[90] Zhizhong Wang, Lei Zhao, Zhiwen Zuo, Ailin Li, Haibo Chen, Wei Xing, and Dongming Lu. Microast: Towards super-fast ultra-resolution arbitrary style transfer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 2742–2750, 2023. 2

[91] Wayne Wu, Kaidi Cao, Cheng Li, Chen Qian, and Chen Change Loy. Disentangling content and style via unsupervised geometry distillation. *arXiv preprint arXiv:1905.04538*, 2019. 2

[92] Xiaolei Wu, Zhihao Hu, Lu Sheng, and Dong Xu. Styleformer: Real-time arbitrary style transfer via parametric style composition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14618–14627, 2021. 2

[93] Zijie Wu, Zhen Zhu, Junping Du, and Xiang Bai. Ccpl: Contrastive coherence preserving loss for versatile style transfer. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2022. 2

[94] Saining Xie and Zhuowen Tu. Holistically-nested edge detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 1395–1403, 2015. 4

[95] Yabin Zhang, Minghan Li, Ruihuang Li, Kui Jia, and Lei Zhang. Exact feature distribution matching for arbitrary style transfer and domain generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8035–8045, 2022. 1, 2, 6, 7, 14, 21, 22

[96] Yuxin Zhang, Fan Tang, Weiming Dong, Haibin Huang, Chongyang Ma, Tong-Yee Lee, and Changsheng Xu. Domain enhanced arbitrary image style transfer via contrastive learning. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–8, 2022. 2

[97] Yu Zhang, Peter Tiňo, Aleš Leonardis, and Ke Tang. A survey on neural network interpretability. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2021. 1

[98] Yexun Zhang, Ya Zhang, and Wenbin Cai. Separating style and content for generalized style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8447–8455, 2018. 1, 2

[99] Zhiwen Zuo, Lei Zhao, Shuobin Lian, Haibo Chen, Zhizhong Wang, Ailin Li, Wei Xing, and Dongming Lu. Style fader generative adversarial networks for style degree controllable artistic style transfer. In *Proc. Int. Joint Conf. on Artif. Intell.(IJCAI)*, pages 5002–5009, 2022. 2

# Supplementary Material

## A. Societal Impact

**Positive Impact.** There may be three positive impacts of the proposed method. (1) The proposed method may help the workers engaged in artistic creation or creative design improve the efficiency and quality of their work. (2) The proposed method may inspire researchers in similar fields to design more effective and superior approaches in the future. (3) The proposed method may help the common users obtain more satisfactory creative results.

**Negative Impact.** The proposed method may be used for generating counterfeit artworks. To mitigate this, further research on identification of generated content is needed.

## B. Used Assets

We used the following assets to (1) conduct the comparison experiments [1.-10.] and (2) train the proposed style transfer networks [11.-12.]. To the best of our knowledge, these assets have no ethical concerns.

1. Gatys [19]: `https://github.com/leongatys/PytorchNeuralStyleTransfer`, MIT License.
2. EFDM [95]: `https://github.com/YBZh/EFDM`, MIT License.
3. StyTr$^2$ [12]: `https://github.com/diyiiyiii/StyTR-2`, No License.
4. ArtFlow [1]: `https://github.com/pkuanjie/ArtFlow`, No License.
5. AdaAttN [56]: `https://github.com/Huage001/AdaAttN`, No License.
6. IECAST [6]: `https://github.com/HalbertCH/IEContraAST`, MIT License.

7. MAST [13]: `https://github.com/diyiiyiii/Arbitrary-Style-Transfer-via-Multi-Adaptation-Network`, No License.
8. TPFR [79]: `https://github.com/nnaisense/conditional-style-transfer`, View License in repository.
9. Johnson [35]: `https://github.com/abhiskk/fast-neural-style`, MIT License.
10. LapStyle [55]: `https://github.com/PaddlePaddle/PaddleGAN/blob/develop/docs/en_US/tutorials/lap_style.md`, Apache-2.0 License.
11. ADM [14]: `https://github.com/openai/guided-diffusion`, MIT License.
12. ImageNet [73]: `https://image-net.org/`, Unknown License.

## C. Details of Style Removal Process

As detailed in Algorithm 1, the style removal process consists of two steps. In the first step, we remove the color of the input image $I$ using a color removal operation $\mathcal{R}_{color}$ (*e.g.*, the commonly used ITU-R 601-2 luma transform [20]), obtaining grayscale image $I'$. In the second step, we use the pre-trained diffusion model $\epsilon_\theta$ and adopt the deterministic DDIM forward and reverse processes to gradually remove the style information. To accelerate the process without sacrificing much performance, we use fewer discretization steps $\{t_s\}_{s=1}^{S_{for}}$ such that $t_1 = 0$ and $t_{S_{for}} = T_{remov}$. We set $S_{for} = 40$ for forward process and $S_{rev} = 40$ for reverse process in all experiments. While using larger $S_{for}$ or $S_{rev}$ could reconstruct the high-frequency details better, we found the current setting is enough for our task. *For more details about their effects, we suggest the readers refer to [41].* After $K_r$ iterations (we set $K_r = 5$ for all experiments) of forward and reverse processes, the style characteristics of $I'$ will be dispelled, and thus we obtain the content $I^c$ of the input image.

