

# StyleSSP: Sampling StartPoint Enhancement for Training-free Diffusion-based Method for Style Transfer

Ruojun Xu  
Zhejiang University, Dcar  
Hangzhou, China  
ruojunxu@zju.edu.cn

Weijie Xi  
Dcar,  
Beijing, China  
xiweijie@bytedance.com

XiaoDi Wang  
Dcar,  
Beijing, China  
wangxiaodi.00@bytedance.com

Yongbo Mao  
Dcar,  
Beijing, China  
maoyongbo@bytedance.com

Zach Cheng  
Dcar,  
Beijing, China  
chengyi.2024@bytedance.com

## Abstract

*Training-free diffusion-based methods have achieved remarkable success in style transfer, eliminating the need for extensive training or fine-tuning. However, due to the lack of targeted training for style information extraction and constraints on the content image layout, training-free methods often suffer from layout changes of original content and content leakage from style images. Through a series of experiments, we discovered that an effective startpoint in the sampling stage significantly enhances the style transfer process. Based on this discovery, we propose **StyleSSP**, which focuses on obtaining a better startpoint to address layout changes of original content and content leakage from style image. StyleSSP comprises two key components: (1) **Frequency Manipulation**: To improve content preservation, we reduce the low-frequency components of the DDIM latent, allowing the sampling stage to pay more attention to the layout of content images; and (2) **Negative Guidance via Inversion**: To mitigate the content leakage from style image, we employ negative guidance in the inversion stage to ensure that the startpoint of the sampling stage is distanced from the content of style image. Experiments show that StyleSSP surpasses previous training-free style transfer baselines, particularly in preserving original content and minimizing the content leakage from style image.*

## 1. Introduction

Recently, Diffusion Models (DMs) have yielded high-quality results in various areas such as text-to-image generation [27, 37, 40] and image or video editing [3, 7, 8, 14]. As part of image editing, diffusion-based style transfer meth-

ods [4, 12, 31, 49] have garnered widespread attention. These methods enable condition-guided image generation that transfers the style of one image onto another while maintaining the original content.

Previous diffusion-based style transfer methods [20, 21, 31, 55] leverage the generative capability of pre-trained DMs using inference-stage optimization, yet they are either time-consuming or fail to fully utilize the generative ability of large-scale diffusion models. Based on these challenges, training-free methods [4, 20, 45, 46] have been proposed. Although these methods have shown promising results, they still encounter two key issues: (1) **Content preservation problem**. Due to the lack of constraints directly imposed on the content of generated images during training, training-free methods often struggle to maintain the original semantic and structural content [26]. Although additional modules like ControlNet [53] can be used as content constraints, experiments have shown that these methods still risk failure (as shown in supplementary materials Sec. 7.1). This issue largely arises from the diffusion model’s imbalanced preference for different conditions when multiple conditions are injected into the U-Net during the sampling stage [9, 13]; (2) **Content leakage from the style image**. Without targeted training for style extraction, training-free methods struggle to effectively decouple the style and content. Therefore, when the style image is directly injected into the pre-trained diffusion model, the generation process is inevitably influenced by the content of style image. Fig. 1 illustrates examples of these two problems.

To address these challenges, we begin by noting recent advancements in image synthesis tasks with DMs. These studies reveal the significant influence of the initial noise (referred to here as the “startpoint”) on the generated out-

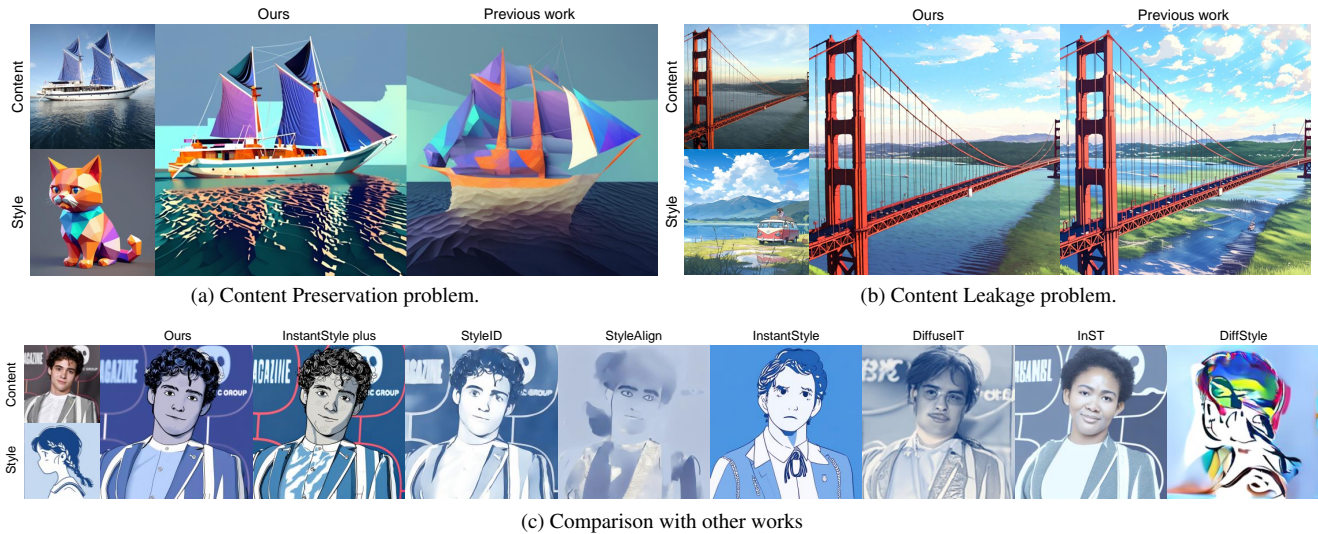


Figure 1. Current problems for style transfer and our improvements. (a) **Original content changes** in previous work (right) even with ControlNet as an additional content controller. (b) **Content leakage** from style image in previous work (right), where the river from original image is covered by a lawn that shouldn’t exist. (c) Given a style image and content image, StyleSSP is capable of synthesizing new images that achieve the best style transfer effect while preserving the details of original content.

come. For example, FreeNoise [32] analyzes the impact of startpoint within video diffusion models, emphasizing the importance of initialization with a sequence of long-range correlated noises. Similarly, FlexiEdit [16] enhances the startpoint by reducing low-frequency components, improving the fidelity to editing prompts. While the significance of startpoint selection is increasingly acknowledged in generation and editing tasks, it remains largely unexplored in style transfer. StyleID [4] does incorporate startpoint manipulation, but only by rescaling the startpoint to offset the pre-trained model’s tendency to generate images with median colors, without fully investigating its role in style transfer.

Inspired by the aforementioned methods, we creatively argue that refining the sampling startpoint is an effective strategy for improving style transfer. Our supplementary materials Sec. 7.1 further demonstrates this point. In these experiments, we show the startpoint’s significant impact on content preservation and tonal adjustment. Based on these findings, we propose **StyleSSP (Style transfer method via Sampling StartPoint enhancement)**, a training-free approach that refines the sampling startpoint in diffusion models. To the best of our knowledge, this work is the first to highlight the importance of selecting an effective sampling startpoint to improve style transfer in a training-free, diffusion-based framework.

First, we propose frequency manipulation to improve original content preservation in style transfer. Inspired by FlexiEdit [16], which highlights that high-frequency components are more closely associated with image layout (e.g., contours and details) than low-frequency components,

we improve detail preservation by reducing low-frequency components of the DDIM latent, which serves as the sampling startpoint. This refinement enhances the model’s ability to retain the image layout during style transfer.

