# DreamStyler: Paint by Style Inversion with Text-to-Image Diffusion Models

**Namhyuk Ahn**[1]   **Junsoo Lee**[1]   **Chunggi Lee**[1,2]
**Kunhee Kim**[3]   **Daesik Kim**[1]   **Seung-Hun Nam**[1]   **Kibeom Hong**[4]

[1] NAVER WEBTOON AI   [2] Harvard University   [3] KAIST   [4] SwatchOn

## Abstract

Recent progresses in large-scale text-to-image models have yielded remarkable accomplishments, finding various applications in art domain. However, expressing unique characteristics of an artwork (*e.g.* brushwork, colortone, or composition) with text prompts alone may encounter limitations due to the inherent constraints of verbal description. To this end, we introduce DreamStyler, a novel framework designed for artistic image synthesis, proficient in both text-to-image synthesis and style transfer. DreamStyler optimizes a multi-stage textual embedding with a context-aware text prompt, resulting in prominent image quality. In addition, with content and style guidance, DreamStyler exhibits flexibility to accommodate a range of style references. Experimental results demonstrate its superior performance across multiple scenarios, suggesting its promising potential in artistic product creation. Project page: https://nmhkahn.github.io/dreamstyler/.

## Introduction

"*Painting is silent poetry.*" — Simonides, Greek poet

Recent text-to-image models have shown unprecedented proficiency in translating natural language into compelling visual imagery (Saharia et al. 2022; Ramesh et al. 2022; Rombach et al. 2022). These have emerged in the realm of art, providing inspiration and even assisting in crafting tangible art pieces. In the AI-assisted art production workflow, artists typically utilize various descriptive prompts that depict the style and context to generate their desired image. However, the unique styles of a painting, its intricate brushwork, light, colortone, or composition, cannot be easily described in a single word. For instance, dare we simplify the entirety of Vincent Van Gogh's lifelong artworks as just one word, 'Gogh style'? Text descriptions cannot fully evoke his unique style in our imagination — his vibrant color, dramatic light, and rough yet vigorous brushwork.

Beyond text description, recent studies (Gal et al. 2022; Ruiz et al. 2023) embed specific attributes of input images into latent space. While they effectively encapsulate a novel *object*, we observed that they struggle to personalize *style* of a painting. For instance, model optimization-based methods (Ruiz et al. 2023; Kumari et al. 2023) are highly susceptible to overfitting and often neglect inference prompts,
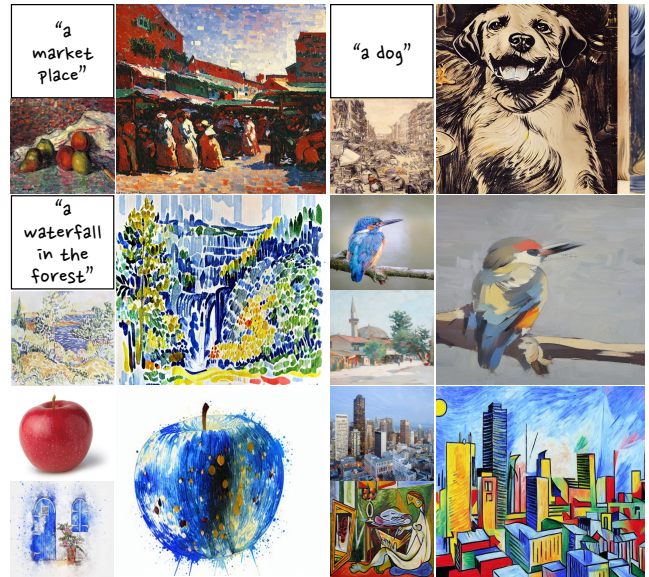
Figure 1: **DreamStyler** synthesizes outputs based on a given context along with a style reference. Note that each model is trained on a single style image shown in this figure.

which is not ideal for real-world production (please refer to the Suppl. for more details). Textual inversion-based methods (Gal et al. 2022; Voynov et al. 2023), in contrast, effectively reflect the inference prompt but fail to replicate style, possibly due to the limited capacity of the learned embeddings. This is because capturing style, from global elements (*e.g.* colortone) to local details (*e.g.* detailed texture), is challenging when relying solely on a single embedding token.

In this work, we present **DreamStyler**, a novel single (one-shot) reference-guided artistic image synthesis framework designed for the text-to-image generation and style transfer tasks (Figure 1). We encapsulate the intricate styles of artworks into CLIP text space. DreamStyler is grounded in textual inversion (TI), chosen for the inherent flexibility that stems from its prompt-based configuration. To overcome the limitations of TI, we introduce an extended textual embedding space, $\mathcal{S}$ by expanding textual embedding into the denoising timestep domain (Figure 2). Based on this space, we propose a multi-stage TI, which maps the textual

information into the $\mathcal{S}$ space. It accomplishes by segmenting the entire diffusion process into multiple *stages* (a chunk of timesteps) and allocating each textual embedding vector to the corresponding stage. The exploitation of the timestep domain in textual inversion significantly improves the overall efficacy of artistic image synthesis. This enhancement stems from the increased capacity of the personalized module, as well as the utilization of prior knowledge suggesting that different denoising diffusion steps contribute differently to image synthesis (Balaji et al. 2022; Choi et al. 2022).

We further propose a context-aware prompt augmentation that simply yet proficiently decouples the style and context information from the reference image. With our approach, the personalization module can embed style features solely into its textual embeddings, ensuring a more faithful reflection of the reference's style. To further refine the artistic image synthesis, we introduce a style and context guidance, inspired by classifier-free guidance (Ho and Salimans 2022). Our guidance bisects the guidance term into style and context components, enabling individual control. Such a guidance design allows users to tailor the outputs based on their preferences or intricacy of the reference image's style.

We validate the effectiveness of DreamStyler through a broad range of experiments. DreamStyler not only demonstrates advanced artistic image synthesis but also paves the new way of applying text-to-image diffusion models to the realms of artistic image synthesis and style transfer tasks.