## D. Details of StyleDiffusion Fine-tuning

Similar to [41] and detailed in Algorithm 2, we first precompute the content latents $\{x^{ci}\}_{i=1}^N$ using the deterministic DDIM forward process of the pre-trained diffusion model $\epsilon_\theta$. *The precomputed content latents can be stored and reused for fine-tuning other styles.* In our experiments, we fine-tune the diffusion models for all styles using the same precomputed latents of 50 content images sampled from ImageNet [73]. Fine-tuning with more content images may improve the results but also increases the time cost. Thus, we made a trade-off and found the current setting could work well for most cases. To accelerate the pro-

**Algorithm 1:** Style Removal Process.

**Input:** pre-trained model $\epsilon_\theta$, input image $I$, return step $T_{remov}$, forward step $S_{for}$, reverse step $S_{rev}$, iteration $K_r$

**Output:** input image's content $I^c$

// Remove color

1   $I' = \mathcal{R}_{color}(I)$

// Diffusion-based style removal

2   Compute $\{t_s\}_{s=1}^{S_{for}}$ s.t. $t_1 = 0, t_{S_{for}} = T_{remov}$

3   $x_0 \leftarrow I$

4   **for** $k = 1 : K_r$ **do**

5     **for** $s = 1 : S_{for} - 1$ **do**

6       $x_{t_{s+1}} \leftarrow$
        $\sqrt{\bar{\alpha}_{t_{s+1}}} f_\theta(x_{t_s}, t_s) + \sqrt{1 - \bar{\alpha}_{t_{s+1}}} \epsilon_\theta(x_{t_s}, t_s)$

7     **end**

8     $x_{t_{S_{rev}}} \leftarrow x_{t_{S_{for}}}$

9     **for** $s = S_{rev} : 2$ **do**

10      $x_{t_{s-1}} \leftarrow$
        $\sqrt{\bar{\alpha}_{t_{s-1}}} f_{\hat{\theta}}(x_{t_s}, t_s) + \sqrt{1 - \bar{\alpha}_{t_{s-1}}} \epsilon_{\hat{\theta}}(x_{t_s}, t_s)$

11    **end**

12 **end**

13 $I^c \leftarrow x_0$

---

cess, we use fewer discretization steps $\{t_s\}_{s=1}^{S_{for}}$ such that $t_1 = 0$ and $t_{S_{for}} = T_{trans}$. We set $S_{for} = 40$ for forward process and $S_{rev} = 6$ for reverse process in all experiments. We found $S_{rev} = 6$ is enough to reconstruct clear content structures during style transfer.

In the second step, we precompute the style latent $x^s$ with the same process as above. The style latent will be used to optimize the style reconstruction loss.

In the third step, we copy $\epsilon_\theta$ to $\epsilon_{\hat{\theta}}$ and start to update $\epsilon_{\hat{\theta}}$ in two substeps. In the first substep, we feed the style latent $x^s$ and generate the stylized image $I_{ss}$ through the deterministic DDIM reverse process. The model is updated under the guidance of the style reconstruction loss $\mathcal{L}_{SR}$. The first substep is repeated $K_s$ times (we set $K_s = 50$ for all experiments) until converged. In the second substep, we feed each content latent in $\{x^{ci}\}_{i=1}^N$ and generate the stylized image through the deterministic DDIM reverse process. The model is updated under the guidance of the style disentanglement loss $\mathcal{L}_{SD}$. At last, we repeat the whole third step $K$ epochs (we set $K = 5$ for all experiments) until converged.

## E. Timing and Resource Information

Here, we provide more details on the timing and resource information of our StyleDiffusion using an Nvidia Tesla A100 GPU when stylizing $512 \times 512$ size images.

**Style Removal.** When we use the default setting $(S_{for}, S_{rev}) = (40, 40)$, the forward and reverse processes

---

**Algorithm 2:** StyleDiffusion Fine-tuning.

**Input:** pre-trained model $\epsilon_\theta$, content images' contents $\{I_{ci}^c\}_{i=1}^N$, style image's content $I_s^c$, style image $I_s$, return step $T_{trans}$, forward step $S_{for}$, reverse step $S_{rev}$, fine-tuning epoch $K$, style reconstruction iteration $K_s$

**Output:** fine-tuned model $\epsilon_{\hat{\theta}}$

// Precompute content latents

1   Compute $\{t_s\}_{s=1}^{S_{for}}$ s.t. $t_1 = 0, t_{S_{for}} = T_{trans}$

2   **for** $i = 1 : N$ **do**

3     $x_0 \leftarrow I_{ci}^c$

4     **for** $s = 1 : S_{for} - 1$ **do**

5       $x_{t_{s+1}} \leftarrow$
        $\sqrt{\bar{\alpha}_{t_{s+1}}} f_\theta(x_{t_s}, t_s) + \sqrt{1 - \bar{\alpha}_{t_{s+1}}} \epsilon_\theta(x_{t_s}, t_s)$

6     **end**

7     Save the latent $x^{ci} \leftarrow x_{t_{S_{for}}}$

8   **end**

// Precompute style latent

9   Compute $\{t_s\}_{s=1}^{S_{for}}$ s.t. $t_1 = 0, t_{S_{for}} = T_{trans}$

10 $x_0 \leftarrow I_s^c$

11 **for** $s = 1 : S_{for} - 1$ **do**

12   $x_{t_{s+1}} \leftarrow \sqrt{\bar{\alpha}_{t_{s+1}}} f_\theta(x_{t_s}, t_s) + \sqrt{1 - \bar{\alpha}_{t_{s+1}}} \epsilon_\theta(x_{t_s}, t_s)$

13 **end**

14 Save the latent $x^s \leftarrow x_{t_{S_{for}}}$

// Fine-tune the diffusion model

15 Initialize $\epsilon_{\hat{\theta}} \leftarrow \epsilon_\theta$

16 Compute $\{t_s\}_{s=1}^{S_{rev}}$ s.t. $t_1 = 0, t_{S_{rev}} = T_{trans}$

17 **for** $k = 1 : K$ **do**

// Optimize the style reconstruction loss

18    **for** $i = 1 : K_s$ **do**

19      $x_{t_{S_{rev}}} \leftarrow x^s$

20      **for** $s = S_{rev} : 2$ **do**

21        $x_{t_{s-1}} \leftarrow$
         $\sqrt{\bar{\alpha}_{t_{s-1}}} f_{\hat{\theta}}(x_{t_s}, t_s) + \sqrt{1 - \bar{\alpha}_{t_{s-1}}} \epsilon_{\hat{\theta}}(x_{t_s}, t_s)$

22        $I_{ss} \leftarrow f_{\hat{\theta}}(x_{t_s}, t_s)$

23        $\mathcal{L} \leftarrow \mathcal{L}_{SR}(I_{ss}, I_s)$

24        Take a gradient step on $\nabla_{\hat{\theta}} \mathcal{L}$

25      **end**

26    **end**

// Optimize the style disentanglement loss

27    **for** $i = 1 : N$ **do**

28      $x_{t_{S_{rev}}} \leftarrow x^{ci}$

29      **for** $s = S_{rev} : 2$ **do**

30        $x_{t_{s-1}} \leftarrow$
         $\sqrt{\bar{\alpha}_{t_{s-1}}} f_{\hat{\theta}}(x_{t_s}, t_s) + \sqrt{1 - \bar{\alpha}_{t_{s-1}}} \epsilon_{\hat{\theta}}(x_{t_s}, t_s)$

