

# StyleStudio: Text-Driven Style Transfer with Selective Control of Style Elements

Mingkun Lei<sup>1</sup> Xue Song<sup>2</sup> Beier Zhu<sup>1,3</sup> Hao Wang<sup>4</sup> Chi Zhang<sup>1\*</sup>

<sup>1</sup> AGI Lab, Westlake University <sup>2</sup> Fudan University <sup>3</sup> Nanyang Technological University

<sup>4</sup> The Hong Kong University of Science and Technology (Guangzhou)

{leimingkun, chizhang}@westlake.edu.cn xuesong21@m.fudan.edu.cn

beier.zhu@ntu.edu.sg haowang@hkust-gz.edu.cn

<https://stylestudio-official.github.io/>



Figure 1. Results of our text-driven style transfer model. Given a style reference image, our method effectively reduces style overfitting, generating images that faithfully align with the text prompt while maintaining consistent layout structure across varying styles.

## Abstract

*Text-driven style transfer aims to merge the style of a reference image with content described by a text prompt. Recent advancements in text-to-image models have improved the nuance of style transformations, yet significant challenges remain, particularly with overfitting to reference styles, limiting stylistic control, and misaligning with textual content. In this paper, we propose three complementary strategies to address these issues. First, we introduce a cross-modal Adaptive Instance Normalization (AdaIN) mechanism for better integration of style and text features, enhancing alignment. Second, we develop a Style-based Classifier-Free Guidance (SCFG) approach that enables selective control over stylistic elements, reducing irrelevant influences. Finally, we incorporate a teacher model during early generation stages to stabilize spatial layouts and mitigate artifacts. Our extensive evaluations demonstrate significant improvements in style transfer quality and alignment with textual prompts. Furthermore, our approach can be integrated into existing style transfer frameworks without fine-tuning.*

\* denotes Corresponding author

## 1. Introduction

Text-driven style transfer is an important task in image synthesis, aiming to blend the style of a reference image with the content aligned to a text prompt. Recent advancements in text-to-image generative models, such as Stable Diffusion [6, 23, 27], have enabled nuanced style transformations pertaining to the reference image while preserving content fidelity. This technique holds significant practical value, particularly in domains such as digital painting, advertising, and game design.

Nevertheless, modern style transfer techniques still fall short of expectations due to the inherent ambiguity in defining “style.” A style image encompasses various elements, including color palettes, textures, lighting, and brush strokes, all of which shape its overall aesthetic. Existing models often replicate all these elements, which can inadvertently lead to overfitting, where the generated output overly mirrors the characteristics of the reference style image. This over-replication of details not only diminishes the aesthetic flexibility of the generated image but also restricts its adaptability to different stylistic or content-based requirements. Therefore, an ideal style transfer approach would thus allow for more selective stylistic adjustments,

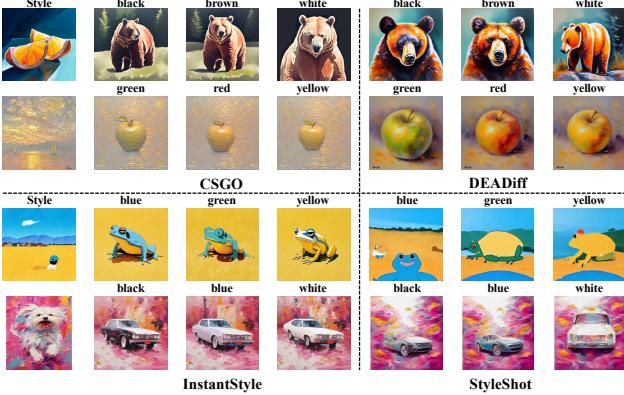


Figure 2. Illustration of overfitting issues in text-to-image generation, where the model tends to follow dominant colors or patterns from the style image rather than aligning precisely with the text prompt. Each prompt follows the format “A <color> <object>.” From top to bottom, the objects are: bear, apple, frog, and car.



Figure 3. Illustration of the checkerboard artifact encountered in the CSGO [37] method during inference. The leftmost column shows the results generated by SDXL [23]. The prompts, from top to bottom, are “A red apple” and “A pink cup.” All generated results use the same initial noise latent.

granting the user the flexibility to emphasize or omit specific stylistic components to achieve a balanced, intentional transformation.

Another challenge arising from overfitting is the difficulty in maintaining text alignment accuracy during text-to-image generation. As shown in Fig. 2, current models often prioritize dominant colors or patterns from the style image, even if they contradict instructions specified in the text prompt. This rigidity undermines the model’s ability to interpret and incorporate nuanced textual guidance, resulting in a decreased capacity for precision and customization in the generated output.

Finally, style transfer can introduce undesirable artifacts, destabilizing the underlying text-to-image generation models. One common artifact, as shown in Fig. 3, is the layout instability (*e.g.*, checkerboard effect), wherein repetitive patterns inadvertently emerge throughout the generated image, irrespective of user instructions. This highlights the unique challenges introduced by the additional complexity of style transfer.

In this paper, we propose three complementary strategies to address these challenges. First, to mitigate conflicts

between the text prompt and the style reference during generation, we introduce a mechanism where the style image features are integrated into the text features using Adaptive Instance Normalization (AdaIN) before merging them with the image features. This adaptive integration creates a more cohesive guiding feature that subsequently guides the final image generation, aligning the stylistic features with the text-based instructions more harmoniously. Second, to disentangle and selectively control various elements within style images, we develop a style-based classifier-free guidance (SCFG) inspired by text-based classifier-free guidance in diffusion models. Specifically, we employ a layout-controlled generation model, such as ControlNet, to produce a comparable “negative” image that lacks the target style we aim to transfer. This negative image functions similarly to a “null” prompt in diffusion models, allowing the guidance to focus exclusively on the target style element and filter out extraneous stylistic features. Finally, to enhance spatial stability, we incorporate a “teacher model” into the early stages of generation. The tutor model, based on the original text-to-image model, simultaneously performs denoising generation with the same textual prompt and shares its spatial attention maps with the style model at each time step. This method ensures stable and consistent spatial distribution, effectively mitigating issues like the checkerboard artifact. Additionally, this approach enables consistent spatial layouts across different style reference images for the same text prompt, facilitating more straightforward comparisons and evaluations of stylistic transformations.

In summary, our contributions are as follows:

- AdaIN-based Integration: We develop cross-modal Adaptive Instance Normalization (AdaIN) to harmoniously integrate style and text features, improving alignment during generation.
- Style-based Classifier-Free Guidance (SCFG): We introduce a style-guided approach to focus on the target style and reduce unwanted stylistic features.
- Teacher Model for Layout Stability: We incorporate a teacher model to share spatial attention maps during early generation, ensuring stable layouts and mitigating artifacts like the checkerboard effect.
- Extensive evaluations conducted on a wide range of styles and prompts, as shown in Fig. 1, demonstrate that our method significantly improves the alignment of the generated images. Furthermore, our approach is versatile and can be integrated into various existing style transfer frameworks while remaining fine-tuning-free.

## 2. Related Work

**Text-to-image generation.** Text-conditioned image generative models [1, 3, 5, 23, 26, 27, 30] have demonstrated remarkable capabilities in generating high-quality images.

Notably, models in the Stable Diffusion [6, 23, 27] series have achieved impressive advancements due to structural modifications and optimizations in their text encoder components. These improvements have led to significant enhancements in both the visual quality of generated images and the model’s ability to align with complex textual prompts. The robust performance of text-to-image (T2I) diffusion models has also catalyzed the development of various related visual generation tasks, including text-based image editing [10, 18, 20, 21, 31, 32], subject-driven image generation [4, 8, 15, 22, 29], and other applications that leverage their strong generative and interpretative abilities.

**Stylized Image Generation.** Style transfer applies the style of a reference image to a target content image and can be broadly categorized into image-driven and text-driven approaches. Image-driven methods, such as InST [40], StyleID [2], and InjectFusion [14], focus on preserving content while injecting style, employing techniques like stochastic inversion, query preservation, and feature blending in the h-space [17]. Text-driven methods aim to mitigate content leakage, where excessive style application distorts content features. B-LoRA [7] achieves style-content separation via LoRA weight optimization, while DEADiff [24] enhances text-image consistency through joint cross-attention learning. StyleAlign [36] and Visual Style Prompt [13] modify self-attention mechanisms by swapping either query-key (Q-K) or key-value (K-V) pairs, ensuring style alignment. Adapter-based methods provide an efficient alternative by injecting style features into pre-trained diffusion models. IP-Adapter [38] leverages a decoupled cross-attention mechanism for zero-shot style transfer and mitigates content leakage via weight tuning. InstantStyle [34] improves stylization by selectively injecting style features into Stable Diffusion’s UNet [23, 27, 28]. StyleShot [9] extracts multi-level style features for fine-grained control, while CSGO [37] enhances adapter-based stylization by training on a curated style dataset for effective style-content decoupling.

### 3. Method

In this section, we detail the three complementary strategies we propose to address the challenges inherent in text-driven style transfer. Each strategy builds upon existing models CSGO [37] but introduces novel mechanisms to overcome the inherent limitations of current approaches in style transfer tasks.

#### 3.1. Preliminaries

Before delving into our specific methods, we begin by providing some background on the key components that underlie our approach.

**Latent Diffusion Models.** Latent Diffusion Models (LDMs) [27] represent a powerful framework for efficient

image generation. Operating in the latent space of a Variational Auto-Encoder (VAE) [16], LDMs optimize computational efficiency while retaining high-quality image generation capabilities. Stable Diffusion (SD) [6, 23, 27], one of the most prominent models in this family, generates high-fidelity images by denoising noisy latent representations conditioned on text prompts. This denoising process can be formalized as:

$$\mathcal{L}(\theta) = \mathbb{E}_{t,z,c,\epsilon \sim \mathcal{N}(0,1)} \left[ \|\epsilon - \epsilon_\theta(z_t, t, c)\|_2^2 \right], \quad (1)$$

where  $\epsilon$  is noise sampled from a standard Gaussian distribution,  $\epsilon_\theta$  is the noise predicted by the U-Net architecture,  $z_t$  is the noisy latent representation at timestep  $t$ , and  $c$  is the conditioning text prompt. The goal of this loss function is to guide the model in iteratively reducing noise, which enables the model to generate high-quality images consistent with the input text.

**Attention Mechanisms.** To enhance the generation process, Stable Diffusion incorporates several self-attention and cross-attention layers [33]. Recent studies show that self-attention captures dependencies within the latent representations, allowing the model to effectively gather context across different parts of the feature map. Cross-attention facilitates the integration of conditioning text embeddings with latent representations, ensuring that the generated output aligns with the input text prompt. This dual attention mechanism enhances the expressiveness and contextual relevance of the output, achieving a high level of quality and stylistic consistency. The attention mechanism  $A(\cdot)$ , in its general form, can be defined as:

$$Q = W_Q(f), \quad K = W_K(f'), \quad V = W_V(f'), \\ A(Q, K, V) = MV, \text{ where } M = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right). \quad (2)$$

$f$  and  $d$  denote the input features and the dimensionality of head. For self-attention,  $f' = f$ , allowing the model to focus on intra-feature relationships. For cross-attention,  $f'$  corresponds to the conditioning input, such as text embeddings, allowing the latent feature map  $f$  to be modulated by the conditioning information. The attention map  $M$  determines the focus of the model during generation.