Second, we introduce negative guidance in the DDIM inversion stage to alleviate content leakage from style images. This approach ensures that the sampling startpoint is further “distanced” from the content of style image. Our experiments (Fig. 5) show that, compared to traditional negative guidance [28] applied during the sampling stage, applying guidance in the inversion stage yields superior results by mitigating multi-condition control failures [9, 13]. Additionally, we use the pre-trained IP-Instruct model [39] as our style and content extractor, providing negative guidance in the inversion stage for a better startpoint.

In summary, our main contributions are as follows:

- We propose a novel sampling startpoint enhancement method for training-free diffusion-based style transfer, addressing issues of content leakage from style images and changes in original content. To the best of our knowledge, we are the first work to highlight the importance of the startpoint in this area.
- We propose frequency manipulation to reduce the low-frequency components of the DDIM latent, which serves as the sampling startpoint, thereby enhancing original content preservation.
- We propose negative guidance via inversion to distance the sampling startpoint from the content of style image, thus alleviating content leakage from style image.
- Extensive experiments on the style transfer dataset vali-

date that the proposed method significantly outperforms previous works both quantitatively and qualitatively.

## 2. Related Work

### 2.1. Diffusion-Based Text-to-Image Generation

Recently, diffusion models have achieved significant success in image generation. Diffusion Probabilistic Models (DPMs) [42] are proposed to transform random noise into high-resolution images through a sequential sampling process. Many previous diffusion-based image generation works have demonstrated strong generative capabilities. Latent Diffusion Models (LDMs) [36, 38] further revolutionize this approach by operating in a compressed latent space, using a pre-trained auto-encoder [33, 34] to enhance computational efficiency and yield high-resolution images from textual descriptions. This transition to latent space not only accelerates the generation process but also improves the quality and coherence of the generated images. As text-to-image (T2I) diffusion models [18, 51] continue to grow in influence within the field of image generation, it has become evident that texts offer limited control over spatial and textural aspects of images. This has promoted the development of using more conditions from a reference image based on the T2I diffusion model [51, 53]. One of these particular conditions is style, which is the key focus of this paper.

### 2.2. Style Transfer with T2I Models

Style transfer is a condition-guided image generation task that applies the style of one image to another while preserving the original content. Early neural style transfer was extensively explored in deep convolutional [6], generative adversarial [11, 25, 56], and transformer-based networks [15, 29], marking substantial progress over traditional methods based on signal processing [5, 24]. This evolution has enabled numerous applications, particularly in advertising and marketing. With the powerful generative capacity of the T2I diffusion model, neural style transfer increasingly relies on pre-trained diffusion models to achieve style transfer. Previous methods [20, 21, 31, 55] have relied on paired datasets with shared content but different styles to learn style concepts through reconstruction. For instance, DEADiff [31] trains an additional image encoder guided by textual descriptions to separate style and content in the reference image. Although these approaches have demonstrated impressive style transfer capabilities, they are often time-consuming or fail to fully exploit the generative potential of large-scale diffusion models.

Training-free style transfer methods are gaining popularity due to their generalization and convenience. DiffStyle [12] leverages h-space and adjusts skip connections to effectively convey style and content information, respec-

tively. InstantStyle [45] integrates features from a reference style image into style-specific layers, enhancing the style transfer process. However, these approaches often encounter challenges in preserving the original image layout. Methods like StyleID [4] and InstantStyle plus [46] have underscored the importance of inversion in content preservation, designing fusion operations for intermediate features between user-provided style reconstructions and other image streams. Nonetheless, these methods still face content leakage issues from style images.

To address these issues, we propose a novel, training-free method based on the sampling startpoint enhancement by frequency manipulation and negative guidance via inversion, which avoids content leakage from style images while ensuring strong content preservation.

## 3. Preliminary

### 3.1. Diffusion Model

Stable Diffusion (SD) [38] is a type of latent diffusion model designed to map a random noise vector  $z_t$  and a text prompt  $\mathcal{P}$  to an output image  $I_0$ , aligning with the given conditioning prompt via cross-attention. The objective of this process is defined as:

$$L = \mathbb{E}_{z_0, \epsilon \sim N(0, I), t \sim Uniform(1, T)} \|\epsilon - \epsilon_\theta(z_t, t, \mathcal{C})\|_2^2, \quad (1)$$

where  $\mathcal{C} = \varphi(\mathcal{P})$  is the embedding of text prompt generated by the text encoder  $\varphi$ .  $\epsilon$  and  $\epsilon_\theta$  represent the actual and predicted noise, respectively. The noise is gradually removed by sequentially predicting it using pre-trained diffusion model.

**Classifier-Free Guidance (CFG)** [10] enhances image generation quality by using a null-text embedding  $\emptyset$ , which corresponds to the embedding of a null text “”, as a reference for unconditional predictions during sampling. The modified noise prediction is expressed as:

$$\tilde{\epsilon}_\theta(z_t, t, \mathcal{C}, \emptyset) = \epsilon_\theta(z_t, t, \emptyset) + \omega (\epsilon_\theta(z_t, t, \mathcal{C}) - \epsilon_\theta(z_t, t, \emptyset)), \quad (2)$$

where the guidance scale  $\omega \geq 0$  adjusts the strength of the conditional prediction  $\epsilon_\theta(z_t, t, \mathcal{C})$  against to the unconditional prediction  $\epsilon_\theta(z_t, t, \emptyset)$ .

**DDIM Inversion:** The Denoising Diffusion Implicit Model (DDIM) [43] is a generative model that improves image synthesis efficiency and quality through a non-Markovian diffusion process, reducing the number of steps needed to generate samples. Within SD model, deterministic DDIM sampling uses a denoiser network  $\epsilon_\theta$ , described by:

$$z_{t-1} = \sqrt{\frac{\alpha_{t-1}}{\alpha_t}} z_t + \left( \sqrt{\frac{1}{\alpha_{t-1}} - 1} - \sqrt{\frac{1}{\alpha_t} - 1} \right) \cdot \epsilon_\theta(z_t, t), \quad (3)$$

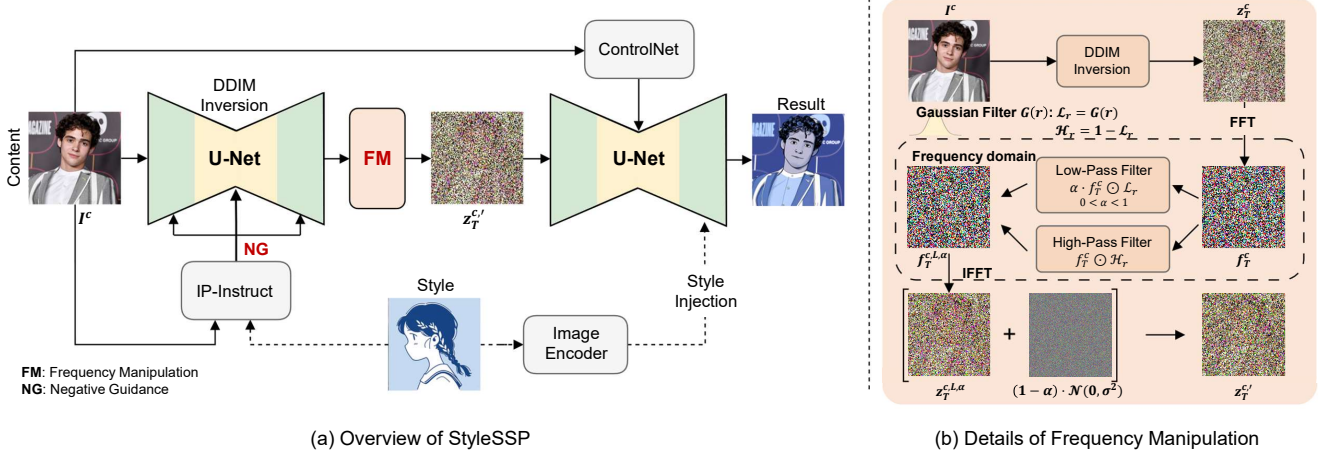


Figure 2. **Overall Framework.** (Left) Illustration of the proposed style transfer method. First, we invert the content image  $I^c$  into the latent noise space as  $z_T^c$ . During this process, we use negative guidance (Sec. 4.2) to ensure that  $z_T^c$  diverges from the content information of the style image. We then apply frequency manipulation (Sec. 4.1) to  $z_T^c$ , obtaining a low-frequency reduced latent  $z_T^{c,\alpha}$  as the startpoint for the sampling stage. During sampling, we follow InstantStyle’s approach by injecting style features exclusively into the style-specific block and utilizing the ControlNet model to further preserve original content. (Right) Detailed explanation of frequency manipulation. We reduce the low-frequency components by a factor  $\alpha$ , while adding Gaussian noise proportional to  $1 - \alpha$ .