## Related Work

**Personalized text-to-image synthesis.** Since the latent-based text conditional generation has been explored (Rombach et al. 2022), following studies (Saharia et al. 2022; Ramesh et al. 2022; Li et al. 2022) have further contributed to enhancing text-to-image synthesis with CLIP (Radford et al. 2021) guidance. Furthermore, Textual inversion (Gal et al. 2022), DreamBooth (Ruiz et al. 2023) and CustomDiffusion (Kumari et al. 2023) introduced approaches that leverage 3-5 images of the subject to personalize semantic features. Recently, Voynov et al. (2023) proposed $\mathcal{P}+$ space, which consists of multiple textual conditions, derived from per-layer prompts. Although they showed promising results in penalization of diffusion models, there are still limitations to fully capturing precise artistic style representations. In contrast, DreamStyler considers the denoising timestep to accommodate temporal dynamics in the diffusion process, achieving high-quality artistic image generation.

**Paint by style.** Neural style transfer renders the context of a source with a style image. Since Gatys, Ecker, and Bethge (2016), studies have been devoted to enhancing the transfer networks for more accurate and convincing style transfer. Notably, AdaIN (Huang and Belongie 2017) and AdaAttN (Liu et al. 2021) investigated matching the second-order statistics of content and style images. AesPA-Net (Hong et al. 2023) and StyTr$^2$ (Deng et al. 2022) adopted recent architectures such as attention or transformer for high-fidelity neural style transfer. Recently, InST (Zhang et al. 2023) utilized the diffusion models by introducing the image encoder to inverse style images into CLIP spaces.
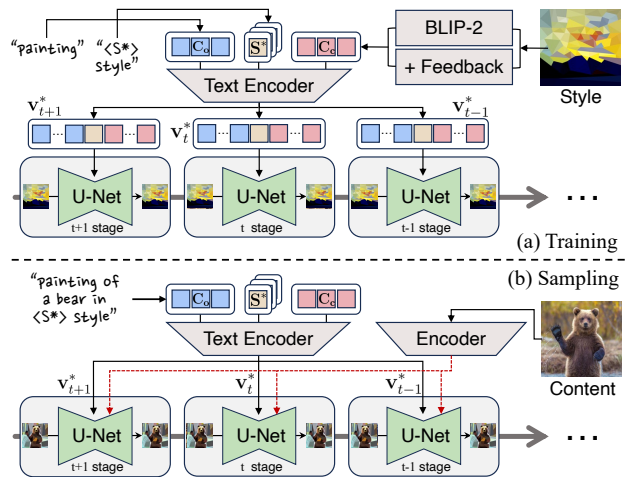


Figure 2: **Model overview. (a)** DreamStyler constructs training prompt with an opening text $C_o$, multi-stage style tokens $\mathbf{S}^*$, and a context description $C_c$, which is captioned with BLIP-2 and human feedback. DreamStyler projects the training prompt into multi-stage textual embeddings $\mathbf{v}^* = \{v_1^*, \ldots, v_T^*\}$, where $T$ is # stages (a chunk of the denoising timestep). As a result, the denoising U-Net provides distinct textual information at each stage. **(b)** DreamStyler prepares the textual embedding using a provided inference prompt. For style transfer, DreamStyler employs ControlNet to comprehend the context information from a content image.

## Method

**Preliminary: Stable Diffusion (SD).** DreamStyler is built upon SD (Rombach et al. 2022). SD projects an input image $x$ into a latent code, $z = E(x)$ using an encoder $E$, while decoder $D$ transforms the latent code back into pixel space, *i.e.* $x' = D(z')$. The diffusion model creates a new latent code $z'$ by conditioning on additional inputs such as a text prompt $y$. The training objective of SD is defined as:

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

At each timestep $t$, the denoising network $\epsilon_\theta$ reconstructs the noised latent code $z_t$, given the timestep $t$ and a conditioning vector $c(y)$. To generate $c(y)$, each token from a prompt is converted into an embedding vector, which is then passed to the CLIP text encoder (Radford et al. 2021).

**Preliminary: Textual Inversion (TI).** Gal et al. (2022) proposed a method to personalize a pre-trained text-to-image model by incorporating a novel embedding representing the intended concept. To personalize the concept, they initialize a word token $S^*$ and its corresponding vector $v^*$, situated in the textual conditioning space $\mathcal{P}$, which is the output of the CLIP text encoder. Instead of altering any weights in SD models, they optimize $v^*$ alone using Eq. (1). To create images of personalized concepts, the inclusion of $S^*$ in the prompts (*e.g.* a photo of $S^*$ dog) is the only required step.

### Multi-Stage Textual Inversion

In some cases, TI fails to sufficiently represent the concept due to the inherent capacity limitations associated with us-
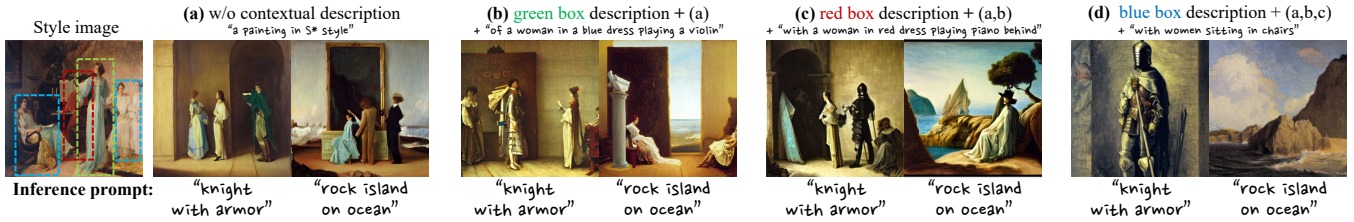
Figure 3: **How does training prompt affect?** Given a style image, we construct training prompts with contextual descriptions (b∼d). **(a)** Training without contextual description in the prompt; *i.e.* trains the model with "a painting in $S^*$ style". The model tends to generate the images that contains objects and compositions from the style image (*e.g.* standing and sitting audiences), instead of attributes depicted in the inference prompt. **(b, c)** Training with partial contextual descriptions (the green and red boxes displayed in the style image, respectively). Such a tendency is significantly reduced, yet the model still synthesizes some objects from the style image (*e.g.* sitting people in the blue box). **(d)** Training with full contextual descriptions. The model produces outputs that fully reflect the inference prompt without introducing any non-style attributes from the style image.