**Style Transfer with an Adapter.** Recent advancements in style transfer have explored the use of adapter-based methods to inject style-specific features into text-to-image models. These adapter networks are designed to modify pre-trained models to accommodate new styles inputs. One notable approach is the IP-Adapter [38], which leverages a dual cross-attention mechanism to integrate both text and image conditions. The CSGO [37] shares the same architecture for text-driven style transfer, with the key distinction being that it is specifically trained on a style transfer dataset, enabling it to better capture and transfer stylistic features. In

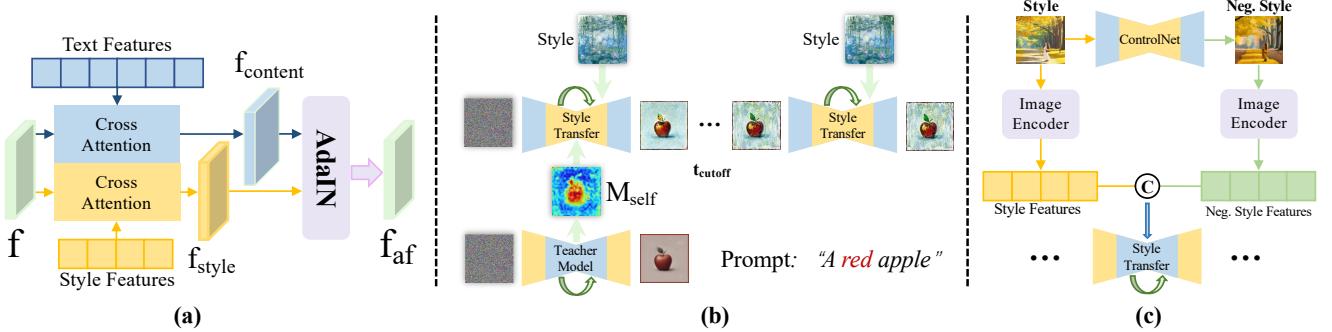


Figure 4. The illustration of our proposed Cross-Modal AdaIN, Teacher Model, Style-Based CFG.

this design, one cross-attention layer extracts features from the text prompt, while the other extracts features from the image. The two sets of features are then combined through a weighted sum, allowing the model to generate images that adhere to both the stylistic attributes specified by the text and the visual characteristics of the image. This combination process can be formalized as:

$$f_{ip} = A(Q, K_t, V_t) + \lambda A(Q, K_i, V_i), \quad (3)$$

where  $\lambda$  is a weight balancing the contributions from text-based and image-based attention. In this context,  $Q$  is derived from the feature map  $f$  of the latent representation,  $K_t$  and  $V_t$  are projections from the text embeddings  $f_t$ , and  $K_i$  and  $V_i$  are projections from the image feature map  $f_i$ .

While the IP-Adapter provides a straightforward and effective method for adapting models to perform style transfer, it has notable limitations. Specifically, the additive fusion of text and image conditions can lead to overfitting to the image style, particularly when there is conflicting information between the text and image conditions. In such cases, the model may overly prioritize the image style at the expense of the text prompt, resulting in outputs that are not aligned with the intended style or semantic content of the text. Furthermore, simply reducing the weight of the image condition can weaken the visual style, leading to suboptimal results. This introduces a challenge in determining an appropriate hyperparameter for balancing the contributions of both conditions.

In this paper, we introduce a new fusion strategy specifically tailored to style transfer tasks. Our method avoids the limitations of additive fusion by providing a more efficient and stable approach for combining text and image conditioning.

### 3.2. Cross-Modal AdaIN of Text and Style Conditioning

We propose a novel method for text-driven style transfer that better integrates both text and image conditioning. Our approach aims to achieve a balanced fusion of these two conditioning inputs, ensuring that they complement each other effectively.

In typical text-driven style transfer tasks, text conditioning primarily serves to define the content, dictating the semantic structure of the output, while image conditioning predominantly encodes the stylistic features, such as texture and color palette. However, directly combining these two conditioning inputs through weighted summation, as in existing methods [9, 34, 37, 38], forces them to assume similar roles in the fusion process. This can create conflicts, particularly when the text and image provide divergent information about the content and style. For example, a text prompt may describe a scene in one way, while the image may impose a stylistic choice that conflicts with this description, leading to suboptimal results. To address this challenge, we revisit Adaptive Instance Normalization (AdaIN) [12], a widely recognized technique in style transfer. AdaIN operates by normalizing the content input  $x$  based on the statistical properties (mean and standard deviation) of the style input  $y$ , integrating style characteristics while preserving the essential structure of the content. The AdaIN process is defined as follows:

$$\text{AdaIN}(x, y) = \sigma(y) \left( \frac{x - \mu(x)}{\sigma(x)} \right) + \mu(y), \quad (4)$$

where  $\mu(x)$  and  $\sigma(x)$  denote the mean and standard deviation of the content input  $x$ , and  $\mu(y)$  and  $\sigma(y)$  represent the mean and standard deviation of the style input  $y$ . By adjusting the content features to reflect the style statistics, AdaIN effectively fuses style into the content in a controlled manner, preserving the content's alignment with the text description while ensuring stylistic consistency.

As shown in Fig. 4(a), building on AdaIN, we develop a Cross-Modal AdaIN mechanism that integrates text and style conditioning in a way that respects their distinct roles. The Adapter-Based methods [37, 38] employed a weighted sum approach for feature fusion, ensuring that feature maps operated within the same embedding space. We discovered that AdaIN could directly replace the weighted sum strategy. This substitution enables effective feature fusion without the need for additional training, making it particularly advantageous when integrated with Adapter-Based methods [34, 37]. Specifically, we first leverage the cross-attention layers within the U-Net architecture to query both

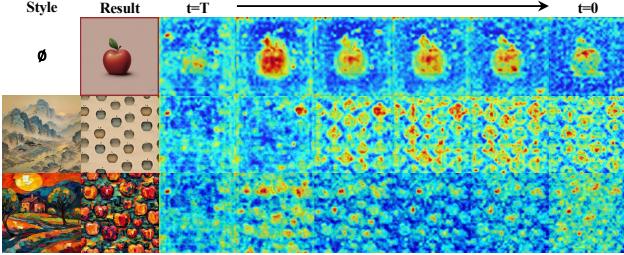


Figure 5. Visualization of the Cross-Attention Map for the word “apple” in the prompt “A red apple” during the generation process. When artifacts appear, the attention tends to scatter as well.

the style and text features based on the U-Net feature map. This step results in two separate grid feature maps—one for the style and one for the text—that share the same spatial resolution. Next, we apply AdaIN between two feature maps, where the feature maps queried by text conditions are normalized by style feature maps. Finally, the normalized text feature maps are fused into the raw U-Net features by simple addition. This fusion is constructed in a residual design, similar to the approach used in the raw cross-attention layers. Mathematically, this adaptive fusion process can be expressed as follows:

$$f_{\text{style}} = A(Q, K_{\text{style}}, V_{\text{style}}), f_{\text{text}} = A(Q, K_{\text{text}}, V_{\text{text}}), \hat{f}_{\text{af}} = \gamma_{\text{style}} \cdot \left( \frac{f_{\text{text}} - \mu_{\text{text}}}{\sigma_{\text{text}}} \right) + \beta_{\text{style}}, \quad (5)$$

where the  $\gamma_{\text{style}}$  and  $\beta_{\text{style}}$  come from the style image feature map, in the same way of run AdaIN. By adaptively balancing the influence of text and style, our method effectively minimizes potential conflicts between the two inputs, eliminating the need for setting a tricky hyperparameter.

### 3.3. Layout Stabilization with Teacher Model

In image generation, the layout is a crucial component of visual aesthetics. As shown in Fig. 3, we observe instances of artifacts during generation, such as checkerboard patterns. Upon analyzing the data presented in Fig. 5, we observed that these instabilities are correlated with a lack of aggregation in the core generative regions within the Cross-Attention mechanism. In the unstable generation examples, the layout instability reveal that the model fails to properly attend to regions associated with the word “apple” leading to visual distortions and compositional issues. This behavior diverges significantly from what is observed in the raw SDXL [23] model at different timesteps.

In image generation, Self-Attention plays a crucial role in maintaining the layout and spatial structure of the original content [18]. Self-Attention mechanism in Stable Diffusion [23, 27] captures high-level spatial relationships, which effectively stabilize the foundational layout during generation. The preserved layout information, encapsulated within the Self-Attention AttnMaps, serves as a structural frame-

work that guides the composition and distribution of elements across the image.

In the context of Text-Driven Style Transfer, maintaining a stable layout is crucial for ensuring that the generated image accurately reflects the structure and composition described by the textual prompt. Specifically, layout refers to the spatial arrangement of objects, elements, and the overall scene composition in the image. Without stable layout alignment, these elements may become misaligned, resulting in images where the content is distorted, or the intended focus, perspective, and balance outlined in the prompt are lost. To address this, we propose a method that stabilizes the layout by selectively replacing certain Self-Attention AttnMaps [23, 27, 33] in the stylized image with those from the original diffusion model. This selective replacement helps preserve the spatial relationships and arrangement of key features in the image, ensuring that the core layout remains consistent throughout the denoising process. By doing so, we retain the structural coherence of the original image while still applying the desired stylistic transformation, leading to a more coherent and faithful alignment with the textual prompt. In Fig. 4(b), provides a visual overview of this process, showcasing the integration of layout stabilization within the generative framework.

Unlike conventional image editing approaches [32] that replace Self-Attention maps (AttnMaps) across all timesteps and decoder layers, our method selectively replaces AttnMaps only during the initial denoising timesteps but across all layers of the UNet. Applying full-timestep replacement, as conventional methods do, risks excessive loss of stylistic details, limiting the effectiveness of style transfer.

### 3.4. Style-Based CFG

A particularly challenging scenario in style transfer arises when the reference style image contains multiple stylistic elements, such as a combination of cartoon style and nighttime aesthetics. In such cases, the model faces style ambiguity, where various style features are present, but the focus is intended to be on just one specific element. Current methods struggle to effectively disentangle these different styles and selectively emphasize the desired one. To address this challenge, a flexible method is required that can selectively emphasize the desired style elements while filtering out irrelevant or conflicting features. Inspired by the concept of classifier-free guidance (CFG) [11], commonly used in diffusion models for text-guided image generation, we propose a Style-Based CFG design to provide controlled adjustments in the style transfer process.

**Classifier-Free Guidance Mechanism.** In standard CFG, the model generates outputs conditioned on a given text prompt, as well as outputs generated without any conditioning (i.e., the unconditional model output). The final output

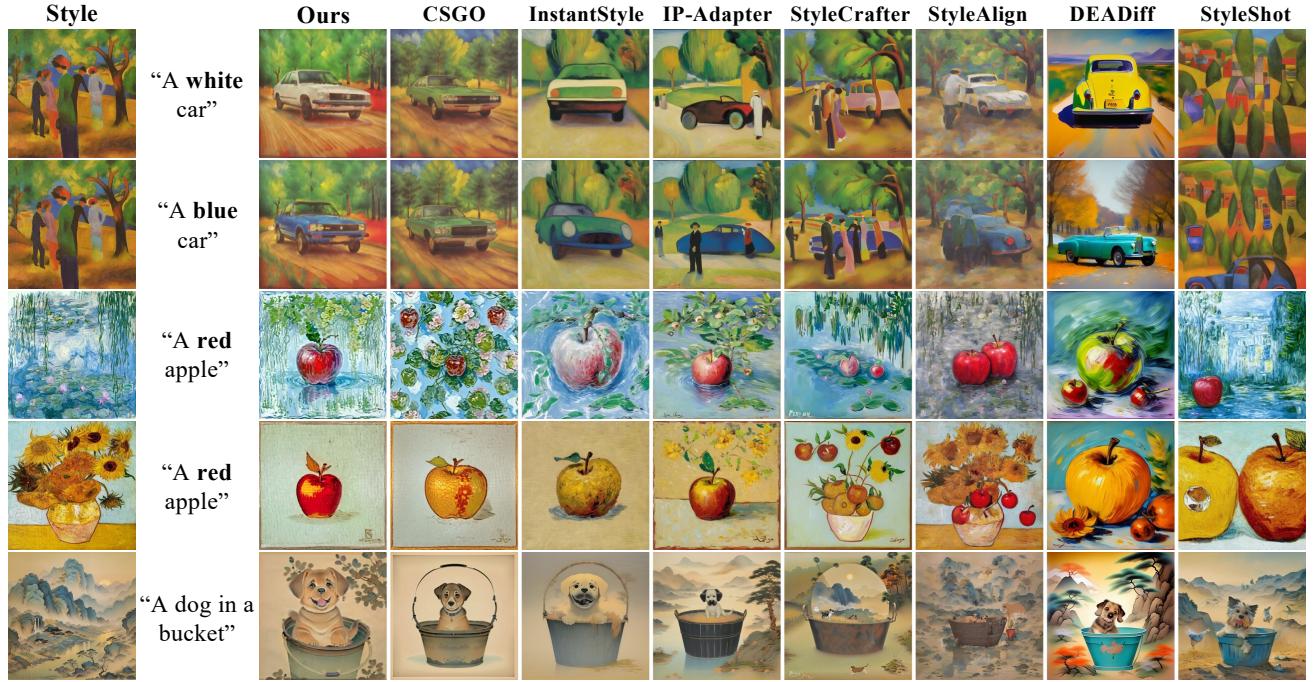


Figure 6. Qualitative comparison with state-of-the-art methods. Our approach effectively preserves image style while accurately adhering to text prompts for generation.

is a weighted combination of these two, where the conditional output steers the generation process in alignment with the prompt, while the unconditional output acts as a form of negative conditioning, helping to prevent the model from producing undesirable features. The unconditional output in the Style-Based CFG can be replaced with conditioning on a negative prompt, such as "blur" or "artifact" to actively discourage the generation of undesirable features. Therefore, the CFG [11] mechanism can be formalized as:

$$\hat{\epsilon}_\theta(z_t, t, y) = (1 + w) \cdot \epsilon_\theta(z_t, t, y_{\text{cond}}) - w \cdot \epsilon_\theta(z_t, t, y_{\text{neg}}), \quad (6)$$

where  $y_{\text{cond}}$  is the positive condition and  $y_{\text{neg}}$  is the negative condition, and  $w$  is a weight controlling the balance between these outputs.