31        $I_{cs} \leftarrow f_{\hat{\theta}}(x_{t_s}, t_s)$

32        $\mathcal{L} \leftarrow \mathcal{L}_{SD}(I_{ci}^c, I_{cs}, I_s^c, I_s)$

33        Take a gradient step on $\nabla_{\hat{\theta}} \mathcal{L}$

34      **end**

35    **end**

36 **end**

each takes around 4.921 seconds. Therefore, the whole style removal process takes around $2 \times 4.921 \times 5 = 49.21$ seconds. It requires about 11GB of GPU memory to run at resolution $512 \times 512$ pixels.

**Fine-tuning.** As illustrated in Algorithm 2, the StyleDiffusion fine-tuning process consists of a latent precomputing stage and a model updating stage. The latent precomputing stage is carried out just once and can be reused for fine-tuning other styles. When we use $S_{for} = 40$ as default, the forward process takes around 4.921 seconds. Therefore, when we precompute the latents from 50 images, it takes around $50 \times 4.921 = 246.05$ seconds and requires about 11GB GPU memory. For the model updating stage, when the batch size is 1 and $S_{rev} = 6$, the first substep (optimizing the style reconstruction loss $\mathcal{L}_{SR}$) takes around 2.092 seconds for each repeat, and the second substep (optimizing the style disentanglement loss $\mathcal{L}_{SD}$) takes around 3.351 seconds for each content latent. Therefore, one epoch with 50 repeated first substep and 50 precomputed content latents for the second substep takes around $50 \times 2.092 + 50 \times 3.351 = 272.15$ seconds. When we fine-tune the model with 5 epochs, it takes around 23 minutes in total. The fine-tuning process requires about 26GB of GPU memory.

**Inference.** When we use the default setting $(S_{for}, S_{rev}) = (40, 6)$, the forward process takes around 4.921 seconds and the reverse process takes around 0.691 seconds. Therefore, the total inference time is $4.921 + 0.691 = 5.612$ seconds. The inference process requires about 13GB of GPU memory.

Currently, *we have not optimized the model size and GPU memory consumption here*. We believe there is substantial room for improvement, and we would like to elabo-

rate on that in future work.

## F. More Ablation Study and Analyses

**Quantitative Analyses of C-S Disentanglement.** Here, we provide more quantitative results to analyze the C-S Disentanglement achieved by our StyleDiffusion. As shown in Fig. 9, we observe that the style is well disentangled from the content in the style image by adjusting the return step $T_{remov}$ of the style removal module. As such, when more style information is removed (blue line), it will be transferred to the corresponding stylized result (orange line). The quantitative analyses are consistent with the qualitative
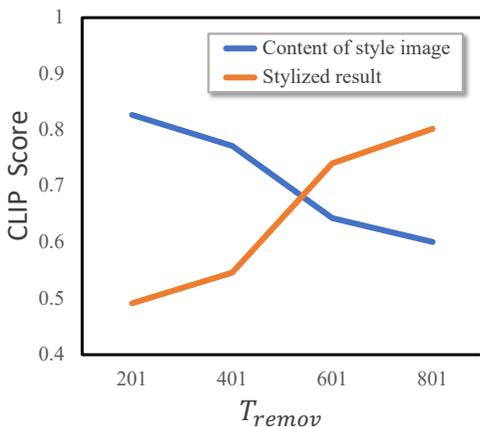


Figure 10. **C-S trade-off** achieved by adjusting the return step $T_{trans}$ of the *style transfer module* at the **training** stage while fixing $T_{trans} = 301$ at the testing stage. SSIM and CLIP score (averaged on 384 image pairs) measure the content similarity and the style similarity, respectively.



Figure 9. **C-S disentanglement of style image** achieved by adjusting the return step $T_{remov}$ of the *style removal module*. CLIP score (averaged on 384 image pairs) measures the style similarity with the style image. When more style information is removed (blue line), it will be transferred to the stylized result (orange line).
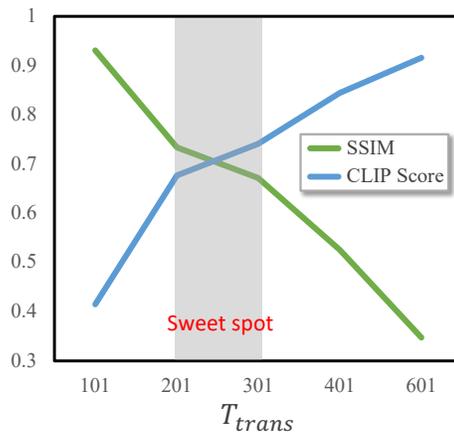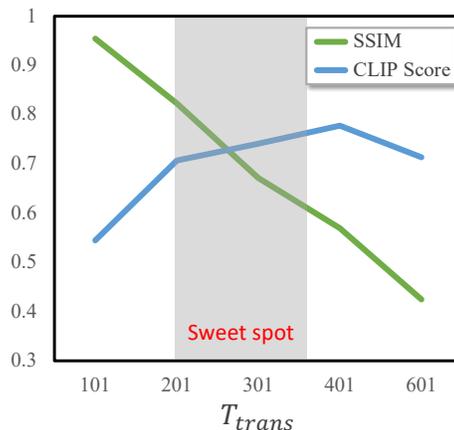


Figure 11. **C-S trade-off** achieved by adjusting the return step $T_{trans}$ of the *style transfer module* at the **testing** stage while fixing $T_{trans} = 301$ at the training stage. SSIM and CLIP score (averaged on 384 image pairs) measure the content similarity and the style similarity, respectively.