ing a single embedding token. Moreover, this single embedding strategy is inappropriate for accommodating the changing process of diffusion models. As explored in Balaji et al. (2022); Choi et al. (2022), diffusion models display intriguing temporal dynamics throughout the process, necessitating different capacities at various diffusion steps. In light of this, managing all denoising timesteps with a single embedding potentially has limitations due to the spectrum of local to global expressions embodied in paintings. Thus, articulating paintings is intricately related to the denoising timesteps, which operate in a coarse-to-fine synthesis manner (Balaji et al. 2022). To address these challenges, we introduce a *multi-stage TI* that employs multiple embeddings, each corresponding to specific diffusion stages (Figure 2).

We first propose an extended textual embedding space $\mathcal{S}$. The premise of the $\mathcal{S}$ space is to decompose the entire diffusion process into multiple distinct *stages*. To implement this, we split the denoising timesteps into $T$ chunks and denote each chunk as a stage. Based on the $\mathcal{S}$ space, the multi-stage TI prepares the copies of the initial style token ($S^*$) as a multi-stage token set $\mathbf{S}^* = \{S_1^*, \ldots, S_T^*\}$. In this way, the multi-stage TI projects a style image into $T$ style tokens, contrasting the TI that embeds it into a single token. The token set is then encoded by a CLIP text encoder to form stagewise embedding vectors, denoted as $\mathbf{v}^* = \{v_1^*, \ldots, v_T^*\}$. Lastly, the multi-stage TI optimizes these embeddings following the subsequent equation.

$$\mathbf{v}^* = \arg\min_{\mathbf{v}} \mathbb{E}_{z,\mathbf{v},\epsilon,t}[||\epsilon - \epsilon_\theta(z_t, t, c(v_t))||_2^2]. \quad (2)$$

The application of multi-stage TI significantly enhances the representation capacity beyond that of vanilla TI, which we will illustrate in a series of experiments. Furthermore, this method enables the fusion of multiple tokens, each originating from different styles, at a specific stage $t$. Consequently, it facilitates the creation of unique and novel styles tailored to the user's individual preferences.

## Context-Aware Text Prompt

While the multi-stage TI enhances representational capacity, it still faces fundamental problems when training with a style reference; the style and context of the image may become entangled during the optimization of the embeddings.

This problem mainly arises from attempts to encapsulate all features of the image into $S^*$, not just the style aspect. As depicted in Figure 3, without contextual information in the training prompt, the model overlooks the context of inference prompt. However, when we inject contextual descriptions into the training prompt, the model better disentangles the style from the context. In our observations, such a phenomenon occurs more frequently as the representational capacity increases, likely due to the model's increased efforts to accommodate all information within its capacity.

Hence, we construct training prompts to include contextual information about the style image. Let $C = [C_o, \mathbf{S}^*]$ be the vanilla prompt used in multi-stage TI training, where $C_o$ is the opening text (*e.g.* "a painting"), and $\mathbf{S}^*$ is multi-stage style token set, described above. In the proposed strategy, we incorporate a contextual descriptor $C_c$ (*e.g.* "of a woman in a blue dress") into the middle of the prompt (Figure 2), *i.e.* $C = [C_o, C_c, \mathbf{S}^*]$. We annotate all the non-style attributes (*e.g.* objects, composition, and background) from the style image to form the contextual descriptor. When we caption non-style attributes, BLIP-2 (Li et al. 2023) is employed to aid in the automatic prompt generation.

Although a context-aware prompt significantly reinforces style-context decoupling, for some style images with complicated contexts (Figure 3), BLIP-2 might not capture all details, which could limit the model's disentanglement capability. In such cases, we further refine caption $C_c$ based on human feedback (e.g., caption by humans). This human-in-the-loop strategy is straightforward yet markedly improves the model's ability to disentangle styles. Since our goal is one-shot model training, the time spent refining the caption is minimal; typically less than a minute. With the context-aware prompt, the text-to-image models can now distinguish style elements from contextual ones and specifically embed these into the (multi-stage) style embeddings $\mathbf{v}^*$. The motivation for augmenting the training prompt is also suggested in StyleDrop (Sohn et al. 2023), a current personalization approach in the text-to-image diffusion model.