where  $\alpha = (\alpha_1, \dots, \alpha_T) \in \mathbb{R}_{\geq 0}^T$  are hyper-parameters defining noise scales at  $T$  diffusion steps. In this work, we use the publicly available SD model [30], where the diffusion forward process is applied to a latent image encoding  $z_0 = E(I_0)$ , and an image decoder is employed at the end of the diffusion backward process  $I_0 = D(z_0)$ .

By representing the DDIM sampling equation as an ordinary differential equation (ODE), the forward process can be expressed in terms of  $\epsilon_\theta(z_t, t)$  by inverting the reverse diffusion process (DDIM Inversion) as follows:

$$z_{t+1}^* = \sqrt{\frac{\alpha_{t+1}}{\alpha_t}} z_t^* + \sqrt{\alpha_{t+1}} \left( \sqrt{\frac{1}{\alpha_{t+1}} - 1} - \sqrt{\frac{1}{\alpha_t} - 1} \right) \cdot \epsilon_\theta(z_t^*, t). \quad (4)$$

In Eq. 4,  $z_t^*$  denotes latent features during the DDIM Inversion process. Therefore, we obtain the DDIM Inversion trajectory, denoted as  $[z_t^*]_{t=0}^T$ . Recent works [4, 55] have shown that initiating DDIM sampling from  $z_T = z_T^*$  benefits to original content preservation. These findings highlight the importance of a proper startpoint for the sampling stage (denoted as  $z_T$ ), motivating our approach to guide the inversion stage and manipulate the DDIM latent  $z_T$ .

### 3.2. Frequency Analysis

Inspired by FlexiEdit [16], which highlights that high-frequency components play a more significant role in forming the object’s layout than low-frequency components, we conduct a frequency analysis on the DDIM latent  $z_T$  to explore frequency-domain operations that benefit to preserve

the original content in style transfer. Our method separates the DDIM latent  $z_T$  into high- and low-frequency components in the frequency domain as follows:

$$f_T^{L,\alpha} = \alpha * f_T \odot \mathcal{L}_r + f_T \odot \mathcal{H}_r, \text{ where } \alpha \in [0, 1], \quad (5)$$

$$f_T^{H,\alpha} = f_T \odot \mathcal{L}_r + \alpha * f_T \odot \mathcal{H}_r, \text{ where } \alpha \in [0, 1], \quad (6)$$

$$z_T^{L,\alpha} = IFFT(f_T^{L,\alpha}), \quad z_T^{H,\alpha} = IFFT(f_T^{H,\alpha}), \quad (7)$$

where  $FFT(\cdot)$  and  $IFFT(\cdot)$  denote the 2D Fast Fourier Transform and its inverse, respectively;  $f_T$  represents the frequency spectrum of  $z_T$ ;  $\mathcal{L}_r$  is a low-pass filter (e.g., Gaussian, Butterworth, or Chebyshev), and  $\mathcal{H}_r = 1 - \mathcal{L}_r$  is the corresponding high-pass filter. Here,  $\odot$  denotes element-wise multiplication.

Since  $\alpha \in [0, 1]$ ,  $z_T^{L,\alpha}$  and  $z_T^{H,\alpha}$  represent low- and high-frequency reduced latents, respectively, with reduction degrees adjusted by the scale  $\alpha$ . In Fig. 3, we observe that as  $\alpha$  increases in reconstructions from  $z_T^{H,\alpha}$ , content preservation effects improve significantly. Conversely, reconstructions from  $z_T^{L,\alpha}$  consistently maintain layout accuracy across varying  $\alpha$  values, indicating that high-frequency components in  $z_T$  are more crucial in determining the layout of image.

### 4. Method

Based on our discoveries (shown in supplementary materials Sec. 7.1), which highlight the importance of a better sampling startpoint, we propose a sampling startpoint enhancement method called **StyleSSP** for training-free diffusion-based style transfer, shown in Fig. 2. Focus-

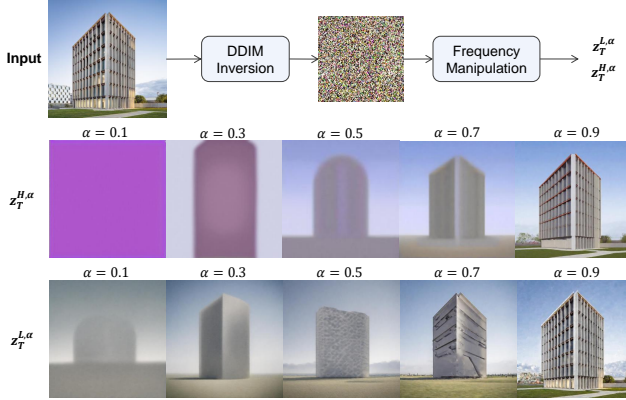


Figure 3. Reconstruction results with varying  $\alpha$  values, demonstrating that high-frequency components play a critical role in the image layout, while low-frequency components contribute less to layout preservation.

ing on the problems of original content changes and content leakage from style images in current training-free methods, StyleSSP proposes two main components: (1) Frequency Manipulation (Sec. 4.1) and (2) Negative Prompt Guidance via Inversion (Sec. 4.2).

Let  $I^c$  be a given content image whose text prompt  $\mathcal{P}$  is generated by BLIP [19]. Our goal is to modify the style of  $I^c$  to that of style image  $I^s$ . The generated styled image  $I^{cs}$  will maintain the content of  $I^c$  while its style is consistent with  $I^s$ . In the following sections, we refer to the content, style, and stylized images as their encoded counterparts  $z_0^c$ ,  $z_0^s$ , and  $z_0^{cs}$ , respectively.

#### 4.1. Frequency Manipulation

Frequency analysis in Sec. 3.2 indicates that high-frequency components within DDIM latent  $z_T^c$  of content image are more crucial in determining the layout of original image than low-frequency components. Based on this, we manipulate the frequency components of DDIM latent  $z_T^c$  by a high-pass filter, which can achieve better preservation of original layout, resulting in improvement of details representation in the generated image.

To this end, we first obtain the latent of content image with DDIM Inversion, and then filter the DDIM latent  $z_T^c$  to get the low-frequency reduced latent  $z_T^{c,L,\alpha}$ , which more tightly bound with the layout of image.

$$z_T^c = \text{DDIM-Inv}(z_0^c), \quad (8)$$

$$z_T^{c'} = z_T^{c,L,\alpha} + \mathcal{N}(0, \sigma^2) * (1 - \alpha), \quad (9)$$

where the definition of  $z_T^{c,L,\alpha}$  is given in Eq. 7, denoting the low-frequency reduced DDIM latent of content image. This procedure selectively reduces the low-frequency components by factor  $\alpha$  and introduces Gaussian noise scaled

by  $1 - \alpha$ , resulting in a manipulated latent  $z_T^{c'}$ . As shown in Fig. 4, we illustrate the importance of frequency manipulation for preserving the background details of image.