**Style-Based Classifier-Free Guidance.** In the context of style transfer, we extend this CFG mechanism to address the challenge of style ambiguity in images. Specifically, we introduce the concept of a negative style image that retains the overall content of the reference image but excludes the target style element. This negative image serves as a counterpart to the original style image and helps the model focus on transferring only the desired style component. To generate the negative style image, we use a layout-controlled generation model, such as ControlNet, which allows us to create an image  $z_t^{\text{neg}}$  that preserves the structural features of the reference image but omits the target style. This negative image functions similarly to a negative prompt in text-guided CFG, effectively guiding the model to emphasize the de-

sired style while suppressing undesired style elements. As shown in Fig. 4(c), the SCFG mechanism operates as follows: 1) Generate a Negative Image: Using ControlNet, generate a negative sample image  $z_t^{\text{neg}}$  that retains the structural elements of the image while omitting the target style. 2) Apply SCFG to Guide Generation: Formulate the Style-Based CFG by modulating the balance between the target style image  $z_t^{\text{style}}$  and the negative style image  $z_t^{\text{neg}}$ . This balance is governed by a weight factor  $w$ , which determines the contribution of each image during generation. The noise prediction is modified as follows:

$$\begin{aligned} \hat{\epsilon}_\theta(z_t, t, y) = & (1 + w) \cdot \epsilon_\theta(z_t, y_{\text{cond}}, y_{\text{cond}}^{\text{style}}) \\ & - w \cdot \epsilon_\theta(z_t, y_{\text{neg}}^{\text{text}}, y_{\text{neg}}^{\text{style}}), \end{aligned} \quad (7)$$

By integrating SCFG, our method refines the generation process by isolating specific style components, filtering out extraneous ones, and thereby focusing style transfer on the desired features. This approach reduces the risk of overfitting to irrelevant style components, allowing the model to perform effective style transfer in complex scenarios with multiple stylistic elements.

## 4. Evaluation and Experiments

**Implementation details.** We have implemented our method on top of the latest Adapter-Based Style Transfer approach, named CSGO [37]. To ensure fairness in comparison and mitigate the influence of random initialization,

Metric	SDXL-based Methods					SD15-based Methods		<b>Ours</b>
	IP-Adapter [38]	InstantStyle [34]	CSGO [37]	StyleAlign [36]	StyleCrafter [19]	StyleShot [9]	DEADiff [24]	
Text Alignment $\uparrow$	0.221	0.229	0.216	0.180	0.189	0.202	0.229	<b>0.235</b>
infer Time (s)	6	6	9	48	4	3	2	17
User-study Text %	7.48	6.46	7.99	5.78	3.06	2.55	1.87	<b>62.92</b>
User-study Style %	6.63	8.67	6.97	7.82	8.67	5.10	5.27	<b>50.85</b>

Table 1. Quantitative comparison with state-of-the-art methods. Our approach achieves the best performance on the text alignment metric and outperforms others in the user study evaluation.

Cross-Modal AdaIN	Teacher Model	Text Alignment $\uparrow$
		0.216
✓	✓	0.223 (+3.2%)
✓	✓	0.228 (+5.5%)
		0.235 (+8.7%)

Table 2. Ablation study evaluating the impact of our proposed methods. Both designs significantly enhance text alignment accuracy.

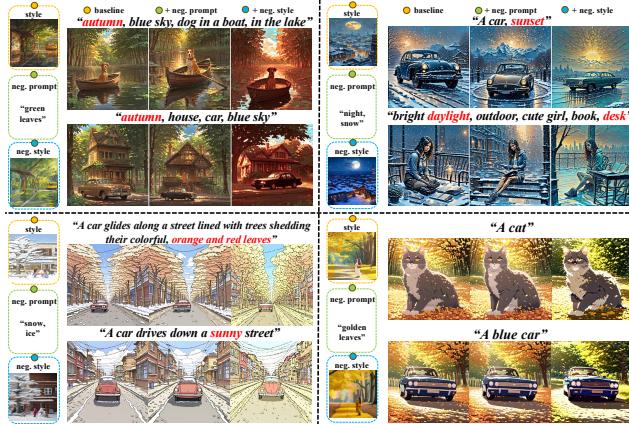


Figure 7. Impact of Style-Based CFG. The proposed Style-Based CFG successfully eliminates unintended style elements such as “green leaves”, “night”, “snow” and “golden leaves”, while text-based CFG fails to address these spurious attributes.

we fixed the initial noise for all methods, as the initial noise can significantly impact the final results [10].

**Evaluation.** To enable a comparison with existing methods, we evaluated both our approach and several previous ones, including CSGO [37], InstantStyle [34], IP-Adapter [38] with weight tuning, StyleCrafter [19], StyleAlign [36], DEADiff [24], StyleShot [9]. To evaluate the performance of these methods in terms of prompt adherence after text-driven style transfer, we constructed a benchmark consisting of 52 prompts and 20 style reference images. These prompts were selected from the settings used in StyleAdapter [36], while the style images are taken from StyleShot [9]. Additionally, using ChatGPT, we generated 30 prompts in the form of “A <color> <object>” which better highlights the issue of style overfitting.

More results of our method can be found in Appendix E, and its integration with other methods is detailed in Appendix F.

## 4.1. Comparison with State-of-the-Arts

**Qualitative Comparisons.** Fig. 6 presents qualitative comparisons with state-of-the-art methods. Existing approaches struggle with prompt alignment due to style overfitting. While CSGO [37], InstantStyle [34], and DEADiff [24] mitigate content leakage, they fail to fully capture prompt-specified details. Conversely, IP-Adapter [38], StyleShot [9] and StyleCrafter [19] exhibit content leakage, leading to images misaligned with the input prompts. In contrast, our method ensures better prompt alignment while maintaining layout stability. More detailed results and explanations could be found in Appendix D.

**Quantitative Comparisons.** To verify the alignment between the generated image and its specified object, we compute the CLIP cosine similarity [25] between the image and the corresponding text description. As shown in Tab. 1, Our method outperforms the others, achieving the highest text alignment capability.

**User Study.** We also conducted a user study to assess users’ evaluations of text alignment and style similarity, with the results presented in Tab. 1. For each method’s generated images, 49 users participated in an anonymous vote, selecting the example they felt was the most aligned with the text description and the closest in style to the reference image. The normalized votes (vote rate) serve as the scores for text alignment and style similarity.

## 4.2. Style-Based CFG

We conducted experiments on Style-Based CFG (SCFG), with the results shown in Fig. 7. In the baseline images, unintended style elements such as snow and golden leaves are present, which do not align with the intended style. Adding a negative text prompt allowed for some control over these elements; however, unintended style elements like snow were not effectively mitigated. By applying a negative style image, we achieved more effective control, successfully removing these unwanted elements. This demonstrates the effectiveness of SCFG in precisely managing unintended style elements in generated images.

## 4.3. Ablation Study

**Cross-Modal AdaIN.** To demonstrate the effectiveness of each component in our method, we conducted an ablation study focusing on a quantitative analysis of text alignment. The same dataset used in the quantitative evaluation was employed for this test. As shown in Tab. 2,

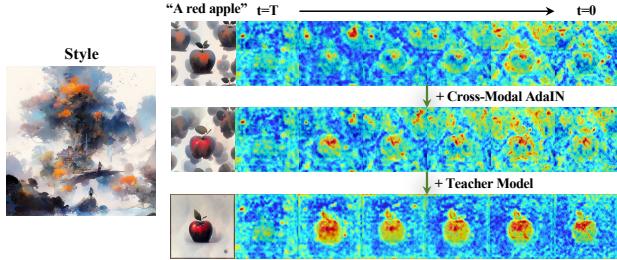


Figure 8. Visualization of cross-attention maps for the word “apple” in the prompt “A red apple” across different models. The proposed teacher model effectively rectifies attention maps, leading to improved image generation quality.

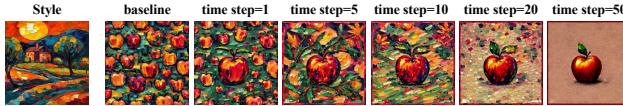


Figure 9. Impact of Teacher Model on Style Image Generation. The term “timestep” refers to the number of denoising steps during which the teacher model is involved.

our method consistently improves text alignment accuracy over the baseline. Incorporating the Teacher Model alone provides an initial enhancement, while cross-modal AdaIN yields a more significant improvement. Finally, combining both the cross-modal AdaIN and the Teacher Model results in the highest level of improvement, indicating complementary effects that enhance alignment performance.

**Teacher Model.** In Fig. 8, shows a visualization of cross-attention maps for the word “apple” in the prompt “A red apple” across different model configurations. The first row represents the baseline model, where the attention map fails to effectively focus on the core concept of “apple” from the text prompt. In the second row, the addition of cross-modal AdaIN provides some improvement, but the attention is still diffuse and lacks precise focus. In the third row, with the Teacher Model added, the attention map becomes more concentrated on the target concept, demonstrating a clear improvement in alignment with the prompt. This comparison indicates that it is the Teacher Model that enables the model to better capture and emphasize key elements from the text prompt, resulting in improved attention quality and image generation accuracy. In Fig. 9, this experiment demonstrates the importance of selecting an appropriate denoising timestep to stop the involvement of the Teacher Model. If the Teacher Model’s influence is removed too early, as seen at lower timesteps, issues with layout stability persist, resulting in compositions that lack coherence and include multiple instances of the target object (“apple”). However, if the Teacher Model is involved throughout the full denoising process (e.g., at timestep 50), there is a noticeable loss of style in the generated image. This suggests that while the Teacher Model is crucial for achieving layout stability, an extended involvement can dilute the style elements, making it essential to find a balanced point to discontinue its

influence for optimal results. More detailed results and explanations can be found in Appendix C.



Figure 10. More results of our text-driven style transfer model. Illustration of the prompt format used: “A [color] bus”.