## Style and Context Guidance

Classifier-free guidance (Ho and Salimans 2022) improves conditional image synthesis. It samples adjusted noise pre-
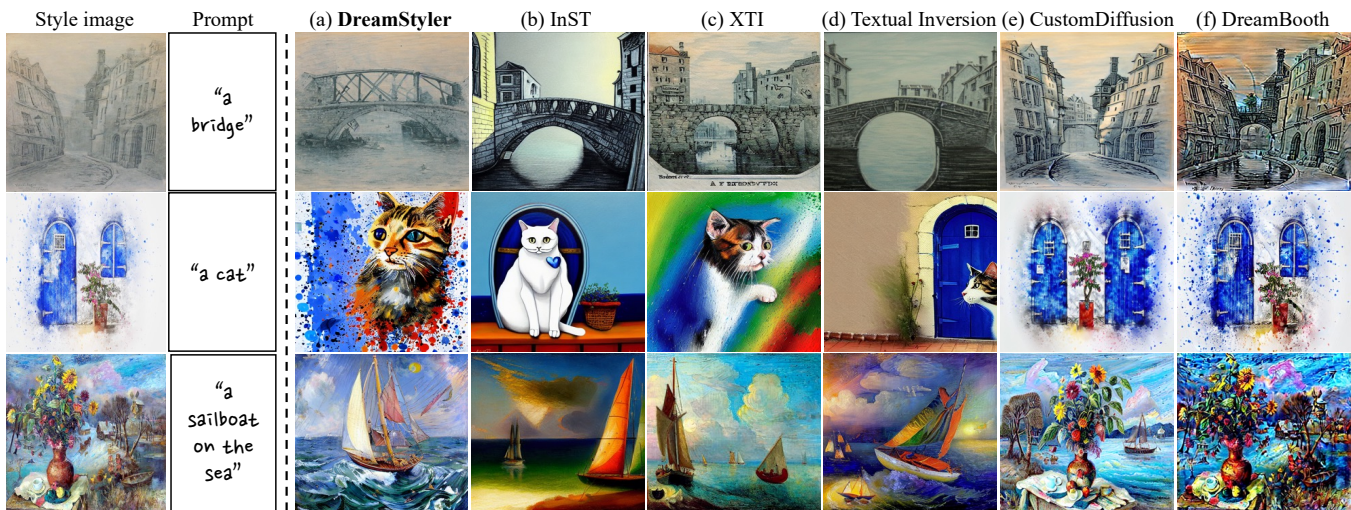
Figure 4: **Qualitative comparison** on the style-guided text-to-image synthesis task.
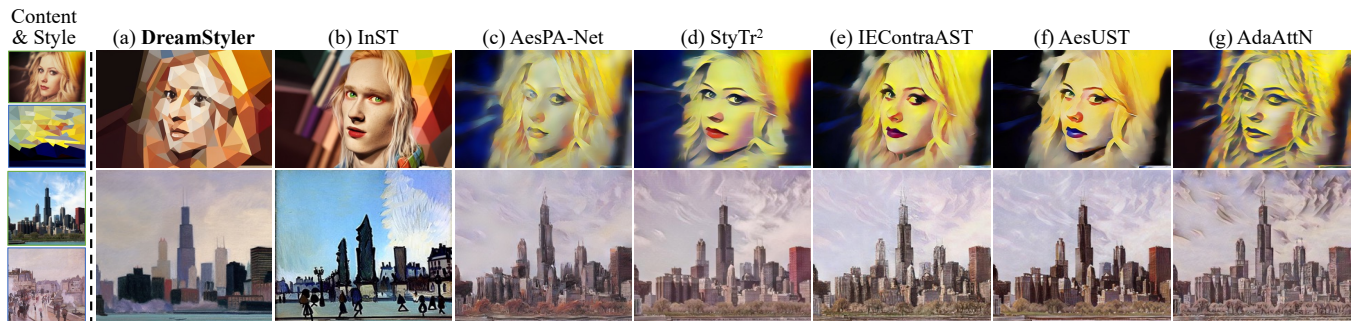


Figure 5: **Qualitative comparison** on the style transfer task.

diction $\hat{\epsilon}(.)$, by leveraging unconditional output under null token $\emptyset$ as: $\hat{\epsilon}(\mathbf{v}) = \epsilon(\emptyset) + \lambda(\epsilon(\mathbf{v}) - \epsilon(\emptyset))$, where, $\lambda$ is the guidance scale and we omit $c(.)$, $z$ and $t$ for brevity.

In style-guided image synthesis, this guidance pushes both style and context uniformly with $\lambda$. The uniform guidance could face limitations since the spectrum of "style" of artistic paintings is wider than that of natural photos. Given this diversity, a more nuanced control mechanism is required. Furthermore, there exist demands to individually control style and context in the art-making process. To this end, we propose style and context guidance as in below.

$$\hat{\epsilon}(\mathbf{v}) = \epsilon(\emptyset) + \lambda_s[\epsilon(\mathbf{v}) - \epsilon(\mathbf{v_c})] + \lambda_c[\epsilon(\mathbf{v_c}) - \epsilon(\emptyset)]$$
$$+ \lambda_c[\epsilon(\mathbf{v}) - \epsilon(\mathbf{v_s})] + \lambda_s[\epsilon(\mathbf{v_s}) - \epsilon(\emptyset)] \quad (3)$$

where, $\mathbf{v}_s, \mathbf{v}_c$ are the embeddings of prompts $C, C_c$, respectively. $\lambda_s, \lambda_c$ denote style and context guidance scale. We derive Eq. (3) by decomposing $\mathbf{v}$ into $\mathbf{v}_s, \mathbf{v}_c$. We employ two paired terms to balance the influence of each guidance. Please refer to Suppl. for detailed derivation and analysis.

By separating the guidance into style and context, users are afforded the flexibility to control these elements individually. Specifically, an increase in $\lambda_c$ increases the model's sensitivity towards context (*e.g.* inference prompt or content image), whereas amplifying $\lambda_s$ leads the model towards a more faithful style reproduction. This flexible design allows users to generate stylistic output tailored to their individual preferences, and it also facilitates the adoption of various styles, each with a range of complexities (Hong et al. 2023).

## Style Transfer

DreamStyler transmits styles by inverting a content image into a noisy sample and then denoising it towards the style domain (Meng et al. 2021). With this approach, however, the preservation of content would be suboptimal (Ahn et al. 2023). To improve this, we inject additional conditions from the content image into the model (Zhang and Agrawala 2023) (Figure 2). This straightforward pipeline well preserves with the structure of the content image, while effectively replicating styles. Moreover, by leveraging a powerful prior knowledge from text-to-image models, the style quality of DreamStyler surpasses that of traditional methods.

## Experiment

**Implementation details.** We use $T = 6$ for multi-stage TI and utilize human feedback-based context prompts by default. Please refer to Suppl. for more details.