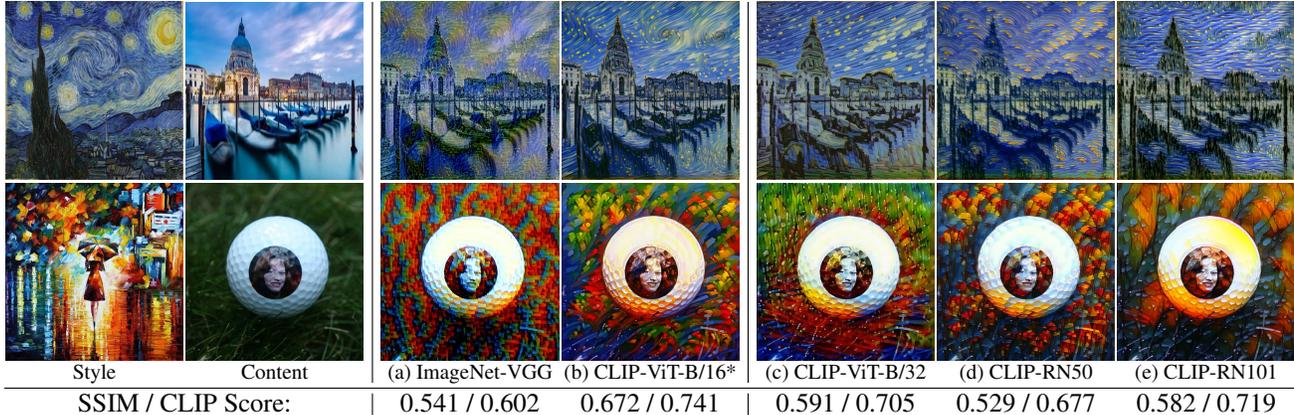| | | (a) ImageNet-VGG | (b) CLIP-ViT-B/16* | (c) CLIP-ViT-B/32 | (d) CLIP-RN50 | (e) CLIP-RN101 |
|---|---|---|---|---|---|---|
| Style | Content | | | | | |
| SSIM / CLIP Score: | | 0.541 / 0.602 | 0.672 / 0.741 | 0.591 / 0.705 | 0.529 / 0.677 | 0.582 / 0.719 |

Figure 12. **Ablation study** on different **disentanglement space** (VGG vs. CLIP, columns (a-b)) and **CLIP image encoders** (columns (b-e)). * denotes our default setting. Zoom-in for better comparison.

results displayed in Fig. 4 of our main paper.

**Quantitative Analyses of C-S Trade-off.** We also provide more quantitative results to analyze the C-S trade-off achieved by our StyleDiffusion. As shown in Fig. 10 and Fig. 11, we can flexibly control the C-S trade-off at both the training stage (Fig. 10) and the testing stage (Fig. 11) by adjusting the return step $T_{trans}$ of diffusion models. The sweet spot areas are highlighted in the figures, which are the most probable for obtaining satisfactory results. Overall, the quantitative analyses are consistent with the qualitative results displayed in Fig. 5 of our main paper.

**CLIP Space vs. VGG Space.** As discussed in our main paper, we leverage the open-domain CLIP [68] space to formulate the style disentanglement. The pre-trained CLIP space integrates rich cross-domain image (and supplementarily, text) knowledge and thus can measure the "style distance" more accurately. As shown in Fig. 12 (a-b), we compare it with the ImageNet [73] pre-trained VGG-19 [76], which has been widely adopted in prior arts [19, 30] to extract the style information. As is evident, using the CLIP space recovers the style information more sufficiently and realistically, significantly outperforming the VGG space (which is also validated by the bottom quantitative scores). It may be attributed to the fact that the VGG is pre-trained on ImageNet and therefore lacks a sufficient understanding of artistic styles. In contrast, the CLIP space encapsulates a myriad of knowledge of not only the photograph domain but also the artistic domain, which is more powerful in depicting the style of an image. Besides, it is worth noting that the CLIP space naturally provides multi-modal compatibility, which can facilitate users to control the style transfer with multi-modal signals, *e.g.*, image and text (see later Sec. G).

**Different CLIP Image Encoders.** We also investigate the effects of different CLIP [68] image encoders to conduct the style disentanglement. As shown in Fig. 12 (b-e), in general, ViTs [15] achieve better visual results than

ResNets (RN) [25], *e.g.*, the brushstrokes are more natural in the top row, and the colors are more vivid in the bottom row. And ViT-B/16 performs better than ViT-B/32 in capturing more fine-grained styles. Interestingly, our findings coincide with the reported performance of these image encoders on high-level vision tasks (*e.g.*, classification) in the original CLIP paper [68]. It indicates that our stylization performance is closely related to the high-level semantic representations learned by the image encoder, which also gives evidence to the correlations between high-level vision tasks and low-level vision tasks.

**Diffusion-based Style Removal vs. AE-based Style Removal.** To demonstrate the superiority of diffusion-based style removal, we compare it with a possible alternative, *i.e.*, Auto-Encoders (AEs), since one may argue that the diffusion model is a special kind of (Variational) Auto-Encoder network [61]. We directly use the AEs released by Li *et al.* [53], which employ the VGG-19 network [76] as the encoders, fix them and train decoder networks for inverting VGG features to the original images. They select feature maps at five layers of the VGG-19, i.e., Relu_X_1 (X=1,2,3,4,5), and train five decoders accordingly, which we denote as AEX (X=1,2,3,4,5) in the following. When used for style removal, we iteratively perform the encoding and decoding processes of AEs for the input images. The comparison results are shown in Fig. 13. As can be observed in the bottom five rows, AE-based style removal cannot plausibly remove the detailed style and often introduces color noises/artifacts and destroys the content structures, which is undesirable for style removal. By contrast, diffusion-based style removal can smoothly remove the style details while preserving the main content structures, significantly outperforming AE-based style removal.