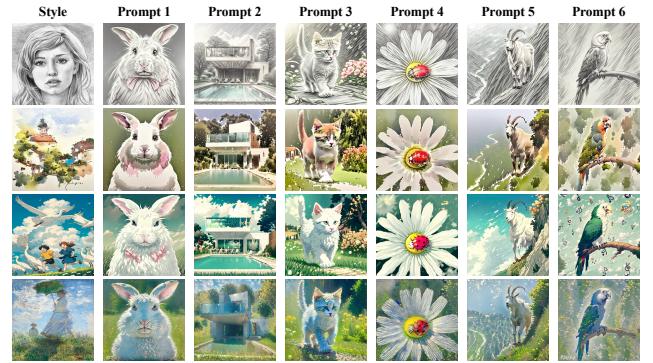


Figure 11. The prompts, from left to right, are: “A fluffy white rabbit with pink ears and nose”, “A modern house with a pool.”, “A lovely kitten walking in a garden.”, “A daisy with a ladybug on it”, “A mountain goat on a cliff”, and “A parrot singing a song”.

**More Results.** Fig. 10 demonstrates that our method effectively mitigates style overfitting by consistently adapting the target color while preserving structural details. Similarly, Fig. 11 shows that our model can handle complex prompts without losing style fidelity or scene composition. More detailed results can be found in Appendix E.

## 5. Conclusion

In conclusion, existing text-driven style transfer faces key issues such as style overfitting and layout instability, which limit the adaptability and coherence of generated images. To address these challenges, we proposed three methods: cross-modal AdaIN for harmonizing style and text features, style-based classifier-free guidance (SCFG) for selective control of stylistic elements, and a Teacher Model to enhance layout stability. Our results confirm that our approach effectively mitigates these issues, improving alignment, stability, and control in style transfer, making it a versatile and robust solution for text-to-image synthesis tasks.

**Limitations and Future Work.** While our method improves artifact removal and layout stability, the Teacher Model slightly increases inference time. Additionally, generating negative-style image requires expertise and manual effort. Future work could focus on improving efficiency and further exploring strategies to mitigate style overfitting, enabling more adaptive and generalizable style transfer across diverse prompts and visual domains.

## References

- [1] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023. 2
- [2] 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*, pages 8795–8805, 2024. 3
- [3] Xiaoliang Dai, Ji Hou, Chih-Yao Ma, Sam Tsai, Jiliang Wang, Rui Wang, Peizhao Zhang, Simon Vandenhende, Xiaofang Wang, Abhimanyu Dubey, et al. Emu: Enhancing image generation models using photogenic needles in a haystack. *arXiv preprint arXiv:2309.15807*, 2023. 2
- [4] Ganggui Ding, Canyu Zhao, Wen Wang, Zhen Yang, Zide Liu, Hao Chen, and Chunhua Shen. Freecustom: Tuning-free customized image generation for multi-concept composition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9089–9098, 2024. 3
- [5] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. *Advances in neural information processing systems*, 34:19822–19835, 2021. 2
- [6] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024. 1, 3
- [7] Yarden Frenkel, Yael Vinker, Ariel Shamir, and Daniel Cohen-Or. Implicit style-content separation using b-lora. In *European Conference on Computer Vision*, pages 181–198. Springer, 2025. 3
- [8] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022. 3
- [9] Junyao Gao, Yanchen Liu, Yanan Sun, Yinhao Tang, Yanhong Zeng, Kai Chen, and Cairong Zhao. Styleshot: A snapshot on any style. *arXiv preprint arXiv:2407.01414*, 2024. 3, 4, 7, 1, 2
- [10] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. 3, 7
- [11] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 5, 6
- [12] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision*, pages 1501–1510, 2017. 4
- [13] Jaeseok Jeong, Junho Kim, Yunjey Choi, Gayoung Lee, and Youngjung Uh. Visual style prompting with swapping self-attention. *arXiv preprint arXiv:2402.12974*, 2024. 3
- [14] Jaeseok Jeong, Mingi Kwon, and Youngjung Uh. Training-free content injection using h-space in diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5151–5161, 2024. 3
- [15] Chen Jin, Ryutaro Tanno, Amrutha Saseendran, Tom Diethe, and Philip Teare. An image is worth multiple words: Learning object level concepts using multi-concept prompt learning. *arXiv preprint arXiv:2310.12274*, 2023. 3
- [16] Diederik P Kingma. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 3
- [17] Mingi Kwon, Jaeseok Jeong, and Youngjung Uh. Diffusion models already have a semantic latent space. *arXiv preprint arXiv:2210.10960*, 2022. 3
- [18] Bingyan Liu, Chengyu Wang, Tingfeng Cao, Kui Jia, and Jun Huang. Towards understanding cross and self-attention in stable diffusion for text-guided image editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7817–7826, 2024. 3, 5
- [19] Gongye Liu, Menghan Xia, Yong Zhang, Haoxin Chen, Jinbo Xing, Yibo Wang, Xintao Wang, Yujiu Yang, and Ying Shan. Stylecrafter: Enhancing stylized text-to-video generation with style adapter. *arXiv preprint arXiv:2312.00330*, 2023. 7, 1, 2, 3, 5
- [20] 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. *arXiv preprint arXiv:2108.01073*, 2021. 3
- [21] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6038–6047, 2023. 3
- [22] Jisu Nam, Heesu Kim, DongJae Lee, Siyoon Jin, Seungryong Kim, and Seunggyu Chang. Dreammatcher: Appearance matching self-attention for semantically-consistent text-to-image personalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8100–8110, 2024. 3
- [23] 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. *arXiv preprint arXiv:2307.01952*, 2023. 1, 2, 3, 5
- [24] 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8693–8702, 2024. 3, 7
- [25] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 7

- [26] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022. 2
- [27] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 2, 3, 5
- [28] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5–9, 2015, proceedings, part III* 18, pages 234–241. Springer, 2015. 3
- [29] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22500–22510, 2023. 3
- [30] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022. 2
- [31] Xue Song, Jiequan Cui, Hanwang Zhang, Jingjing Chen, Richang Hong, and Yu-Gang Jiang. Doubly abductive counterfactual inference for text-based image editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9162–9171, 2024. 3
- [32] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1921–1930, 2023. 3, 5
- [33] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. 3, 5, 2
- [34] Haofan Wang, Matteo Spinelli, 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. 3, 4, 7, 1, 2, 5
- [35] Zhouxia Wang, Xintao Wang, Liangbin Xie, Zhongang Qi, Ying Shan, Wenping Wang, and Ping Luo. Styleadapter: A single-pass lora-free model for stylized image generation. *arXiv preprint arXiv:2309.01770*, 2023. 1, 2
- [36] Zongze Wu, Yotam Nitzan, Eli Shechtman, and Dani Lischinski. Stylealign: Analysis and applications of aligned stylegan models. *arXiv preprint arXiv:2110.11323*, 2021. 3, 7, 2
- [37] Peng Xing, Haofan Wang, Yanpeng Sun, Qixun Wang, Xu Bai, Hao Ai, Renyuan Huang, and Zechao Li. Csgo: Content-style composition in text-to-image generation. *arXiv preprint arXiv:2408.16766*, 2024. 2, 3, 4, 6, 7, 1
- [38] Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. Ipadapter: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023. 3, 4, 7, 1, 2
- [39] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 1
- [40] Yuxin Zhang, Nisha Huang, Fan Tang, Haibin Huang, Chongyang Ma, Weiming Dong, and Changsheng Xu. Inversion-based style transfer with diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10146–10156, 2023. 3

This **Appendix** provide additional details regarding the experimental setup described in the main paper and offer an extended analysis of the contributions of individual components. The content is organized as follows:

- **Details of Experiments.** This section provides additional information about the experiments discussed in the main paper, including specifics on the quantitative evaluations and the user study setup.
- **Ablation Study.** Qualitative comparisons from the ablation experiments are presented, analyzing the impact of the Teacher Model, particularly in terms of timestep selection and the choice of Attention Map.
- **Additional Qualitative Comparisons.** This section presents extensive qualitative comparisons, demonstrating that cross-modal AdaIN effectively prevents style overfitting, while the Teacher Model ensures layout stability and mitigates the occurrence of artifacts.
- **Integration with Other Methods.** This section explores how our approach can be integrated with existing methods, such as InstantStyle [34] and StyleCrafter [19], showcasing its ability to enhance their performance and adaptability.

## A. Implementation Details

We set the random seed to 42 for reproducibility, used 50 inference steps, and applied a uniform guidance scale of 5 across all methods. In the qualitative and quantitative comparison experiments, for the implementation involving the Teacher Model, its participation was specifically limited to the first 20 steps. All experiments were conducted on a single NVIDIA GTX-4090 GPU.

Adapter-based methods [9, 34, 37, 38] are particularly suitable for style transfer. Their fine-tuning-free nature, combined with high-quality style transfer performance, has made them widely adopted. CSGO [37] employs a widely used adapter-based model structure and is the first method trained on a meticulously curated dataset specifically designed for style transfer. This effectively decouples the content and style in style images, enhancing the grasp of style details such as brushstrokes and textures. Therefore, in the experimental section, we selected it as the baseline and implemented specific modifications based on it. The implementation details are as follows:

We only retained the modules in CSGO [37] related to text-driven style transfer, removing irrelevant components, *e.g.*, ControlNet [39]. This optimization reduces potential interference while lowering experimental costs, including memory usage during inference. At the same time, both the Teacher Model and cross-modal AdaIN are optional and can be used based on specific needs. For the quantitative experiments in the main paper, we incorporated both the Teacher Model and the cross-modal AdaIN module to achieve optimal text alignment. In the qualitative and quantitative com-

parison experiments, the Teacher Model participated for the first 20 time steps, with the total number of inference steps set to 50.

---

### Algorithm 1 SDXL-Guided Self-Attention Replacement

---

**Input:**  $P_{\text{dst}}$ : a target prompt;  $I_{\text{ref}}$ : style reference image;  $S$ : random seed; DM: raw Stable Diffusion Model; ST: style transfer Method Model;  $t_{\text{cutoff}}$ : stop replacement time step;

**Output:**  $I_{\text{style}}$ : text-driven stylized image;

```

1:  $z_T \sim \mathcal{N}(0, 1)$ , a unit Gaussian random value sampled with
   random seed  $S$ ;
2:  $z_T^* \leftarrow z_T$ ;
3: for  $t = T, T - 1, \dots, 1$  do
4:   if  $t > t_{\text{cutoff}}$  then
5:      $z_{t-1}, M_{\text{self}} \leftarrow \text{DM}(z_t, P_{\text{dst}}, t)$ ;
6:      $z_{t-1}^* \leftarrow \text{ST}(z_t^*, I_{\text{ref}}, P_{\text{dst}}, t) \{ M_{\text{self}}^* \leftarrow M_{\text{self}} \}$ ;
7:   else
8:      $z_{t-1}^* \leftarrow \text{ST}(z_t^*, I_{\text{ref}}, P_{\text{dst}}, t)$ 
9:   end if
10:  end for
11:  Return  $I_{\text{res}} \leftarrow \text{Decoder}(z_0)$ ;
```

---

## B. Evaluation Settings and User Study

In the quantitative experiments presented in the main paper, the evaluation was conducted using prompts derived from StyleAdapter [35], with specific examples provided in Fig. 12. The style images were randomly sampled from the test set of StyleShot [9], with representative examples shown in Fig. 13. Ultimately, each method generated 1,000 images for the quantitative experiments.

Beyond quantitative evaluations, we conducted a user study to gain subjective insights into the performance of different methods. The study involved 12 pairs of reference images and prompts. For each pair, participants were asked to assess and select the method they found superior based on two criteria: text alignment and style similarity. To ensure a fair assessment, participants were provided with a brief explanation of the task and evaluation criteria beforehand. We collected responses from 49 participants with diverse backgrounds, including individuals with relevant expertise in text-to-image tasks. The specific design of the questionnaire, including example pairs and evaluation guidelines, is shown in Fig. 16.