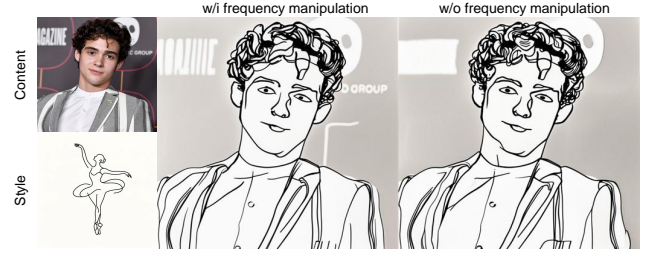


Figure 4. Style transfer results w/o frequency manipulation, representing the detail preservation enhancement of frequency manipulation. Result with frequency manipulation outperforms in keeping the text and lines in the background.

#### 4.2. Negative Guidance via Inversion

To distance the sampling startpoint from the content of style image, we draw from insights in previous negative guidance methods. Negative prompt guidance [28], introduced in conditional generation models such as SD, allows users to specify what to exclude from generated images. This approach has gained significant attention for its effectiveness [1, 47]. Specifically, when the null-text embedding  $\emptyset$  in the unconditional format is replaced with an actual prompt, it represents what to remove from the generated image, leveraging the negative sign. This can be formally expressed as:

$$\hat{\epsilon}_\theta(z_t, t, \mathcal{C}_+, \mathcal{C}_-) = \epsilon_\theta(z_t, t, \mathcal{C}_-) + \omega_i(\epsilon_\theta(z_t, t, \mathcal{C}_+) - \epsilon_\theta(z_t, t, \mathcal{C}_-)), \quad (10)$$

where  $\mathcal{C}_+ = \varphi(\mathcal{P}_+)$  and  $\mathcal{C}_- = \varphi(\mathcal{P}_-)$  are the embedding of positive text prompt  $\mathcal{P}_+$  and negative text prompt  $\mathcal{P}_-$ , respectively.  $\omega_i$  is the negative guidance scale. More details on the principles of negative prompt guidance can be found in supplementary materials Sec. 7.2.

Although negative prompts provide additional control, they may interfere with the original prompt or even be disregarded [2], requiring careful tuning by users. Furthermore, the expressive capacity of text is inherently constrained, particularly for style transfer, where it is nearly impossible to comprehensively capture an image's content or precisely describe its style with words alone. These limitations substantially reduce the effectiveness of negative prompt guidance. To address this issue, we leverage the pre-trained IP-Instruct model [39] as a content and style extractor. The embeddings from this extractor serve as negative guidance, allowing us to overcome the challenges of accurately representing style and content information.

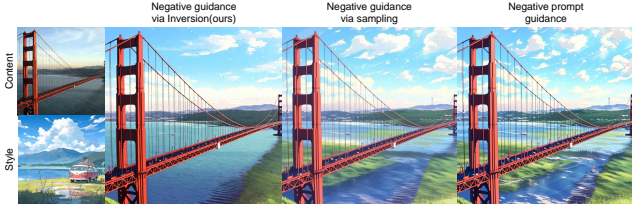


Figure 5. Illustrations of negative guidance via inversion, negative guidance in sampling step and negative prompt guidance results for style transfer. The latter two all face severe content leakage problems (the out-of-place grass on the river), while our method prevents this phenomenon.

$$\hat{\epsilon}_{\theta}(z_t, t, \mathcal{C}_+, \mathcal{E}_-) = \epsilon_{\theta}(z_t, t, \mathcal{E}_-) + \omega_i(\epsilon_{\theta}(z_t, t, \mathcal{C}_+) - \epsilon_{\theta}(z_t, t, \mathcal{E}_-)), \quad (11)$$

where  $\mathcal{E}_- = \text{concat}(\Phi(I^c)^s, \Phi(I^s)^c)$ .  $\Phi(I^s)^c$  denotes the content embedding of style image  $I^s$ , and  $\Phi(I^c)^s$  denotes style embedding of content image  $I^c$ .  $\Phi$  is the IP-Instruct model to extract style and content information.

Notably, based on our significant discovery, which highlights the importance of sampling startpoint for style transfer, we innovatively employ negative guidance during DDIM Inversion. We utilize  $\hat{\epsilon}_{\theta}(z_t, t, \mathcal{C}_+, \mathcal{E}_-)$  in Eq. 11 to replace the  $\epsilon_{\theta}(z_t^*, t)$  in Eq. 4, presenting the predicted noises that are added into the content image gradually. As shown in Fig. 5, the negative guidance via inversion outperforms both the traditional negative prompt guidance and the negative guidance in the sampling stage. This demonstrates that negative guidance via inversion can prevent content leakage by keeping the startpoint away from the content of style image.

### 4.3. Injection & Controlling

**Style Injection:** Previous studies [23, 52] have demonstrated that each layer of a deep network captures different types of semantic information, which informs the style injection strategy. This approach focuses on injecting style solely into the blocks responsible for style generation in the U-Net architecture, thereby preventing content leakage. This strategy is supported by findings from InstantStyle [45], which show that the first upsampling block of U-Net primarily captures style-related features such as color, material, and atmosphere. Consequently, in this work, we concentrate on injecting style features into a specific block to achieve seamless style transfer, in line with the approach used in InstantStyle.

**ControlNet for Content Preservation:** ControlNet has become one of the most widely adopted techniques for spatial conditioning, including for canny edges, depth maps, human poses, and more. In this work, we utilize ControlNet model to help preserve the layout of the content image,

thereby enabling more precise control over the original content during style transfer.

## 5. Experiments

### 5.1. Experimental Settings

We conduct all experiments in pre-trained Stable Diffusion XL [30] and tile ControlNet [53], as well as adopt DDIM inversion and sampling with a total 50 timesteps ( $t = \{1, \dots, 50\}$ ). The negative guidance operates with guidance scale  $\omega_i$  equal to 1.5 while the CFG scale for sampling stage is set to 5.0. We use the Gaussian filter with variance  $\sigma$  equal to 0.3 in frequency manipulation, and determine the scale value  $\alpha$  to be 0.7. We utilize ViT-L/14 from CLIP [35] as the image encoder. All the experiments are conducted on an NVIDIA A100 GPU.

**Dataset:** Our evaluations employ content images from MS-COCO [22] dataset and style image from WikiArt [44] dataset. For quantitative comparison, we randomly selected content and style images from each dataset, generating 800 stylized images.

**Evaluation metric:** We employ the evaluation metric ArtFID [48], LPIPS [54] and FID [41], consistent with StyleID. ArtFID evaluates overall style transfer performances with consideration of both content and style preservation and also is known as strongly coinciding with human judgment, which is computed as  $ArtFID = (1 + LPIPS) \cdot (1 + FID)$ . LPIPS measures content fidelity between the stylized image and the corresponding content image, and FID assesses the style fidelity between the stylized image and the corresponding style image.

### 5.2. Qualitative Results

Fig. 6 presents the superior style transfer results of StyleSSP across various subjects, demonstrating its robustness and versatility in adapting to diverse content and styles. The results show that our method not only performs straightforward color transfer but also captures more distinctive features, such as brush strokes and textures from the style image, leading to visually appealing style transfer effects. Additional results can be found in the supplementary materials Sec. 7.3.