## G. Extensions

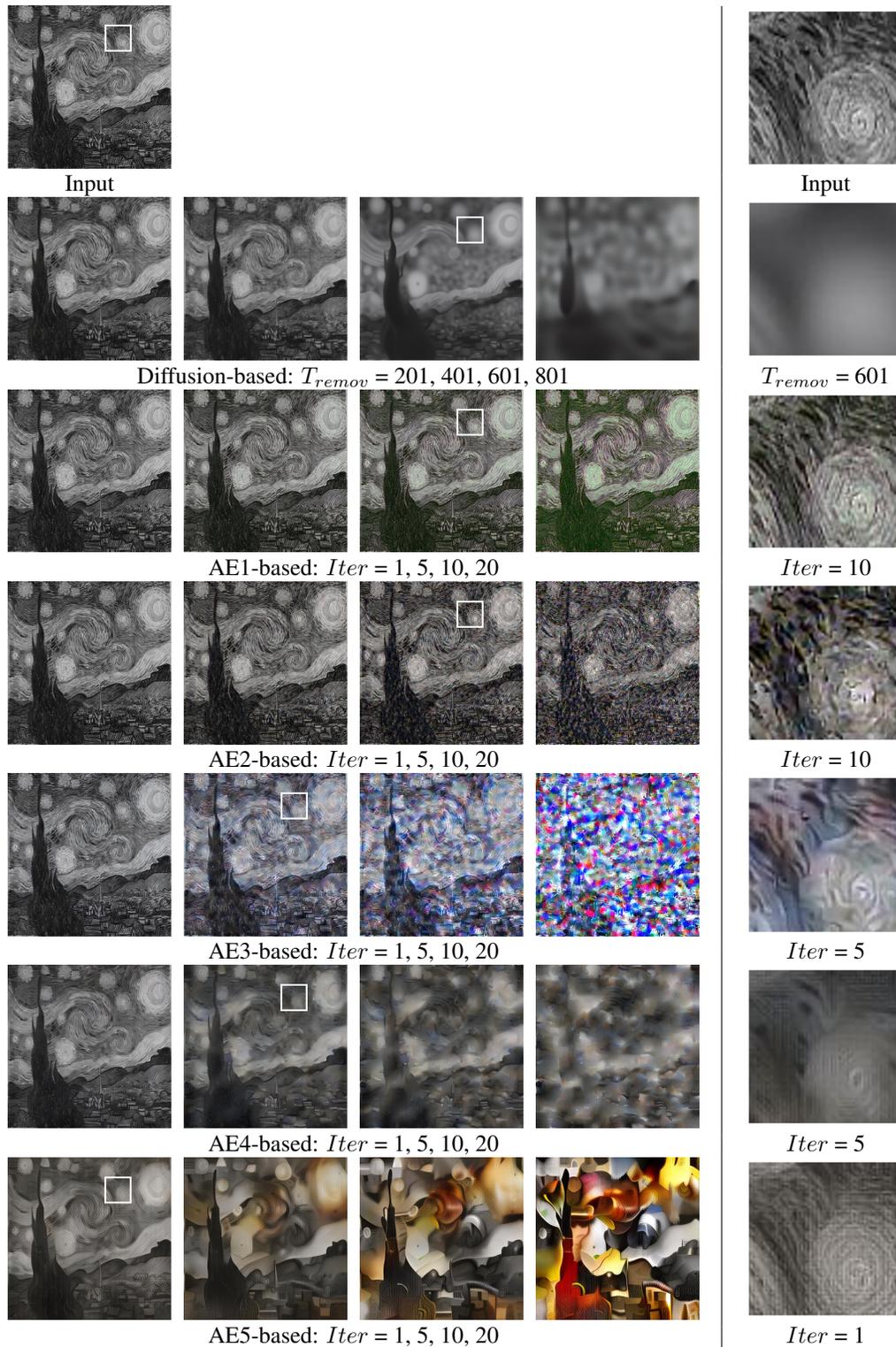**Photo-realistic Style Transfer.** Our StyleDiffusion suc-

Figure 13. **Diffusion-based style removal vs. AE-based style removal.** The last column shows the enlarged areas of the corresponding best style removed results manually selected in each row. As can be observed, diffusion-based style removal can better remove the detailed style of the style image while preserving the main content structures. In contrast, AE-based style removal cannot plausibly remove the detailed style and often introduces color noises/artifacts and destroys the content structures.
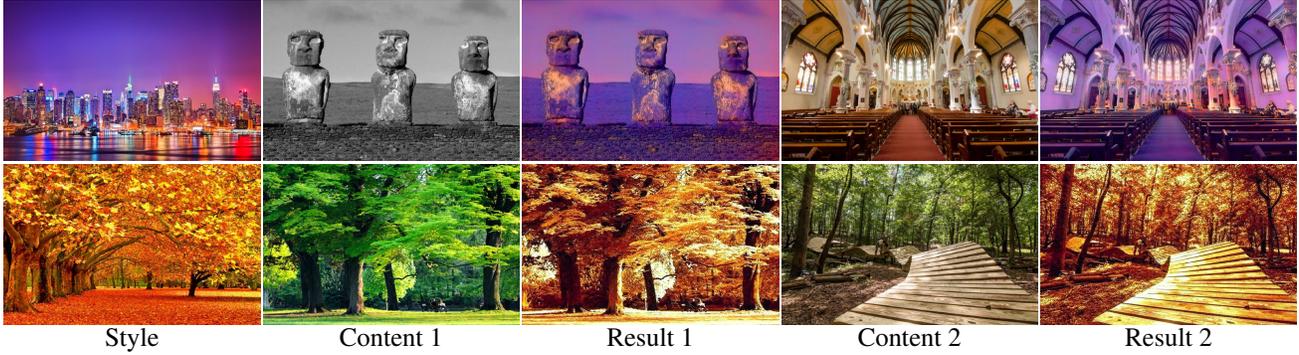
Figure 14. **Photo-realistic style transfer** achieved by our StyleDiffusion. We set $T_{remov} = 401$ and $T_{trans} = 101$ for this task.



| Content | Style | Result | + "Pointillism" | + "Sketch" | + "Cubism" | + "Watercolor" |