## C. Additional Ablation Study

**Qualitative Results of the Ablation Study.** While the main paper presents a quantitative analysis, a qualitative comparison provides a more intuitive understanding of the contributions of each component. By incrementally integrating the corresponding components, we demonstrate their individual effects. Fig. 17 showcases representative visual outcomes from our qualitative experiments. A com-

A robot.	A white rose.
A girl wearing a red dress, she is dancing.	A sunflower smiling at the sun.
A boy wearing glasses, he is reading a thick book.	A cactus wearing a hat.
A little cute boy.	A daisy with a ladybug on it.
A woman wearing a green sportswear, she is running.	A pine tree with a snowman hugging it.
A woman wearing a purple hat and a yellow scarf.	A mushroom in winter.
A man wearing a black leather jacket and a red tie.	A lotus with a frog meditating on it.
A little boy with glasses and a watch.	A cherry blossom.
A smiling little girl.	A palm tree.
A curly-haired boy.	A river with rapids and rocks.
A little girl holding flowers.	A lake with calm water and colorful pebbles.
A lovely kitten walking in a garden.	A waterfall with mist and rainbows.
A puppy sitting on a sofa.	A stone with a face carved on it, standing on a pedestal in a museum.
A fluffy white rabbit with pink ears and nose.	A stone with a hole in it.
A brown puppy with black spots and a red collar.	A stone with a pattern of stripes on it.
A black and white panda.	A snowy mountain peak.
A dog in a bucket.	A mountain goat on a cliff.
A cat wearing a hat.	A red baseball cap.
A cute little fish in an aquarium.	A football on the grass.
A bird in a word.	A motorcycle.
A kitten sleeping on a pillow.	A modern house with a pool.
A parrot singing a song.	A house made of cardboard boxes.
A monkey playing with a banana.	A house covered with ice and snow.
A turtle wearing sunglasses.	
A hamster eating a carrot.	

Figure 12. Details of the Test Set. The prompts used in the quantitative experiments were derived from StyleAdapter [35].

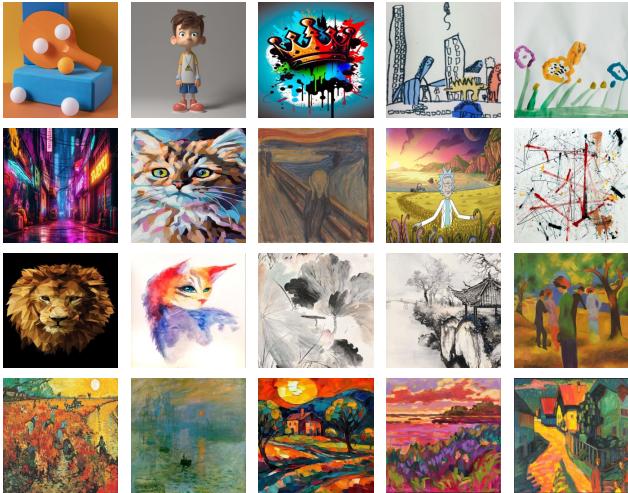


Figure 13. Details of the Test Set. The style images used in the quantitative experiments were randomly sampled from the test set of StyleShot [9].

parison between the second and third columns highlights that cross-modal AdaIN significantly improves text alignment while preserving style similarity. Furthermore, as shown in the green apple example, introducing the Teacher Model not only enhances layout stability but also resolves remaining artifacts, ensuring spatial consistency across different styles.

**Self-Attention map and layout stability.** In the UNet of Stable Diffusion [23, 27], Cross-Attention [33] primarily aligns the prompt with the generated image, determining how textual input influences the overall style and content. Self-Attention [33], on the other hand, focuses on the internal coherence of the image, maintaining spatial relationships and structural consistency. As shown in Fig. 18, swapping the Self-Attention Map ensures layout stability and consistency across different styles of images, whereas replacing the Cross-Attention Map fails to achieve this effect, resulting in noticeable differences in the main layout under varying styles. All experiments were conducted by adding the Teacher Model to the baseline CSGO frame-

work. To objectively evaluate the impact of the Teacher Model, cross-modal AdaIN was not used in these experiments, isolating the Teacher Model’s contribution to layout stability.

**Choice of Teacher Model participation timestep.** To evaluate the impact of the Teacher Model’s participation timestep on the final generation results, we conducted experiments analyzing its effect. The Teacher Model is designed to ensure layout stability while avoiding artifacts, such as checkerboard patterns. To objectively evaluate the impact of the Teacher Model, cross-modal AdaIN was not used in these experiments. As shown in Fig. 19, the term “timestep” refers to the number of denoising steps during which the Teacher Model is active. The results demonstrate that insufficient participation (short timesteps) fails to resolve layout issues, while prolonged involvement (long timesteps) negatively affects the final style fidelity. Rows 3 and 4 illustrate that even small changes in the timestep significantly influence the results, while Rows 5 and 6 show that the optimal timestep can vary across different styles. Based on these findings, a timestep between 10 and 20 strikes a reasonable balance between layout stability and style preservation.

**Compare with image-based style transfer(I2I).** Although our method utilizes the Self-Attention Map provided by the Teacher Model, this does not equate to I2I. As shown in Fig. 20, the I2I approach provided by CSGO [37] fails to preserve the color information of the content image effectively. In contrast, our method can more accurately adhere to the prompt’s description. To ensure fairness, the noise used in our method is identical to that used in generating the content image.

## D. Additional Comparisons

Qualitative experiments are conducted to visually demonstrate the strengths of our method, particularly in capturing style details and ensuring alignment with the given textual descriptions. This allows for a more intuitive comparison with state-of-the-art methods, showcasing the superior performance of our approach in real-world scenarios. We provided additional qualitative comparisons between our method and state-of-the-art approaches to better illustrate the strengths and weaknesses of each method.

In Fig. 24, our method outperforms others in both overall style similarity and the ability to capture fine details, such as textures. Additionally, it achieves the highest accuracy in aligning with the prompt descriptions. For methods based on the Stable Diffusion XL [23], approaches like CSGO [37] and InstantStyle [34] exhibit noticeable style overfitting, while IP-Adapter [38] and StyleCrafter [19] tend to suffer from content leakage. Meanwhile, StyleAlign [36] produces results of relatively lower quality. For methods based on the Stable Diffusion 1.5,

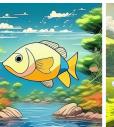
Style	Ours	CSGO	InstantStyle	IP-Adapter	StyleCrafter	StyleAlign	DEADiff	StyleShot
								
<b>"A snowy mountain peak."</b>								
<b>clip style similarity</b>	0.642	0.661	0.639	0.650	<b>0.732</b>	0.720	0.606	0.659
<b>dino style similarity</b>	0.467	0.289	0.395	0.350	0.721	<b>0.723</b>	0.315	0.525
								
<b>"A cute little fish in aquarium."</b>								
<b>clip style similarity</b>	0.530	0.629	0.672	0.829	0.802	<b>0.867</b>	0.667	0.691
<b>dino style similarity</b>	0.184	0.216	0.194	0.433	0.584	<b>0.711</b>	0.311	0.396

Figure 14. We observed that existing metrics generally fail to capture adherence to style. They tend to favor higher semantic similarity to the style image rather than better style transfer, a known issue often referred to as content leakage. A higher semantic similarity score does not indicate better style preservation and can, in fact, weaken the style in the generated results.



Figure 15. More results of Style-Based CFG.

DEADiff [24] struggles with accurately capturing the style, and although StyleShot [9] performs reasonably well in capturing style, it still encounters issues such as content leakage. Content leakage can indeed be seen as a form of overfitting to the style reference, where the model overly relies on the style image, causing elements of the style reference to dominate or intrude on the content representation. This highlights a lack of proper disentanglement between style and content in such cases.

A more nuanced form of style overfitting, as discussed in this paper, arises when text-driven style transfer methods struggle to adapt to nuanced variations in prompt details, such as changes in color. The challenge lies in whether these methods can accurately align with the evolving prompt descriptions while preserving the integrity of the style. This aspect is further validated in Fig. 25. Methods such as CSGO [37], InstantStyle [34], and StyleShot [9] struggle to differentiate the color specifications described in the prompt. Additionally, IP-Adapter [38] and DEADiff [24] face challenges with style dissimilarity, while Style-

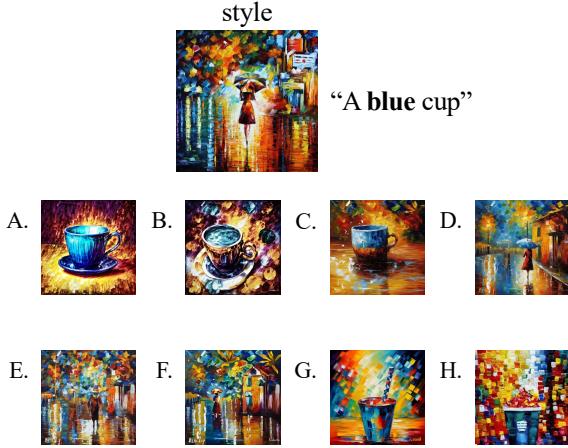
Crafter [19] demonstrates some bias toward the structure of the style reference, particularly evident in the “car” example.

In the main paper, we also focus on the issue of layout stability. Through extensive experiments, as shown in Fig. 26, we demonstrate that our method can effectively ensure layout stability. CSGO [37] frequently exhibits artifacts such as checkerboard patterns, while other methods also encounter issues with layout instability. Notably, content leakage appears to be closely related to layout disruptions. This can be validated from the experimental results of StyleCrafter [19] and IP-Adapter [38]. Although “A red apple” is reflected in the final generated output, the image contains too many unrelated elements from the style reference, making it appear overly cluttered.

## E. More results from our study

In Fig. 27, Fig. 28, and Fig. 15 we provide additional visualization results showcasing the effectiveness and versa-

Choose the image that best matches the **text** description from the 8 images below.



Choose the image that best matches the **style** from the 8 images below.

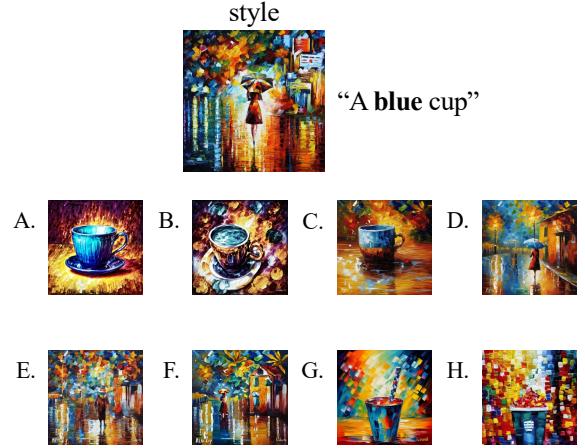


Figure 16. The questionnaire format for the user study. Each option represents the generation result of a method under a given style and prompt.

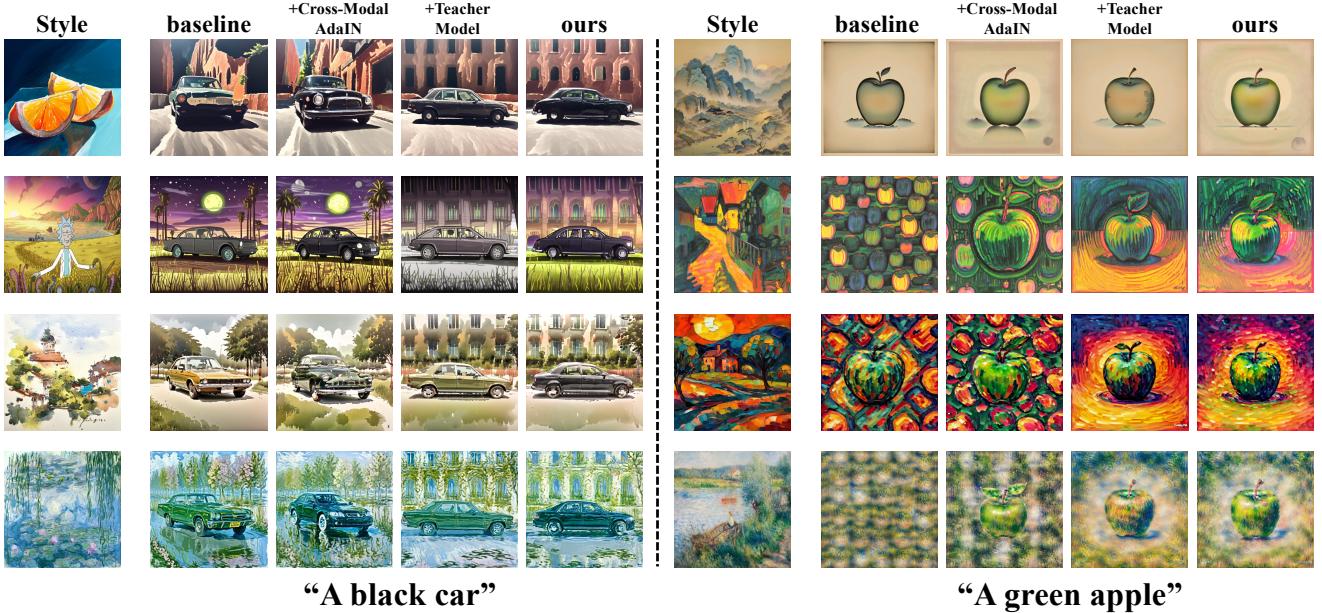


Figure 17. Qualitative results of the ablation study. cross-modal AdaIN enhances text alignment while preserving style similarity, addressing style overfitting issues. Incorporating the Teacher Model improves layout stability and resolves artifacts, ensuring consistent layout arrangements across different styles, as demonstrated in the “A green apple” example.

tility of our method. We have selected a variety of style categories and different color schemes to highlight the alignment effects for text descriptions. Moreover, we achieve excellent layout stability even when using the same prompt.