### 5.3. Comparison with State-of-the-Art Methods

We evaluate our proposed method by comparing it with previous state-of-the-art methods, including training-free diffusion-based methods such as StyleID [4], StyleAlign [20], InstantStyle plus [46], InstantStyle [45], DiffuseIT [17], and DiffStyle [12]. Additionally, we also include the optimization-based method InST [55] in our comparison, based on the experimental settings of StyleID.

**Quantitative Comparisons:** As shown in Tab. 1, our method outperforms previous style transfer methods in



Figure 6. Style transfer results of style and content image pairs. Zoom in for better visualization.



Figure 7. Qualitative comparison with previous work.

terms of ArtFID, FID, and LPIPS, indicating superior style resemblance and content fidelity. Several key observations can be made from this comparison. First, when compared to content preservation methods such as InstantStyle plus, StyleID, and InST, our approach achieves the best LPIPS score, demonstrating a significant improvement in content preservation. Second, our method also achieves the lowest FID, highlighting its superior style transfer performance. In summary, StyleSSP strikes an optimal balance between high-quality style transfer and precise content preservation. **Qualitative Comparisons:** Fig. 7 presents a visual com-

parison between our method and previous works. Overall, our approach achieves the best visual balance between enhancing stylistic effects and preserving the original content, while effectively preventing the content leakage from style image. Several key observations can be made from this figure. First, methods without inversion exhibit significant limitations in content preservation, particularly in the background details, as shown in the 1<sup>st</sup> row. Second, although inversion-based methods such as StyleID, InST, and InstantStyle plus present some content preservation ability, they fail to fully decouple style and content information.

Metric	Ours	StyleID	InstantStyle plus	StyleAlign	InstantStyle	DiffuseIT	InST	DiffStyle
ArtFID↓	<b>21.499</b>	28.801	25.886	36.269	37.524	40.721	40.633	41.464
FID↓	<b>13.448</b>	18.131	16.097	20.338	21.817	23.065	21.571	20.903
LPIPS↓	<b>0.4881</b>	0.5055	0.5140	0.6997	0.6446	0.6921	0.8002	0.8931

Table 1. Quantitative comparison with diffusion model baselines

This results in visible content leakage in some synthesized images, especially in the 4<sup>th</sup> row of Fig. 7. If users interpret the waves in the 4<sup>th</sup> row as part of the style, we show in Sec. 5.5 that content leakage can be controlled by adjusting the negative guidance scale  $\omega_i$ , allowing users to customize the result according to their preferences. Additional results are provided in the supplementary materials Sec. 7.3.

#### 5.4. Ablation Study

Configuration	ArtFID↓	FID↓	LPIPS↓
Baseline	26.683	16.205	0.5509
+ FM	24.112	15.103	0.4973
+ NG	26.542	16.128	0.5496
StyleSSP	<b>21.499</b>	<b>13.448</b>	<b>0.4881</b>

Table 2. Quantitative results from gradually increasing components with StyleSSP.

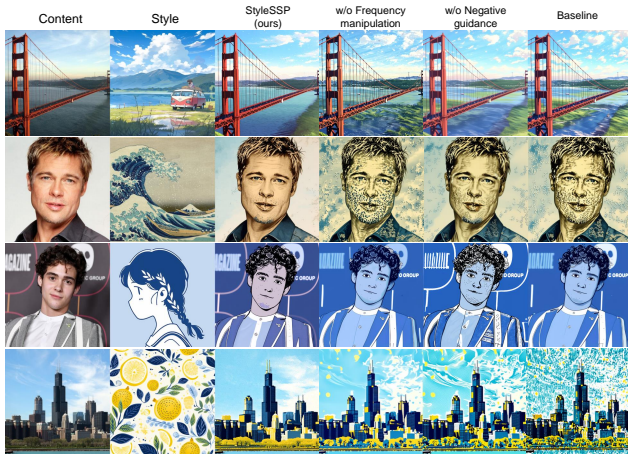


Figure 8. Qualitative comparison with ablation studies.

To validate the effectiveness of the proposed components, we conduct ablation studies from both quantitative and qualitative perspectives. The baseline refers to the method without frequency manipulation (FM) and negative guidance via inversion (NG). Qualitative results, as shown in Fig. 8, illustrate the effects of frequency manipulation for content preservation and negative guidance for preventing content leakage. First, referring to the 3<sup>rd</sup> row of Fig. 8,

frequency manipulation significantly improves the preservation of background details in the content image. Second, referring to the 1<sup>st</sup>, 2<sup>nd</sup>, 4<sup>th</sup> rows, negative guidance effectively prevents the contamination of content by style images in the generated images. By guiding the startpoint distance from the content of style image, negative guidance successfully prevents the contamination of river, human faces, and sky in the original images by the grassland, waves, and yellow dots from style images. Quantitative results shown in Tab. 2 further demonstrate the superior performance of our proposed components. In summary, our method excels in both visual effects and quantitative metrics.

#### 5.5. Additional Analysis



Figure 9. Visualization of the effects of negative guidance scale  $\omega_i$  and frequency manipulation ratio  $\alpha$ .

We investigate the effects of different negative guidance scales  $\omega_i$  and frequency manipulation ratio  $\alpha$ . We observe that the gradual increase of  $\omega_i$  reduces the degree of content leakage from style image, as shown in Fig. 9 (top). This result further implies that negative guidance is effective in mitigating content leakage. In addition, as shown in Fig. 9 (bottom), a lower frequency manipulation ratio  $\alpha$  results in stylized images with clearer contours and more defined layouts, highlighting the importance of reducing low-frequency components in the startpoint for enhancing image structure and detail. This characteristic suggests that users can adjust the degree of contour sharpness and content leakage based on their preferences.

### 6. Conclusion

In this paper, we introduce **StyleSSP**, a novel method for sampling startpoint enhancement in training-free diffusion-based style transfer. To the best of our knowledge, we are



the first to emphasize the importance of the sampling startpoint in style transfer. We identify two key challenges in training-free methods: changes of original content and content leakage from style images. These issues stem primarily from the absence of targeted training for style extraction and constraints on content layout. To address these issues, we propose two components for optimizing the sampling startpoint: (1) frequency manipulation for improved content preservation, and (2) negative guidance via inversion to prevent content leakage. Empirical results demonstrate that StyleSSP effectively mitigates original content changes and content leakage from style image while achieving superior style transfer performance. Comparison experiments show that StyleSSP outperforms previous methods both qualitatively and quantitatively. Future work could explore regionally-aware startpoint manipulation techniques to further enhance objective-level stylization.