Figure 15. **Multi-modal style manipulation.** Our framework is compatible with image and text modulation signals, which provides users with a more flexible way to manipulate the style of images.
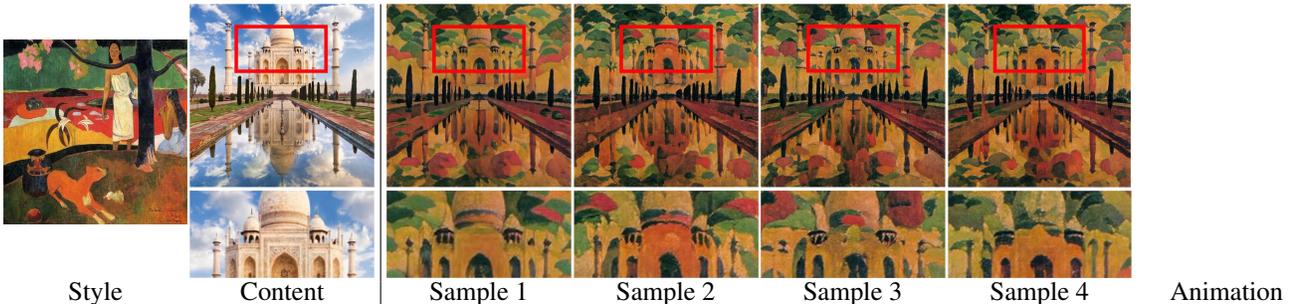


| Style | Content | Sample 1 | Sample 2 | Sample 3 | Sample 4 | Animation |

Figure 16. **Diversified style transfer.** Our framework can easily achieve diversified style transfer during inference by directly adopting the stochastic DDPM [28] forward process. Click on the last image to see animation using Adobe Reader.

cessfully separates style from content in a controllable manner. Thus, it can easily achieve photo-realistic style transfer [60] by adjusting the content extraction of the style removal module. Specifically, since the style of a photo is mainly reflected by the low-level and high-frequency features such as colors and brightness, we reduce $T_{remov}$ to a relatively smaller value, *e.g.*, 401. Moreover, to better preserve the content structures, we adjust the style transfer process and reduce $T_{trans}$ to 101. We show some photo-realistic style transfer results synthesized by our StyleDiffusion in Fig. 14.

**Multi-modal Style Manipulation.** As our framework leverages the open-domain CLIP [68] space to measure the "style distance", it is naturally compatible with image and text modulation signals. By adding a directional CLIP loss term [18, 41] to our total loss, our framework can eas-

ily achieve multi-modal style manipulation, as shown in Fig. 15. *As far as we know, our framework is the first unified framework to achieve both image and text guided style transfer.*

**Diversified Style Transfer.** In the fine-tuning, our style transfer module adopts the deterministic DDIM [78] forward and reverse processes (Eq. (8) and Eq. (9) in the main paper). However, during inference, we can directly replace the deterministic DDIM forward process with the stochastic DDPM [28] forward process (Eq. (2) in the main paper) to achieve diversified style transfer [86], as shown in Fig. 16. The users can easily trade off the diversity and quality by adjusting the return step or iteration of the DDPM forward process. The diverse results can give users endless choices to obtain more satisfactory results [86, 88].
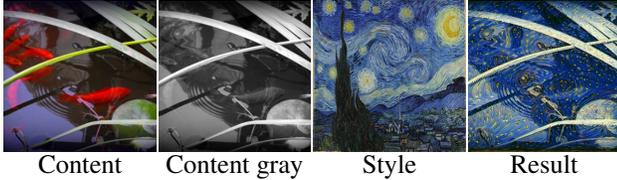
Content    Content gray    Style    Result

Figure 17. **Failure case of type 1: vanishing of salient content.** Some results generated by our method may vanish the salient content of the content image, *e.g.*, the red carps.



Style    Result 1    Result 2    Result 3

Figure 18. **Failure case of type 2: biased color distribution.** Our method may generate results that deviate from the color distribution of the style image.

## H. More Comparison Results

In Fig. 20 and 21, we provide more qualitative comparison results with state-of-the-art style transfer methods.

## I. Additional Stylized Results

In Fig. 22 and 23, we provide additional stylized results synthesized by our proposed StyleDiffusion.

## J. Limitation and Discussion

Except for the limitations we have discussed in the main paper, here we provide some failure cases and analyze the reasons behind them. Further, we also discuss the possible solutions to address them, which may help inspire future improvements to our framework.

**Vanishing of Salient Content.** Some of our generated results may vanish the salient content of the content image, *e.g.*, the red carps in Fig. 17. It can be attributed to the color removal operation used in our style removal module. The commonly used ITU-R 601-2 luma transform [20] may not well preserve the original RGB image's color contrast and color importance, as shown in column 2 of Fig. 17. We adopt it here mainly for its simplicity and fast speed. This problem may be addressed by using more advanced contrast-preserving decolorization techniques, like [58].

**Biased Color Distribution.** As shown in Fig. 18, though our method learns the challenging pointillism style well, the color distribution seems to stray from that of the style image. This problem can be alleviated by increasing the style reconstruction iteration $K_s$ (see Algorithm 2) to inject more style prior, but the training time also increases significantly. One may consider borrowing some ideas from



Style    Style removed    Result 1    Result 2

Figure 19. **Failure case of type 3: inseparable content and style.** Our method is hard to transfer plausible style for style images with inseparable content and style. The second column shows the style removed result of the style image.

existing color transfer approaches [26] to address this problem.