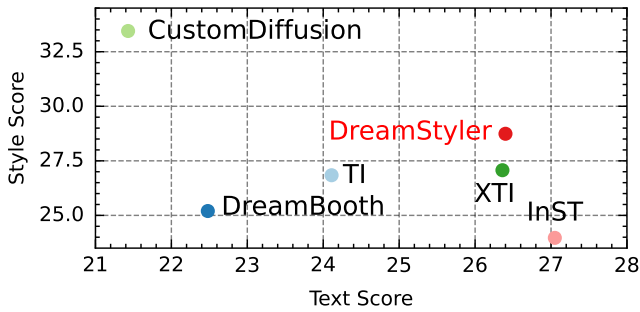**Datasets.** We collected a set of 32 images representing var-

Figure 6: **Performance of text and style scores** in style-guided text-to-image synthesis. DreamStyler effectively balances these metrics and surpasses the majority of methods.

| Method | Text Score | Style Score | User Score |
|---|---|---|---|
| Textual Inversion (Gal et al. 2022) | 24.11 | 26.84 | 2.1% |
| DreamBooth (Ruiz et al. 2023) | 22.48 | 25.20 | 3.9% |
| CustomDiffusion (Kumari et al. 2023) | 21.43 | **33.45** | <u>4.8%</u> |
| XTI (Voynov et al. 2023) | 26.36 | 27.07 | 4.5% |
| InST (Zhang et al. 2023) | **27.05** | 23.97 | 1.8% |
| **DreamStyler (Ours)** | <u>26.40</u> | <u>28.74</u> | **82.9%** |

Table 1: **Quantitative comparison** on the style-guided text-to-image synthesis task. **Bold**: best, <u>underline</u>: second best.

ious artistic styles, following the literature on style transfer (Tan et al. 2019). To evaluate text-to-image synthesis, we prepared 40 text prompts, as described in Suppl.

**Baselines.** In terms of text-to-image synthesis, we compare DreamStyler against diffusion-based personalized methods, ranging from textual inversion to model-optimization approaches. For the style transfer task, we compare our method to state-of-the-art style transfer frameworks. We utilize official codes for all the methods used in the comparison.

**Evaluation.** Text and image scores, based on CLIP, measure the alignment with a given text prompt and style image, respectively. Style score assesses the style consistency by calculating the similarity of Gram features between the style and generated images. More details are provided in Suppl.

### Style-Guided Text-to-Image Synthesis

Table 1 and Figure 6 show quantitative results. DreamStyler delivers a robust performance while managing the trade-off between text and style scores. A tendency is noted that an overemphasis on input text prompts may lead to a compromise in style quality. Despite this, DreamStyler effectively balances these aspects, yielding a performance that goes beyond the trade-off line, indicative of outstanding capability. User score also supports the distinction of DreamStyler.

As shown in Figure 4, previous inversion-based methods (TI, InST, and XTI) effectively preserve the context of text prompts but fall short in adopting the intrinsic artwork of style images. Conversely, the model optimization-based methods (DreamBooth, CustomDiffusion) excel in delivering styles but struggle to adhere to the prompt or introduce



Figure 7: **My object in my style.** Textual inversion faces challenges in accurately capturing both style and context from the reference images. Although CustomDiffusion successfully recreates the object's appearance, it tends to generate objects in a realistic style, which does not entirely match the target style image. On the other hand, DreamStyler excels at synthesizing the object in the user-specified style.

| Method | Text Score | Image Score | User Score |
|---|---|---|---|
| AdaAttN (Liu et al. 2021) | 56.67 | 56.76 | 8.6% |
| AesUST (Wang et al. 2022) | 58.05 | 58.09 | 6.8% |
| IEContraAST (Chen et al. 2021) | 59.38 | 59.42 | 8.6% |
| StyTr$^2$ (Deng et al. 2022) | 56.18 | 56.28 | <u>21.2%</u> |
| AesPA-Net (Hong et al. 2023) | 58.08 | 58.15 | 8.6% |
| InST (Zhang et al. 2023) | <u>65.32</u> | <u>65.37</u> | 2.3% |
| **DreamStyler (Ours)** | **66.04** | **66.05** | **44.1%** |

Table 2: **Quantitative comparison** on the style transfer task.

objects in style images (3rd row). DreamStyler, in contrast, not only faithfully follows text prompts but also accurately reflects the delicate artistic features of style images.

### Style Transfer

As an extended application, DreamStyler also conducts style transfer. As shown in Table 2, we quantitatively compare with previous style transfer studies. Note that since most prior studies have employed Gram loss to boost style quality, we report a CLIP-based image score as an evaluation metric for a more fair comparison. In this benchmark, DreamStyler achieves state-of-the-art performance across text and image scores as well as user preference. Figure 5 also provides evidence of DreamStyler's effectiveness. Our method adeptly captures style features such as polygon shapes or subtle brushwork present in style images. These results highlight the method's capacity to accurately mirror both the thematic intent and the stylistic nuances of the source artwork.

### Stylize My Own Object in My Own Style

Beyond style transfer that stylizes *my image*, one might desire to stylize *my object* (Sohn et al. 2023). In such a scenario, a user leverages both their object and style images. As DreamStyler employs an inversion-based approach, this can be readily accomplished by simply training an additional embedding for the object. Subsequently, the user
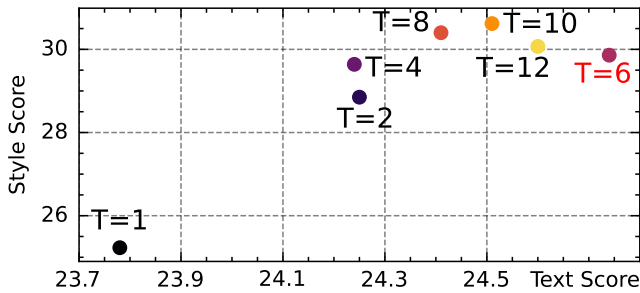
Figure 8: **Study on the number of stages** ($T$) in multi-stage TI. We vary $T$ from 1 to 12 and select $T = 6$ as the final model, considering the trade-off between text and style.
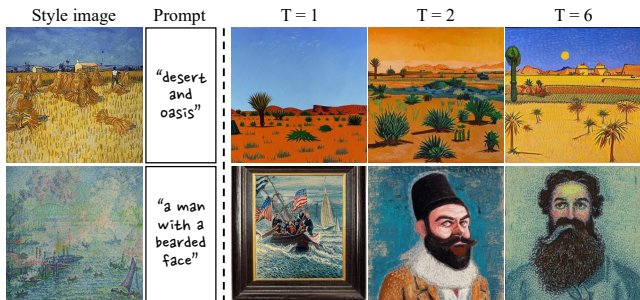