## References

- [1] Mohammadreza Armandpour, Ali Sadeghian, Huangjie Zheng, Amir Sadeghian, and Mingyuan Zhou. Re-imagine the negative prompt algorithm: Transform 2d diffusion into 3d, alleviate janus problem and beyond, 2023. 5
- [2] Yuanhao Ban, Ruochen Wang, Tianyi Zhou, Minhao Cheng, Boqing Gong, and Cho-Jui Hsieh. Understanding the impact of negative prompts: When and how do they take effect?, 2024. 5
- [3] Rumeysa Bodur, Erhan Gundogdu, Binod Bhattarai, Tae-Kyun Kim, Michael Donoser, and Loris Bazzani. iredit: Localised text-guided image editing with weak supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 7426–7435, 2024. 1
- [4] Jiwoo Chung, Sangeek Hyun, and Jae-Pil Heo. Style injection in diffusion: A training-free approach for adapting large-scale diffusion models for style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8795–8805, 2024. 1, 2, 3, 4, 6
- [5] Alexei A. Efros and William T. Freeman. Image quilting for texture synthesis and transfer. In *Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques*, page 341–346, New York, NY, USA, 2001. Association for Computing Machinery. 3
- [6] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. Texture synthesis using convolutional neural networks, 2015. 3
- [7] Zigang Geng, Binxin Yang, Tiankai Hang, Chen Li, Shuyang Gu, Ting Zhang, Jianmin Bao, Zheng Zhang, Han Hu, Dong Chen, and Baining Guo. Instructdiffusion: A generalist modeling interface for vision tasks. *CoRR*, abs/2309.03895, 2023. 1
- [8] Zigang Geng, Binxin Yang, Tiankai Hang, Chen Li, Shuyang Gu, Ting Zhang, Jianmin Bao, Zheng Zhang, Houqiang Li, Han Hu, Dong Chen, and Baining Guo. Instructdiffusion: A generalist modeling interface for vision tasks. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12709–12720, 2024. 1
- [9] Yucheng Han, Rui Wang, Chi Zhang, Juntao Hu, Pei Cheng, Bin Fu, and Hanwang Zhang. Emma: Your text-to-image diffusion model can secretly accept multi-modal prompts, 2024. 1, 2
- [10] Jonathan Ho. Classifier-free diffusion guidance. *ArXiv*, abs/2207.12598, 2022. 3
- [11] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. *CVPR*, 2017. 3
- [12] Jaeseok Jeong, Mingi Kwon, and Youngjung Uh. Training-free content injection using h-space in diffusion models, 2024. 1, 3, 6
- [13] Xuan Ju, Xian Liu, Xintao Wang, Yuxuan Bian, Ying Shan, and Qiang Xu. Brushnet: A plug-and-play image inpainting model with decomposed dual-branch diffusion, 2024. 1, 2
- [14] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *Conference on Computer Vision and Pattern Recognition 2023*, 2023. 1
- [15] Nicholas Kolkin, Jason Salavon, and Greg Shakhnarovich. Style transfer by relaxed optimal transport and self-similarity, 2019. 3
- [16] Gwanhyeong Koo, Sunjae Yoon, Ji Woo Hong, and Chang D Yoo. Flexiedit: Frequency-aware latent refinement for enhanced non-rigid editing. *arXiv preprint arXiv:2407.17850*, 2024. 2, 4
- [17] Gihyun Kwon and Jong Chul Ye. Diffusion-based image translation using disentangled style and content representation, 2023. 6
- [18] Dongxu Li, Junnan Li, and Steven C. H. Hoi. Blip-diffusion: Pre-trained subject representation for controllable text-to-image generation and editing, 2023. 3
- [19] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation, 2022. 5
- [20] Senmao Li, Joost van de Weijer, Taihang Hu, Fahad Shahbaz Khan, Qibin Hou, Yaxing Wang, and Jian Yang. Stylediffusion: Prompt-embedding inversion for text-based editing. *arXiv preprint arXiv:2303.15649*, 2023. 1, 3, 6
- [21] Sijia Li, Chen Chen, and Haonan Lu. Moecontroller: Instruction-based arbitrary image manipulation with mixture-of-expert controllers, 2024. 1, 3
- [22] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2015. 6
- [23] Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting them, 2014. 6
- [24] Yifang Men, Zhouhui Lian, Yingmin Tang, and Jianguo Xiao. A common framework for interactive texture transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6353–6362, 2018. 3
- [25] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *CoRR*, abs/1411.1784, 2014. 3

- [26] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models, 2022. 1
- [27] Alex 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, 2022. 1
- [28] Ryan O'Connor. Stable diffusion 1 vs 2: What you need to know. <https://www.assemblyai.com/blog/stable-diffusion-1-vs-2-what-you-need-to-know>, 2023. 2, 5
- [29] Dae Young Park and Kwang Hee Lee. Arbitrary style transfer with style-attentional networks, 2019. 3
- [30] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis, 2023. 4, 6
- [31] Tianhao Qi, Shancheng Fang, Yanze Wu, Hongtao Xie, Jiawei Liu, Lang Chen, Qian He, and Yongdong Zhang. Deadiff: An efficient stylization diffusion model with disentangled representations. *arXiv preprint arXiv:2403.06951*, 2024. 1, 3
- [32] Haonan Qiu, Menghan Xia, Yong Zhang, Yingqing He, Xintao Wang, Ying Shan, and Ziwei Liu. Freenoise: Tuning-free longer video diffusion via noise rescheduling, 2024. 2
- [33] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 3
- [34] Colin Raffel, Minh-Thang Luong, Peter J. Liu, Ron J. Weiss, and Douglas Eck. Online and linear-time attention by enforcing monotonic alignments, 2017. 3
- [35] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents, 2022. 6
- [36] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents, 2022. 3
- [37] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021. 1
- [38] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2022. 3
- [39] Ciara Rowles, Shimon Vainer, Dante De Nigris, Slava Elizarov, Konstantin Kutsy, and Simon Donné. Ipadapter-instruct: Resolving ambiguity in image-based conditioning using instruct prompts, 2024. 2, 5
- [40] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding, 2022. 1
- [41] Maximilian Seitzer. pytorch-fid: FID Score for PyTorch. <https://github.com/mseitzer/pytorch-fid>, 2020. Version 0.3.0. 6
- [42] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. *CoRR*, abs/1503.03585, 2015. 3
- [43] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv:2010.02502*, 2020. 3
- [44] Wei Ren Tan, Chee Seng Chan, Hernan Aguirre, and Kiyoshi Tanaka. Improved artgan for conditional synthesis of natural image and artwork, 2018. 6
- [45] Haofan Wang, Qixun Wang, Xu Bai, Zekui Qin, and Anthony Chen. Instantstyle: Free lunch towards style-preserving in text-to-image generation. *arXiv preprint arXiv:2404.02733*, 2024. 1, 3, 6
- [46] Haofan Wang, Peng Xing, Renyuan Huang, Hao Ai, Qixun Wang, and Xu Bai. Instantstyle-plus: Style transfer with content-preserving in text-to-image generation. *arXiv preprint arXiv:2407.00788*, 2024. 1, 3, 6
- [47] Max Woolf. Stable diffusion 2.0 and the importance of negative prompts for good results. <https://minimaxir.com/2022/11/stable-diffusion-negative-prompt/>, 2023. 5
- [48] Matthias Wright and Björn Ommer. Artfid: Quantitative evaluation of neural style transfer. *GCPR*, 2022. 6
- [49] Zongze Wu, Yotam Nitzan, Eli Shechtman, and Dani Lischinski. Stylealign: Analysis and applications of aligned stylegan models. *arXiv preprint arXiv:2110.11323*, 2021. 1
- [50] Youcan Xu, Zhen Wang, Jun Xiao, Wei Liu, and Long Chen. Freetuner: Any subject in any style with training-free diffusion, 2024. 2
- [51] Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models, 2023. 3
- [52] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks?, 2014. 6
- [53] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023. 1, 3, 6
- [54] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric, 2018. 6
- [55] Yuxin Zhang, Nisha Huang, Fan Tang, Haibin Huang, Chongyang Ma, Weiming Dong, and Changsheng Xu. Inversion-based style transfer with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10146–10156, 2023. 1, 3, 4, 6
- [56] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networkss. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, 2017. 3

# StyleSSP: Sampling StartPoint Enhancement for Training-free Diffusion-based Method for Style Transfer

## Supplementary Material

### 7. Appendix

#### 7.1. Startpoint Impact Analysis

Given that StyleSSP is specifically designed to enhance the sampling startpoint, we place primary emphasis on the importance of the startpoint in style transfer. We demonstrate how minor modifications to the startpoint can significantly influence style transfer results. As shown in Fig. 10, we present several style transfer results. The titles in the figure — “wi Inversion,” “wo Inversion,” “Noised Latent,” “Shifted Latent,” and “Scaled Latent” — correspond to the startpoints  $z_T$ ,  $z_r$ ,  $z_T^n$ ,  $z_T^{sh}$ , and  $z_T^{sa}$ , respectively. Their formulations are as follows:

$$\begin{aligned} z_r &\sim \mathcal{N}(0, \mathbf{I}) \\ z_T^n &= z_T + \mathcal{N}(0, \mathbf{I}) \\ z_T^{sh} &= z_T + \mathbf{U}(-0.5, 0.5) \\ z_T^{sa} &= z_T \times \mathbf{U}(0.5, 1) \end{aligned} \quad (12)$$

where  $z_T$  is the DDIM latent of the content image,  $\mathcal{N}$  represents a Gaussian distribution, and  $U(-0.5, 0.5)$  and  $U(0.5, 1)$  indicate uniformly random values selected within the ranges -0.5 to 0.5 and 0.5 to 1.0, respectively.