**Inseparable Content and Style.** Our method is hard to achieve plausible style transfer for style images with inseparable content and style, *e.g.*, the simple line art shown in Fig. 19. Since the content of line art is also its style, our framework is hard to separate them properly, as shown in column 2 of Fig. 19. One possible solution is to treat line art as the style only and increase the return step $T_{remov}$ of the style removal module to dispel as much style information as possible, or increase the return step $T_{trans}$ of the style transfer module to learn as sufficient line art style as possible.
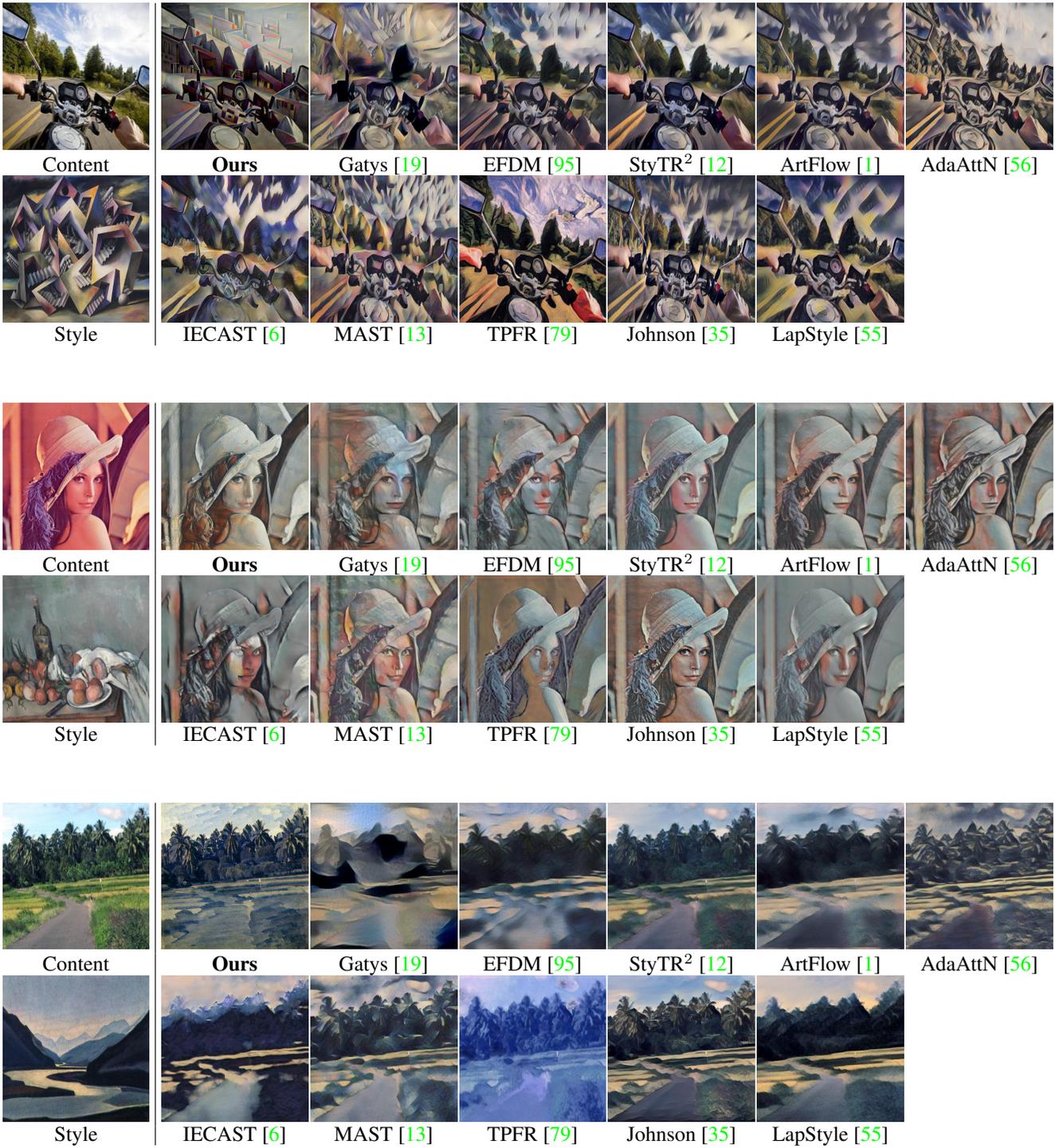
Content  **Ours**  Gatys [19]  EFDM [95]  StyTR$^2$ [12]  ArtFlow [1]  AdaAttN [56]

Style  IECAST [6]  MAST [13]  TPFR [79]  Johnson [35]  LapStyle [55]

Content  **Ours**  Gatys [19]  EFDM [95]  StyTR$^2$ [12]  ArtFlow [1]  AdaAttN [56]

Style  IECAST [6]  MAST [13]  TPFR [79]  Johnson [35]  LapStyle [55]

Content  **Ours**  Gatys [19]  EFDM [95]  StyTR$^2$ [12]  ArtFlow [1]  AdaAttN [56]

Style  IECAST [6]  MAST [13]  TPFR [79]  Johnson [35]  LapStyle [55]

Figure 20. **More qualitative comparison results (set 1)** with state of the art. Zoom-in for better comparison.
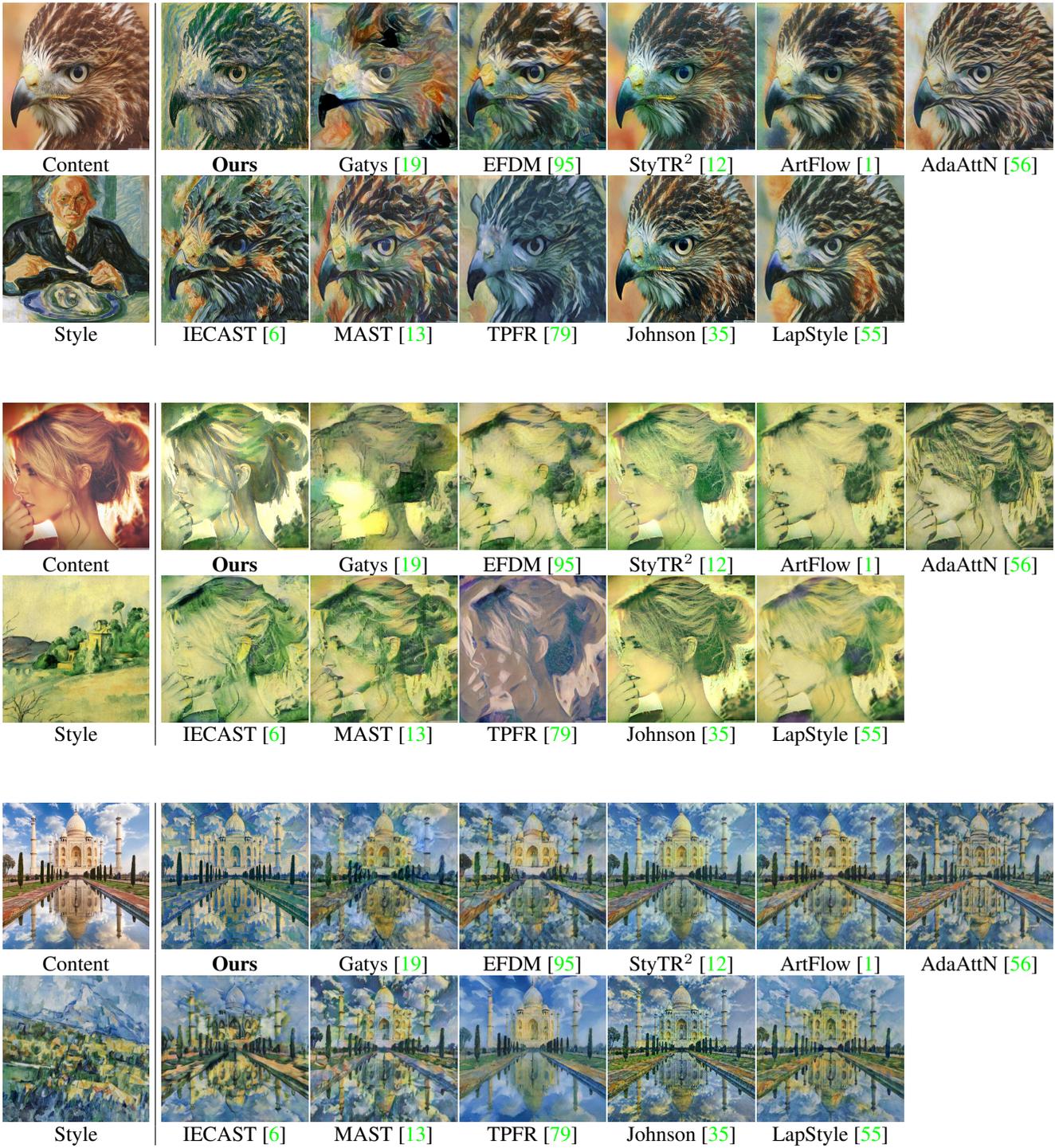
Figure 21. **More qualitative comparison results (set 2)** with state of the art. Zoom-in for better comparison.