Figure 9: **Visual comparison of varying** $T$ in multi-stage TI. At $T = 1$, the model fails in both style replication and prompt understanding. As $T$ increases, the style quality and text alignment are drastically enhanced.

freely merges style and object tokens in the inference prompt to generate images. As depicted in Figure 7, DreamStyler excels in accurately reflecting both the style and object

## Model Analysis

**Ablation study.** In Table 3, we evaluate each component of our method. The usage of multi-stage TI substantially augments both the text and style score, with a marked increase in style quality, accentuating the pivotal role of this module in creating artistic stylization products. A context-aware prompt yields a modest alteration in the quantitative metrics, yet provides a considerable contribution to the qualitative, which we will discuss in the following section. Style and context (S&C) guidance considerably impacts scores, reinforcing its significance in sustaining the comprehensive quality and coherence of the generated outputs.

**Multi-stage TI.** In Figure 8, we delve into the influence of the number of stages ($T$) on performance. A transition from $T = 1$ to 4 results in substantial improvement. Upon reaching $T = 6$, the performance begins to navigate trade-off contours, prompting us to select $T = 6$ for the final model, as we seek to improve the text alignment of the synthesized images. Nevertheless, users have the flexibility to choose a different $T$ value according to their preference. In Figure 9, we provide a visual comparison of the outcomes when $T$ is set to 1, 2, and 6. While $T = 1$ struggles to reflect the artistic features of the style image or comprehend the input prompt,

| Method | Text Score | Style Score |
|---|---|---|
| Baseline (Gal et al. 2022) | 23.78 | 25.23 |
| + Multi-Stage TI | 24.74 | **29.86** |
| + Context-Aware Prompt | 24.65 | 29.50 |
| + S&C Guidance (**Ours**) | **25.38** | 29.62 |

Table 3: **Model ablation study.** Upon the textual inversion baseline (Gal et al. 2022), we attach the proposed components to measure the effectiveness of our method.
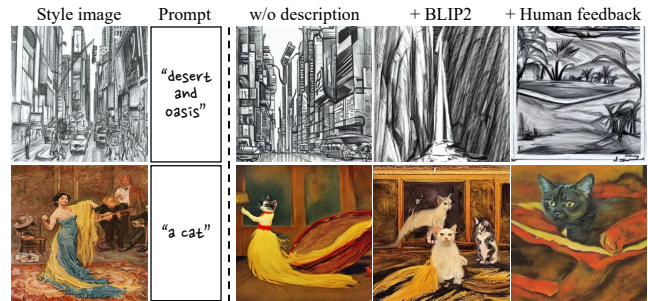


Figure 10: **Comparison of three prompt strategies.** The model trained without contextual description struggles to disentangle style and context from the style image, generating elements present in the style reference (*e.g.* the same composition in 1st row, a yellow dress in 2nd row). The contextual prompt alleviates this issue to some extent, but the BLIP2-based construction cannot completely eliminate it (*e.g.* the same vanishing point in 1st row). The issue is thoroughly addressed when human feedback is utilized.

$T = 2$ uplifts the quality, yet it also falls short of embracing the style. In contrast, $T = 6$ proves proficient at mimicking the style image, effectively replicating delicate brushwork (1st row) or emulating the pointillism style (2nd row).

**Context-aware prompt.** Figure 10 presents a visual comparison of three prompt constructions. Training the model without any contextual description (*i.e.* using "A painting in $S^*$ style.") poses a significant challenge, as it struggles to distinguish style from the context within the style image. Subsequently, this often results in the generation of elements that exist in the style reference, such as objects or scene perspective. The introduction of a contextual prompt considerably alleviates this issue, aiding the model in better separating stylistic elements from context. However, the automatic prompt construction does not fully resolve this, as BLIP-based captions often fail to capture all the details of the style image. The most effective solution is leveraging human feedback in the construction of prompts. This approach effectively tackles the issue, resulting in a more robust separation of style and context in the generated outputs.