As illustrated in Fig. 10, manipulations of the sampling startpoint make a significant impact on the results of style transfer, resulting in notable changes in both the image hue and the content representation. Note that the following results are all conducted with ControlNet as an additional content controller. Several key observations can be made from this figure.

First, referring to the 3<sup>rd</sup> and 4<sup>th</sup> columns in this figure, using the DDIM latent  $z_T$  extracted from the content image as the sampling startpoint results in remarkably better content preservation compared to using random Gaussian noise as the startpoint. This finding motivates us to adopt DDIM inversion as the first step in our method, as is done in many inversion-based methods [4, 46, 55].

Second, we attempted minor modifications to the DDIM latent  $z_T$ . Referring to the 3<sup>rd</sup>, 5<sup>th</sup>, and 6<sup>th</sup> columns in this figure, we observe that these simple manipulations produce significant changes in image tone, and since color variation is a crucial aspect of style transfer, this finding further drives our focus on startpoint enhancement.

Third, by examining the results in the 3<sup>rd</sup> and 5<sup>th</sup> rows, we notice that the startpoint not only affects the tone of generated images but can also influence the content of generated images to some extent, such as the facial outline of the woman in the 3<sup>rd</sup> row and the background in the 5<sup>th</sup> row. This effect has been largely overlooked in previous works, yet it is undeniably critical for style transfer tasks.

In summary, through simple adjustments to the startpoint, we have discovered its substantial impact on style transfer results — affecting content preservation, content modification, and tonal

changes. These insights have driven us to pursue sampling startpoint enhancement for style transfer research. Therefore, our method, StyleSSP, emphasizes guidance during the inversion step and manipulation of the inversion latent space to achieve a more effective sampling startpoint in style transfer issues.

#### 7.2. Principle of Negative Guidance

In this section, we provide a detailed introduction to the principles of negative prompt guidance, starting with conditional generation. For conditional generation, that is, to sample samples from the conditional distribution  $p(x|y)$ . According to the Bayes formula, we can obtain:

$$\begin{aligned} p(x|y) &= \frac{p(y|x)p(x)}{p(y)}, \\ \log p(x|y) &= \log p(y|x) + \log p(x) - \log p(y), \\ \Rightarrow \nabla_x \log p(x|y) &= \nabla_x \log p(y|x) + \nabla_x \log p(x). \end{aligned} \quad (13)$$

In the classifier-guided task, the score-based model with unconditional input is an estimation of  $\nabla_x \log p(x)$ , so in order to obtain  $\nabla_x \log p(x|y)$ , an additional classifier needs to be trained to estimate  $\nabla_x \log p(y|x)$ . At the same time, to control the strength of condition, the guidance scale  $\omega$  is introduced:

$$\nabla_x \log p(x|y) := \omega \nabla_x \log p(y|x) + \nabla_x \log p(x). \quad (14)$$

In classify-free guidance (CFG) tasks, they simultaneously train two score-based models,  $\nabla_x \log p(x)$  and  $\nabla_x \log p(y|x)$ . Since  $\nabla_x \log p(y|x) = \nabla_x \log p(x|y) - \nabla_x \log p(x)$ , it follows that:

$$\nabla_x \log p(x|y) := \omega (\nabla_x \log p(x|y) - \nabla_x \log p(x)) + \nabla_x \log p(x), \quad (15)$$

When negative prompt serves as a condition, the conditions for diffusion model contain two items, one is positive prompt condition  $y$ , and the other is negative prompt condition not  $\tilde{y}$ . Since re-training a score-based model to estimate  $\nabla_x p(x|y, \text{not } \tilde{y})$  is costly, the following simplification is made:

$$\begin{aligned} p(x|y, \text{not } \tilde{y}) &= \frac{p(x, y, \text{not } \tilde{y})}{p(y, \text{not } \tilde{y})} \\ &= \frac{p(y|x)p(\text{not } \tilde{y}|x)p(x)}{p(y, \text{not } \tilde{y})} \\ &\propto \frac{p(x)}{p(y, \text{not } \tilde{y})} \frac{p(y|x)}{p(\tilde{y}|x)}, \end{aligned} \quad (16)$$

so that:

$$\begin{aligned} \nabla_x \log p(x|y, \text{not } \tilde{y}) &\propto \nabla_x \log p(x) \\ &\quad + \nabla_x \log p(y|x) - \nabla_x \log p(\tilde{y}|x). \end{aligned} \quad (17)$$