Figure 22. **Additional stylized results (set 1)** synthesized by our proposed StyleDiffusion. The first row shows content images and the first column shows style images.
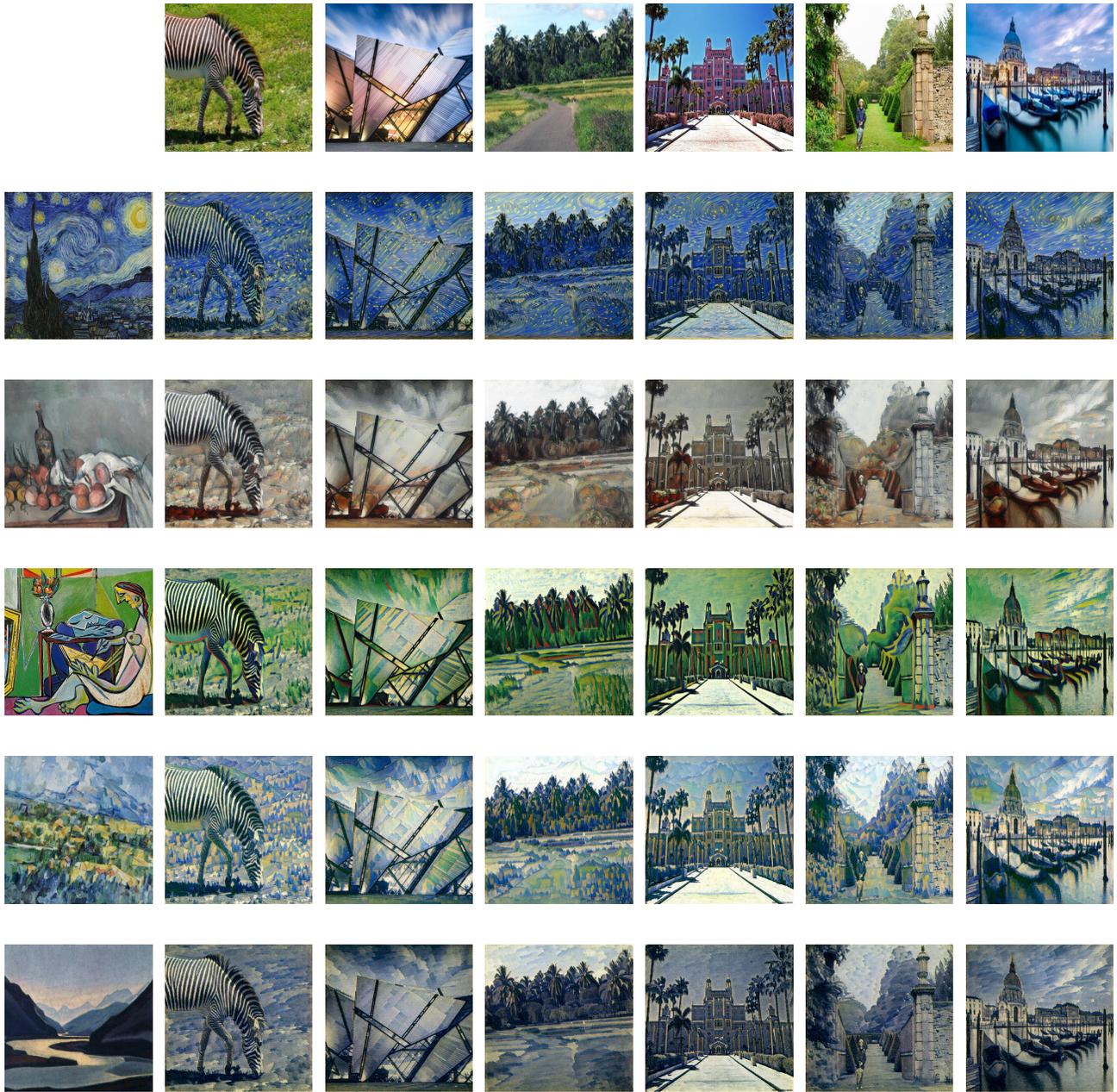
Figure 23. **Additional stylized results (set 2)** synthesized by our proposed StyleDiffusion. The first row shows content images and the first column shows style images.