**Guidance.** In Figure 11, we explore style and context guidance by adjusting the scale parameters. When we amplified the style guidance strength ($\lambda_s$), the model mirrors the style image, illustrating style guidance's capability in managing the image's aesthetics. Yet, overemphasis on style risks compromising the context, leading to outputs that, while stylis-
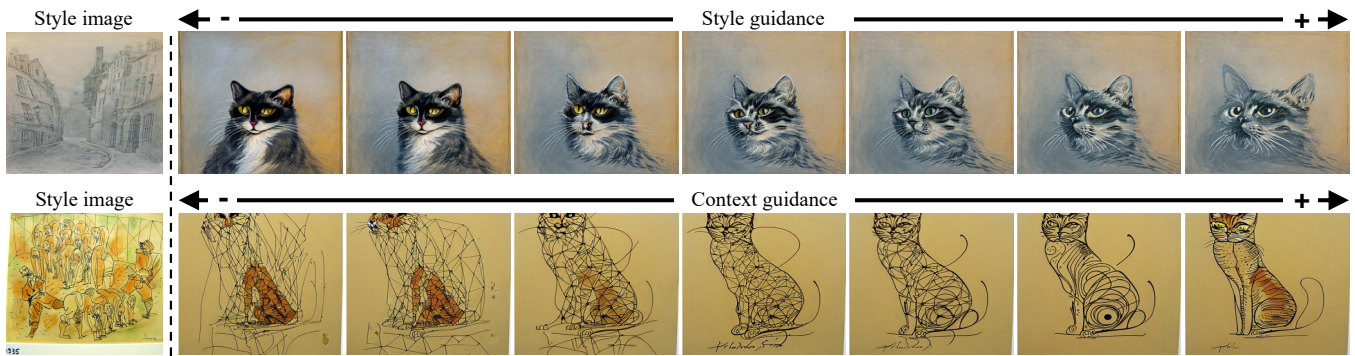
Figure 11: **Study on the style and context guidance.** Inference prompt: "A cat". By adjusting the scale parameters $(\lambda_s, \lambda_c)$, we assess the influence of style and context guidance on the synthesized image. Increasing the style guidance strength causes the model to align more closely with the aesthetics of the style image; however, an excessive emphasis on style could compromise the context. Conversely, increasing the context guidance strength ensures the output corresponds with the inference prompt, but overly strong context guidance could deviate the output from the original style.
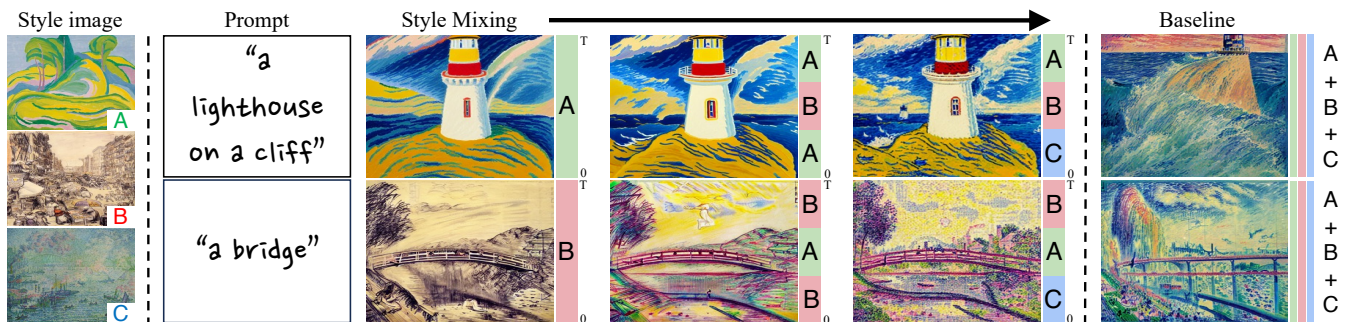


Figure 12: **Style mixing.** Multi-stage TI facilitates style mixing from various style references. A user can customize a new style by substituting style tokens at different stages $t$. For example, the style token closer to $t = T$ tends to influence the structure of the image, while those closer to $t = 0$ have a stronger effect on local and detailed attributes. For comparison, we display the baseline that employs all style tokens at every stage (*i.e.* using "A painting in $S_t^A$, $S_t^B$, $S_t^C$ style" at all stages).

tically congruent, might diverge from the intended context. On the other hand, strengthening context guidance $(\lambda_c)$ ensures the output resembles the inference prompt, highlighting context guidance's essential role in preserving contextual integrity. However, excessively strong context guidance could steer the output away from the original style, underlining the need for a nuanced balance of guidance for generating visually appealing and contextually accurate images. Nevertheless, this offers a new dimension of control over the synthesized image, differing from the classifier-free guidance (Ho and Salimans 2022). The additional control is a crucial element in the workflow of digital art production, considering its delicate and nuanced final outcomes.

**Style mixing.** As shown in Figure 12, multi-stage TI opens up a novel avenue for an intriguing aspect of style mixing from diverse style references. This process empowers users to customize a unique style by deploying different style tokens at each stage $t$. The style tokens close to $t = T$ predominantly impact the structure of the image, akin to broad strokes, while tokens closer to $t = 0$ affect local and detailed attributes, akin to intricate brushwork. To provide a concrete point of comparison, we present a baseline model

that incorporates all style tokens at every stage, using the prompt "A painting in $S_t^A$, $S_t^B$, $S_t^C$ styles". While the baseline produces reasonable style quality, it lacks a control factor for extracting partial stylistic features from the reference. Consequently, the fusion of styles with multi-stage TI underscores the creative and flexible nature of our model, offering users a broad range of applications for artistic creation.

## Conclusion

We have introduced DreamStyler, a novel image generation method with a given style reference. By optimizing multi-stage TI with a context-aware text prompt, DreamStyler achieves remarkable performance in both text-to-image synthesis and style transfer. Content and style guidance provides a more adaptable way of handling diverse style references.

**Limitations.** While DreamStyler exhibits outstanding ability in generating artistic imagery, it is important to acknowledge its limitations within the intricate context of artistic expression. The vast spectrum of artistry, spanning from primitive elements to more nuanced and abstract styles (such as surrealism), demands thorough definition and examination from both artistic and technological perspectives.

# References

Ahn, N.; Kwon, P.; Back, J.; Hong, K.; and Kim, S. 2023. Interactive Cartoonization with Controllable Perceptual Factors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16827–16835.

Balaji, Y.; Nah, S.; Huang, X.; Vahdat, A.; Song, J.; Kreis, K.; Aittala, M.; Aila, T.; Laine, S.; Catanzaro, B.; et al. 2022. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*.

Chen, H.; Wang, Z.; Zhang, H.; Zuo, Z.; Li, A.; Xing, W.; Lu, D.; et al. 2021. Artistic style transfer with internal-external learning and contrastive learning. *Advances in Neural Information Processing Systems*, 34: 26561–26573.

Choi, J.; Lee, J.; Shin, C.; Kim, S.; Kim, H.; and Yoon, S. 2022. Perception prioritized training of diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11472–11481.

Deng, Y.; Tang, F.; Dong, W.; Ma, C.; Pan, X.; Wang, L.; and Xu, C. 2022. StyTr2: Image Style Transfer with Transformers. In *CVPR*, 11326–11336.