## F. Integration with Other Methods

CSGO [37] has been recognized as one of the most effective and state-of-the-art methods for style transfer, which is why it was selected as the primary baseline in the main paper.

To further evaluate the generalizability and robustness of our approach, we additionally explored its application and performance on other models.

### F.1. Integration with InstantStyle [34]

**Cross-Modal AdaIN.** Since InstantStyle [34] is also an adapter-based architecture, it can similarly integrate cross-modal AdaIN to mitigate style overfitting. The results are shown in Fig. 21. Compared to Row 1, Row 2 accurately follows the text description, effectively avoiding errors in

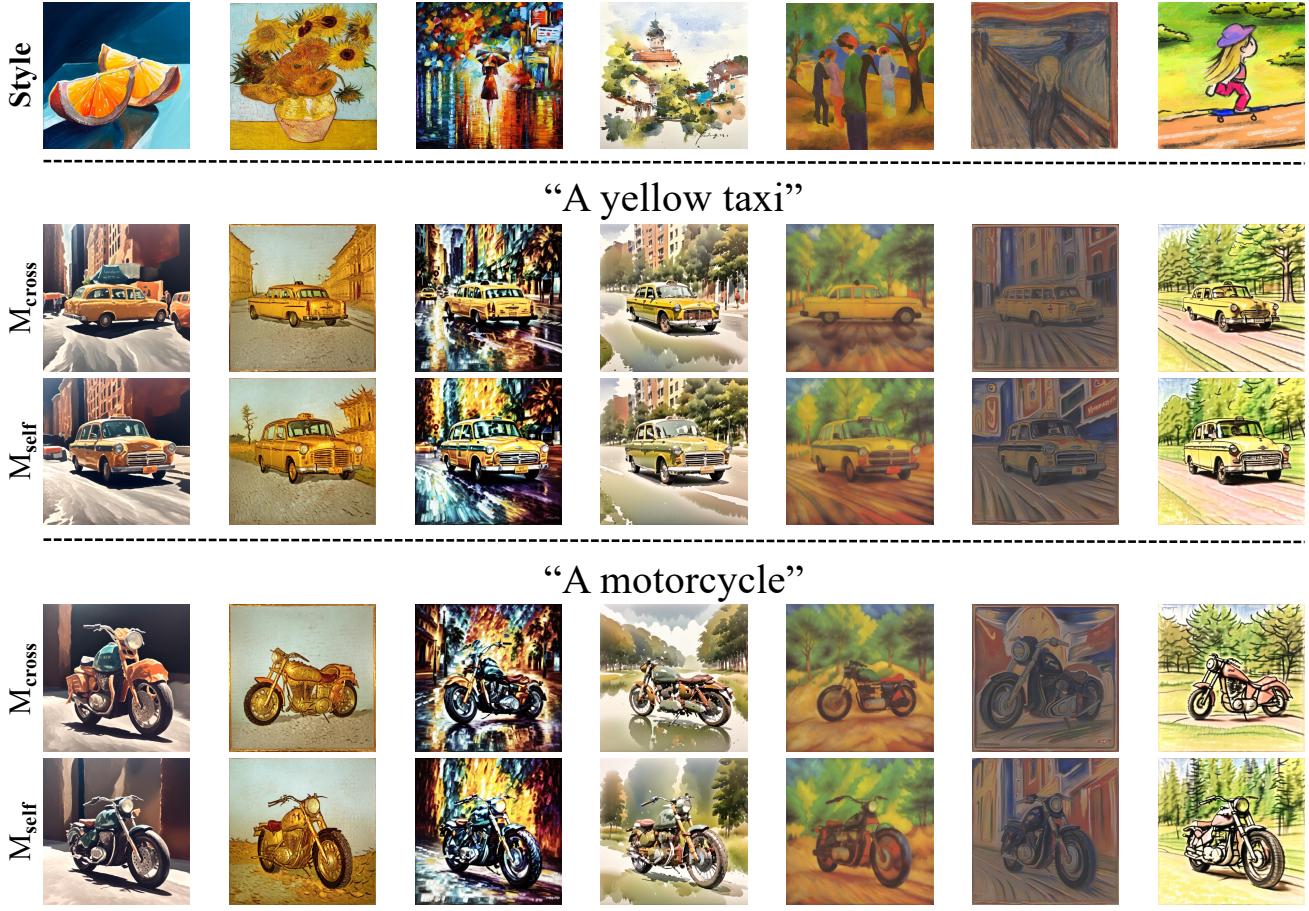


Figure 18. Implementation of the Teacher Model: Comparison of substituting the Self-Attention Map and Cross-Attention Map. The results demonstrate that replacing the Self-Attention Map achieves layout stability and consistency across different styles of images.

the generated output.

**Teacher Model.** InstantStyle [34] also encounters artifacts such as checkerboard patterns. Similar to the previous approach, we investigated the impact of the Teacher Model’s involvement at different timesteps on the results, as shown in Fig. 22. Upon observation, we reached a similar conclusion: if the Teacher Model participates for too many timesteps, it can lead to style loss.

## F.2. Integration with StyleCrafter [19]

**Teacher Model.** A notable issue in StyleCrafter [19] is content leakage, where unrelated content elements from the style image appear in the generated results, ultimately affecting the final output. This phenomenon can lead to generated images that do not align with the descriptions in the prompt. To address this, we incorporated the Teacher Model into the method. As shown in Fig. 23, the inclusion of the Teacher Model significantly mitigates the problem of content leakage, resulting in outputs that maintain stability and consistency across different styles.

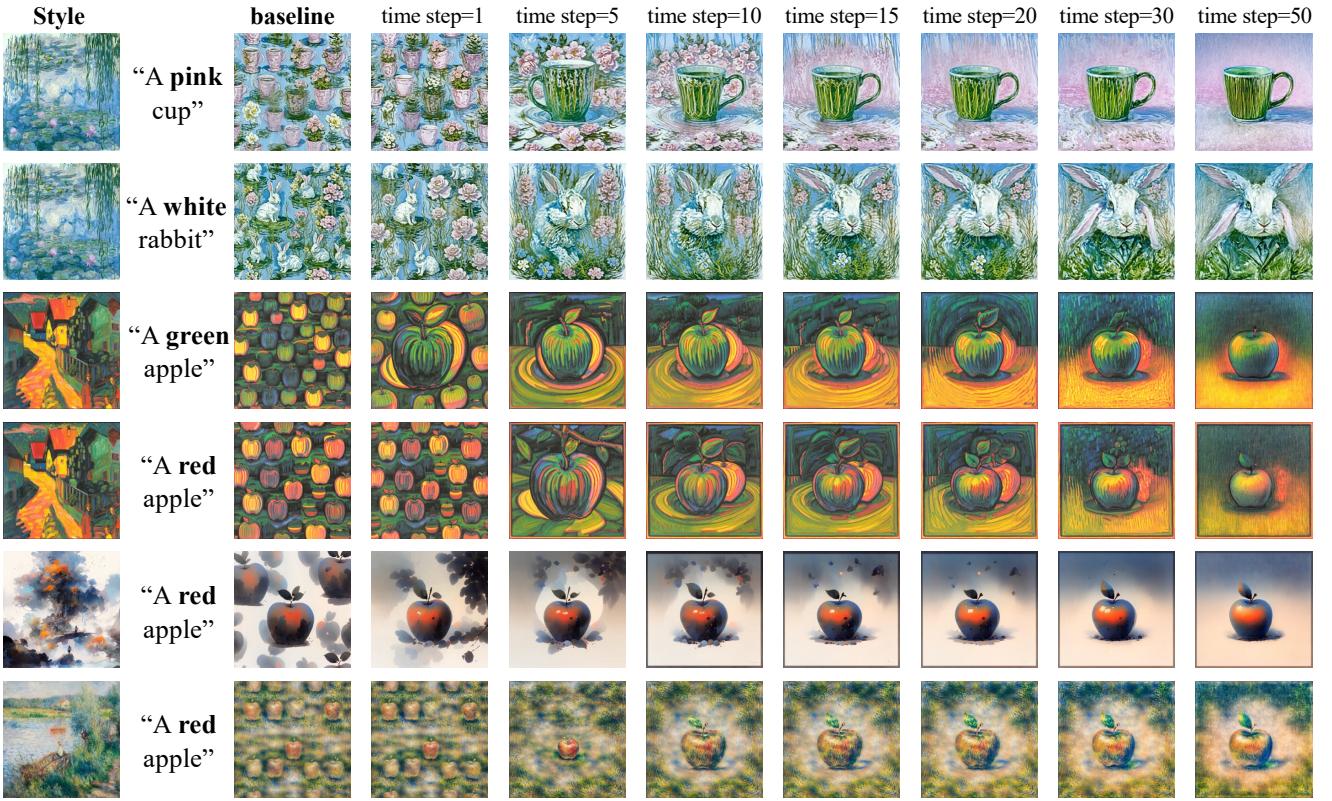


Figure 19. Impact of Teacher Model on Style Image Generation. The term “timestep” refers to the number of denoising steps during which the Teacher Model is involved. Notably, these experiments were conducted without incorporating cross-modal AdaIN to isolate and evaluate the specific impact of the Teacher Model on the generated results.

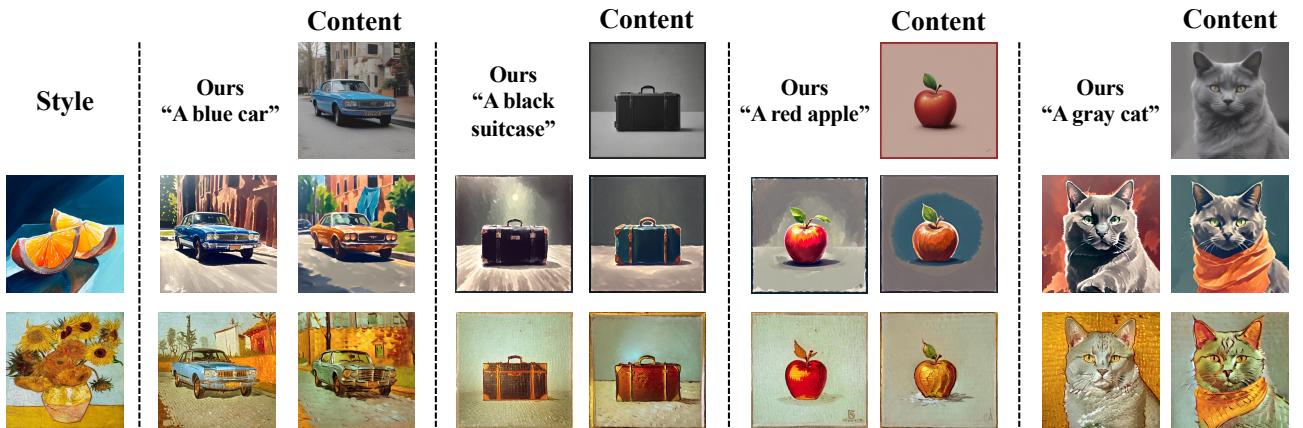


Figure 20. Compared to the image-based style transfer(I2I) provided by CSGO [37], We ensured the use of the same initial noise for both our method and the generation of the content image for I2I. It can be observed that the results obtained using the Teacher Model differ significantly from those of I2I, as I2I fails to preserve the color information of the original image.