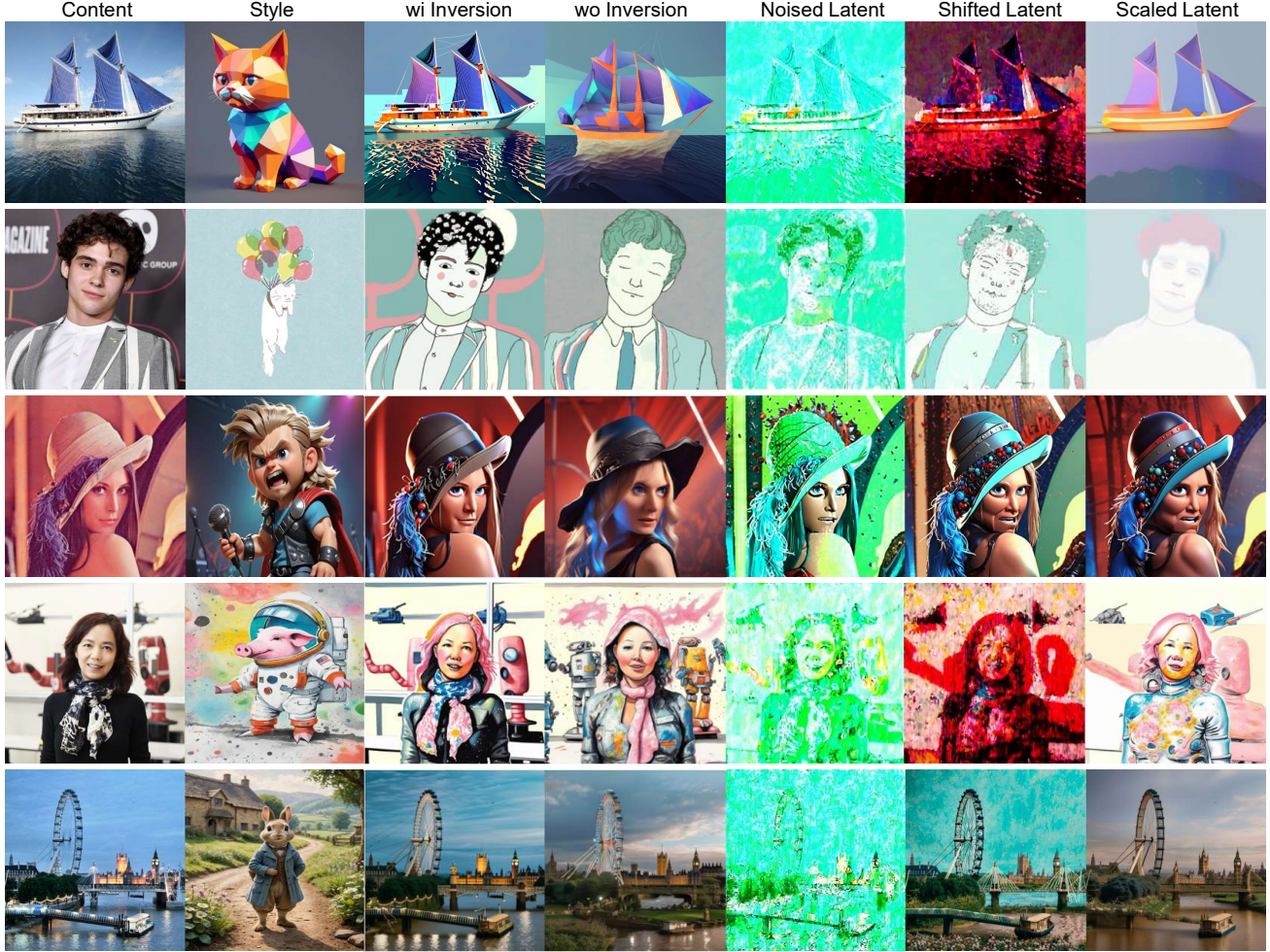


Figure 10. Illustrations of style transfer results based on various startpoints. As shown in this figure, startpoint manipulations yield significant changes in both image hue and content representation, underscoring the crucial role of the sampling startpoint in style transfer. All results are generated with ControlNet as an additional content controller.

The Eq. 16 and Eq. 17 assume that  $x$ ,  $y$  and not  $\tilde{y}$  are mutually independent. Letting  $\omega^+$  be the guidance scale of positive condition and  $\omega^-$  be the guidance scale of negative condition, we have:

$$\begin{aligned} \nabla_x p(x|y, \text{not } \tilde{y}) := & \nabla_x p(x) + \omega^+ (\nabla_x p(x|y) - \nabla_x p(x)) \\ & - \omega^- (\nabla_x p(x|\tilde{y}) - \nabla_x p(x)). \end{aligned} \quad (18)$$

Thus, we can estimate  $\nabla_x p(x|y, \text{not } \tilde{y})$  only by calculating  $\nabla_x p(x)$ ,  $\nabla_x p(x|y)$ ,  $\nabla_x p(x|\tilde{y})$ , and all of these can be obtained through the pre-trained diffusion model.

It should be noted that in the negative guidance method proposed in this paper, IP-Instruct merely exists as a style and content extractor, which can be replaced by any other extractor. Meanwhile, this CFG-based guidance can also be replaced by the gradient-based guidance like FreeTune [50] does. We emphasize that our prominent contribution lies in discovering that by guiding the startpoint of sampling stage to distance from the style image’s content, thereby preventing the content leakage from style image.

### 7.3. Additional Results

We additionally compare the proposed method with the most recent baseline (StyleID) and the baseline with lowest ArtFID (InstantStyle plus). Fig. 11 shows the additionally qualitative comparison of ours with diffusion model baselines.

Also, in Fig. 12, we visualize the style transfer results of various pairs of content and style images, which further demonstrate StyleSSP’s robustness and versatility in adapting to diverse content and style.

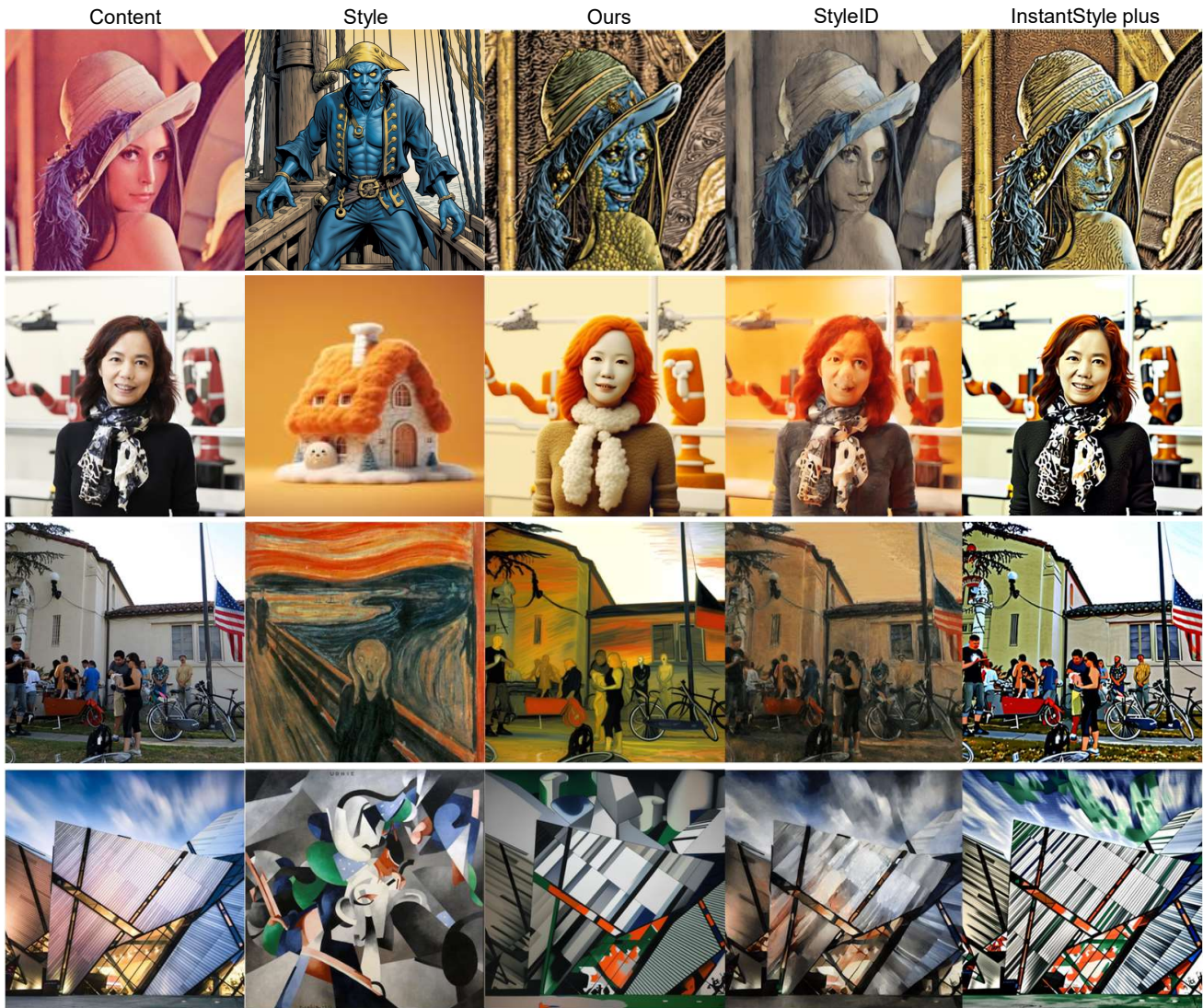


Figure 11. Qualitative comparison with with baselines(StyleID, InstantStyle plus). Zoom in for viewing details.



Figure 12. Style transfer results of style and content image pairs. Zoom in for viewing details.