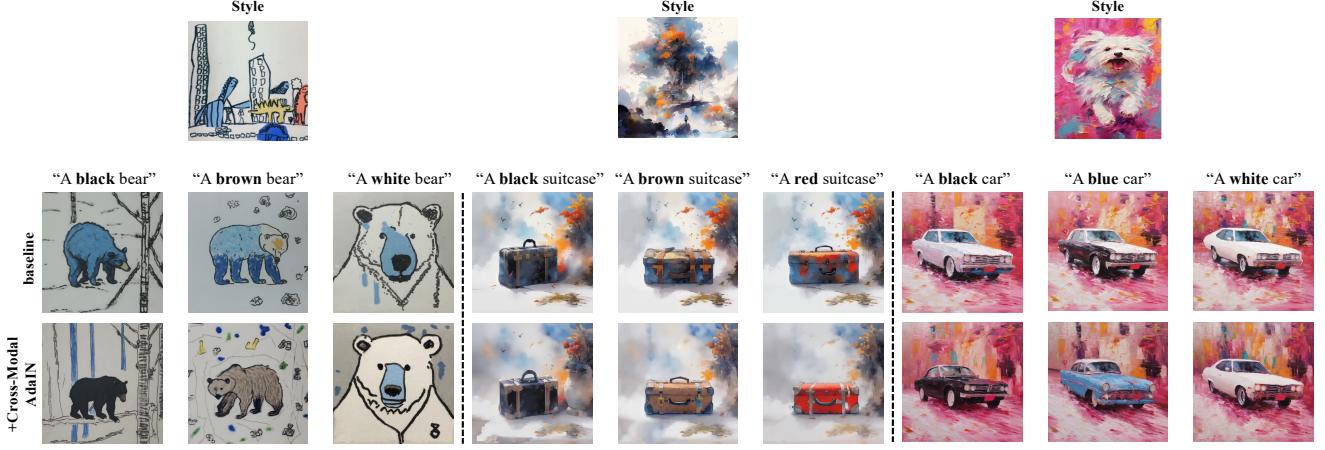


Figure 21. Qualitative results of using cross-modal AdaIN in InstantStyle [34]. The results demonstrate that cross-modal AdaIN effectively prevents style overfitting. The final generated results consistently align with the textual descriptions.

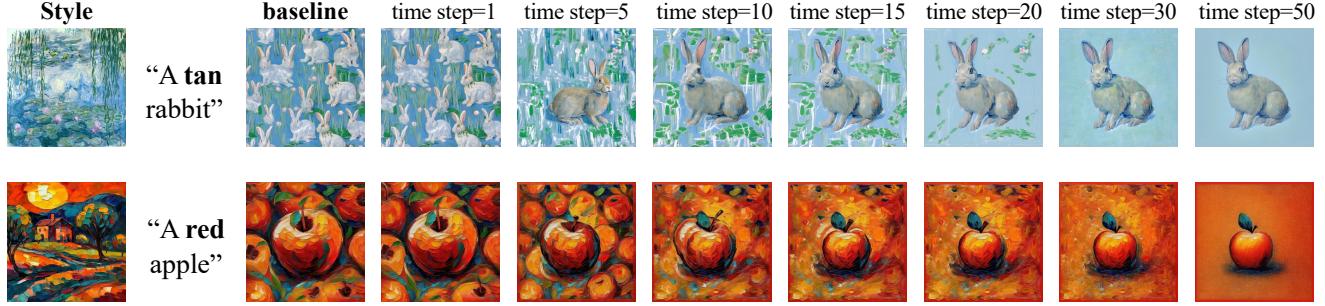


Figure 22. Impact of Teacher Model on InstantStyle [34] Image Generation. The term “timestep” refers to the number of denoising steps during which the Teacher Model is involved. Notably, these experiments were conducted without incorporating cross-modal AdaIN to isolate and evaluate the specific impact of the Teacher Model on the generated results. When the Teacher Model is applied to InstantStyle [34], it helps prevent the generation of artifacts, such as checkerboard patterns.

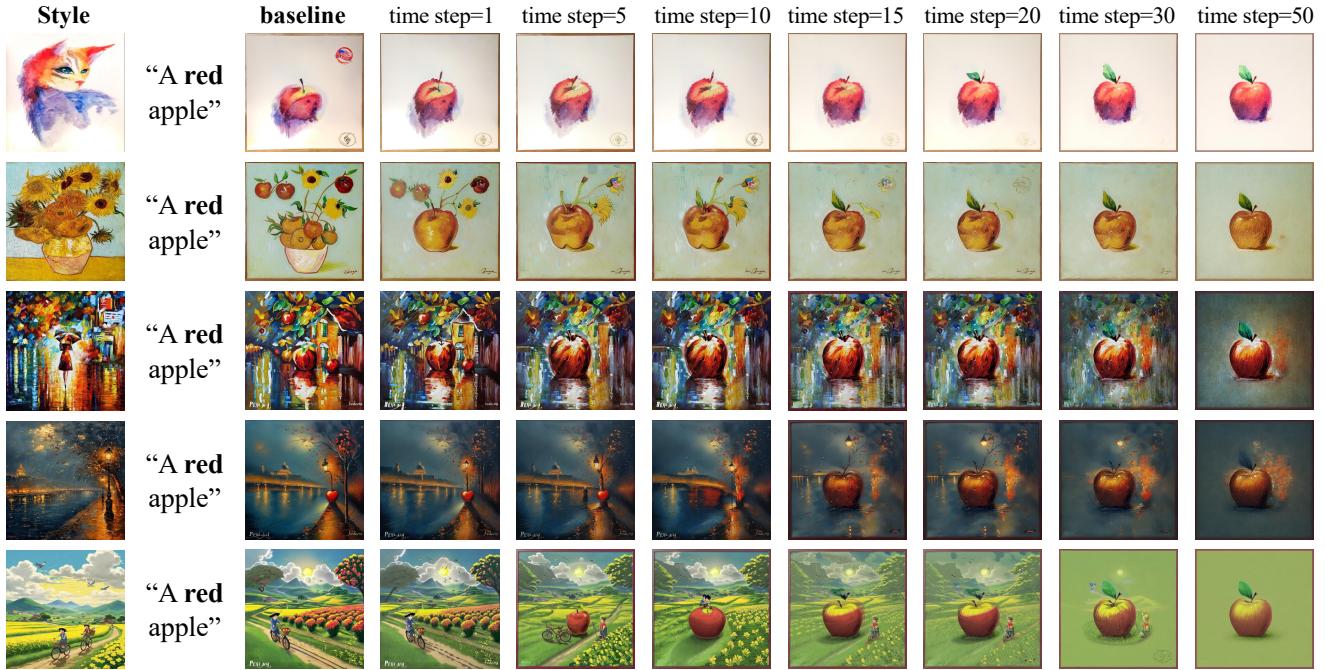


Figure 23. Impact of Teacher Model on StyleCrafter [19] Image Generation. The term “timestep” refers to the number of denoising steps during which the Teacher Model is involved. Notably, these experiments were conducted without incorporating cross-modal AdaIN to isolate and evaluate the specific impact of the Teacher Model on the generated results. In addition to ensuring layout stability, the Teacher Model also effectively reduces the occurrence of content leakage when applied to StyleCrafter [19].

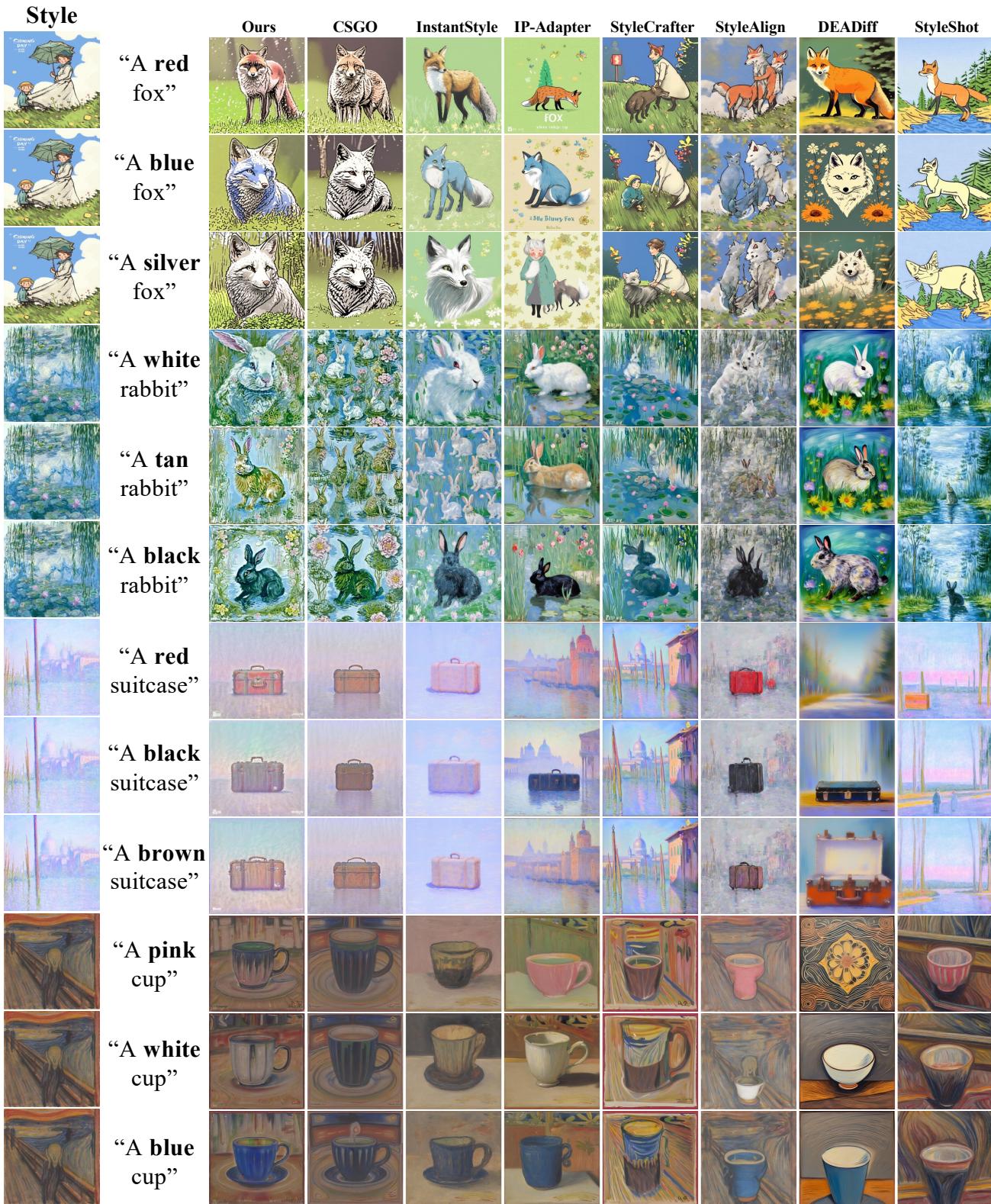


Figure 24. Qualitative comparison with state-of-the-art methods. Our approach effectively preserves image style while accurately adhering to text prompts for generation.



Figure 25. Qualitative comparison with state-of-the-art methods. Our approach effectively preserves image style while accurately adhering to text prompts for generation.

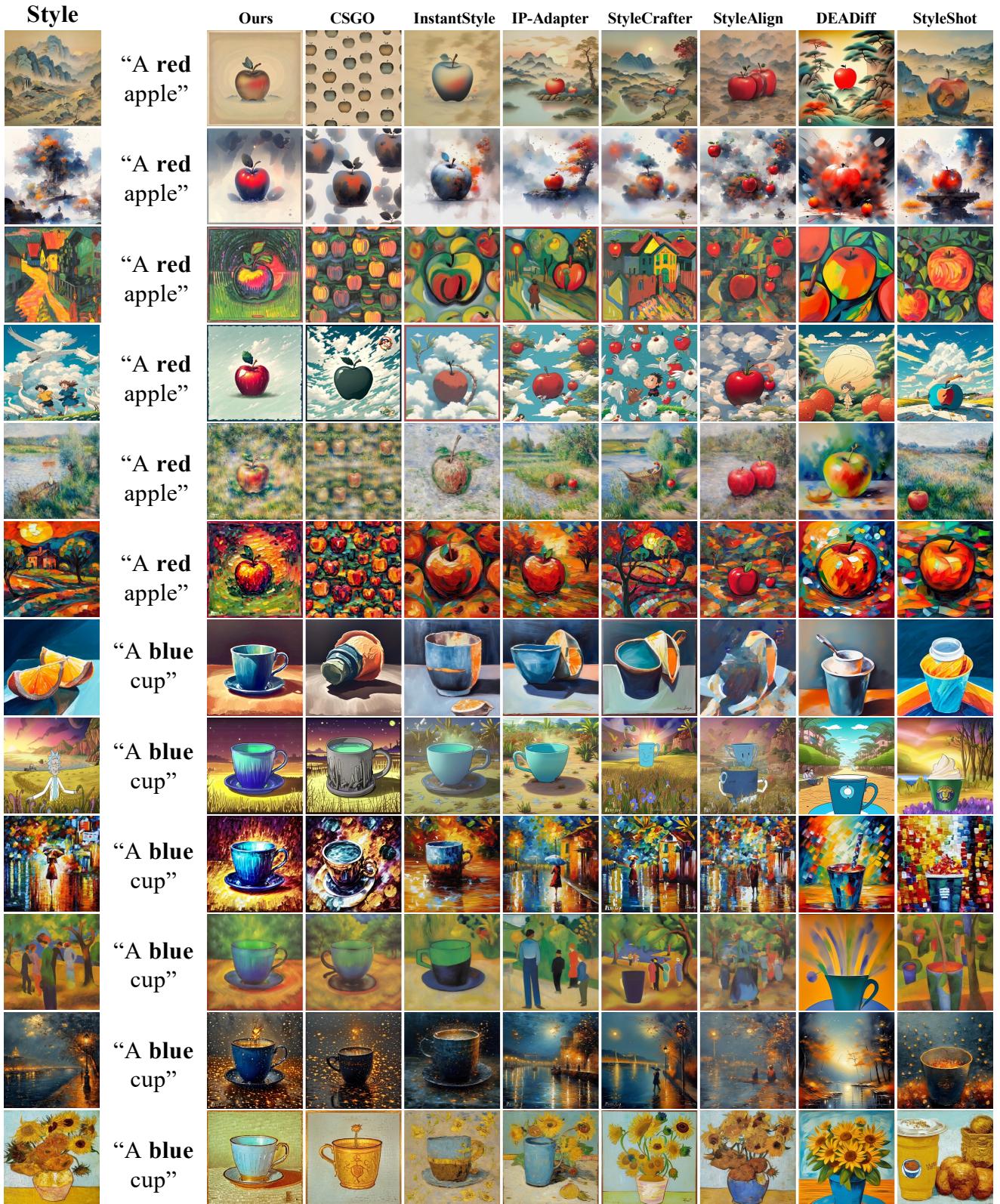


Figure 26. Qualitative comparison with state-of-the-art methods. Our approach effectively maintain layout consistency across different styles under the same prompt.

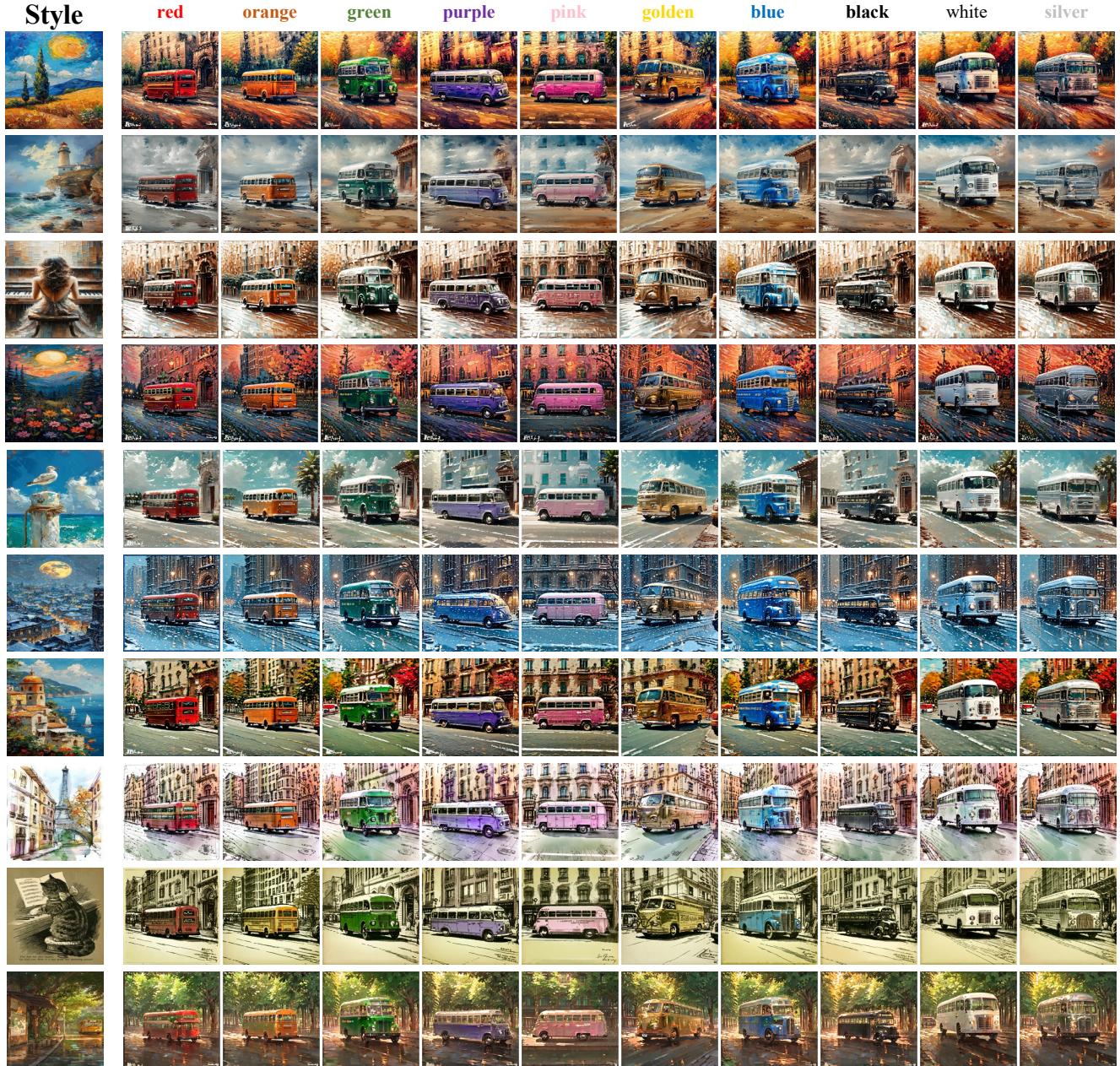


Figure 27. More results of our text-driven style transfer model. Given a style reference image, our method effectively reduces style overfitting, generating images that faithfully align with the text prompt while maintaining consistent layout structure across varying styles. Illustration of the prompt format used: “A [color] bus”.

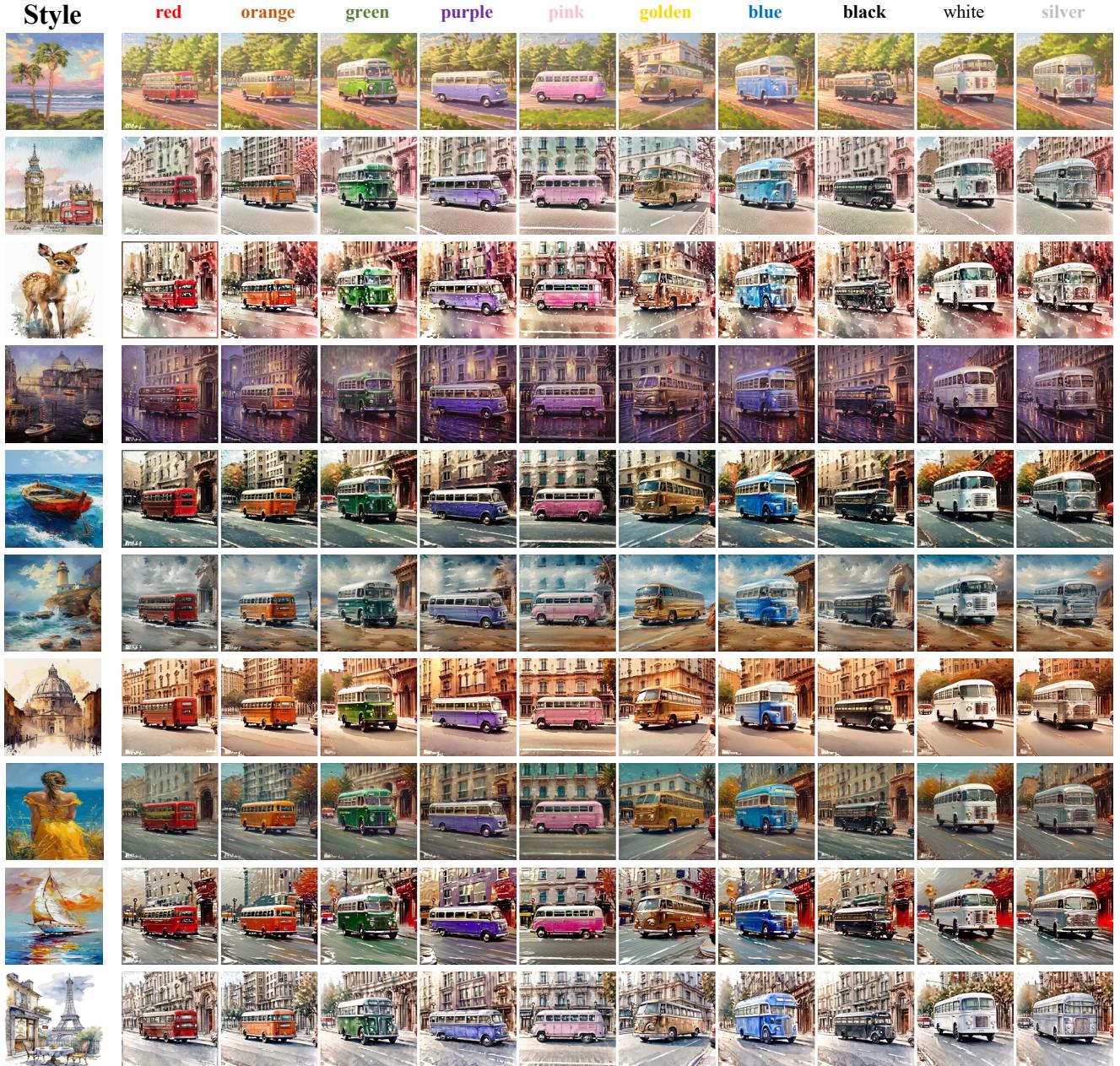


Figure 28. More results of our text-driven style transfer model. Given a style reference image, our method effectively reduces style overfitting, generating images that faithfully align with the text prompt while maintaining consistent layout structure across varying styles. Illustration of the prompt format used: “A [color] bus”.