Gal, R.; Alaluf, Y.; Atzmon, Y.; Patashnik, O.; Bermano, A. H.; Chechik, G.; and Cohen-Or, D. 2022. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*.

Gatys, L. A.; Ecker, A. S.; and Bethge, M. 2016. Image style transfer using convolutional neural networks. In *CVPR*, 2414–2423.

Ho, J.; and Salimans, T. 2022. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*.

Hong, K.; Jeon, S.; Lee, J.; Ahn, N.; Kim, K.; Lee, P.; Kim, D.; Uh, Y.; and Byun, H. 2023. AesPA-Net: Aesthetic Pattern-Aware Style Transfer Networks.

Huang, X.; and Belongie, S. 2017. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, 1501–1510.

Kumari, N.; Zhang, B.; Zhang, R.; Shechtman, E.; and Zhu, J.-Y. 2023. Multi-concept customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1931–1941.

Li, J.; Li, D.; Savarese, S.; and Hoi, S. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*.

Li, W.; Xu, X.; Xiao, X.; Liu, J.; Yang, H.; Li, G.; Wang, Z.; Feng, Z.; She, Q.; Lyu, Y.; et al. 2022. UPainting: Unified Text-to-Image Diffusion Generation with Cross-modal Guidance. *arXiv preprint arXiv:2210.16031*.

Liu, S.; Lin, T.; He, D.; Li, F.; Wang, M.; Li, X.; Sun, Z.; Li, Q.; and Ding, E. 2021. Adaattn: Revisit attention mechanism in arbitrary neural style transfer. In *Proceedings of the IEEE/CVF international conference on computer vision*, 6649–6658.

Meng, C.; He, Y.; Song, Y.; Song, J.; Wu, J.; Zhu, J.-Y.; and Ermon, S. 2021. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*.

Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PMLR.

Ramesh, A.; Dhariwal, P.; Nichol, A.; Chu, C.; and Chen, M. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*.

Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10684–10695.

Ruiz, N.; Li, Y.; Jampani, V.; Pritch, Y.; Rubinstein, M.; and Aberman, K. 2023. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 22500–22510.

Saharia, C.; Chan, W.; Saxena, S.; Li, L.; Whang, J.; Denton, E. L.; Ghasemipour, K.; Gontijo Lopes, R.; Karagol Ayan, B.; Salimans, T.; et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35: 36479–36494.

Simonyan, K.; and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.

Sohn, K.; Ruiz, N.; Lee, K.; Chin, D. C.; Blok, I.; Chang, H.; Barber, J.; Jiang, L.; Entis, G.; Li, Y.; et al. 2023. StyleDrop: Text-to-Image Generation in Any Style. *arXiv preprint arXiv:2306.00983*.

Tan, W. R.; Chan, C. S.; Aguirre, H.; and Tanaka, K. 2019. Improved ArtGAN for Conditional Synthesis of Natural Image and Artwork. *IEEE Transactions on Image Processing*, 28(1): 394–409.

Voynov, A.; Chu, Q.; Cohen-Or, D.; and Aberman, K. 2023. $P+$: Extended Textual Conditioning in Text-to-Image Generation. *arXiv preprint arXiv:2303.09522*.

Wang, Z.; Zhang, Z.; Zhao, L.; Zuo, Z.; Li, A.; Xing, W.; and Lu, D. 2022. AesUST: towards aesthetic-enhanced universal style transfer. In *Proceedings of the 30th ACM International Conference on Multimedia*, 1095–1106.

Zhang, L.; and Agrawala, M. 2023. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*.

Zhang, Y.; Huang, N.; Tang, F.; Huang, H.; Ma, C.; Dong, W.; and Xu, C. 2023. Inversion-based style transfer with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10146–10156.

# Appendix

## Style and Context Guidance

In this section, we derive how we obtain the style and context guidance depicted in Eq. (3) in the main text.

**Derivation.** Classifier-free guidance (Ho and Salimans 2022) modifies the predicted noise estimation to include the gradient of the log-likelihood of $p(\mathbf{v}|x_t)$, where $\mathbf{v}$ denotes the conditional embedding tokens (of a given text) and $x_t$ is the denoised sample at $t$-th denoising step. Given that $p(\mathbf{v}|x_t) \propto p(x_t|\mathbf{v})/p(x_t)$, it follows that $\nabla_{x_t} \log p(x_t|\mathbf{v}) \propto \nabla_{x_t} \log p(x_t, \mathbf{v}) - \nabla_{x_t} \log p(x_t)$. We parameterize the exact score with the score estimator as $\epsilon_\theta(.) \propto \nabla_{x_t} \log p(.)$, enabling us to derive the classifier-free guidance term as:

$$\hat{\epsilon}_\theta(x_t, \mathbf{v}) = \epsilon_\theta(\emptyset) + \lambda_n[\epsilon_\theta(x_t, \mathbf{v}) - \epsilon_\theta(x_t, \emptyset)]. \quad (4)$$

To derive the proposed style and context guidance, we first decompose text conditional embedding tokens $\mathbf{v}$ into its style and context components as $\mathbf{v} = \mathbf{v_s} \cap \mathbf{v_c}$. As an example, $\mathbf{v}$ is the embedding tokens of inference prompt "A painting of a house in the style of $S^*$", while $\mathbf{v_s}$ and $\mathbf{v_c}$ are the embedding tokens of "A painting of a house" and "in the style of $S^*$" prompts, respectively. Given these partitions, we can rewrite $p(\mathbf{v}|x_t)$ using the terms below:

$$p(\mathbf{v}|x_t) = p(\mathbf{v_s} \cap \mathbf{v_c}|x_t) \quad (5)$$
$$\propto p(\mathbf{v_c}|x_t)p(\mathbf{v_s}|\mathbf{v_c}, x_t) \quad (6)$$
$$\propto p(\mathbf{v_s}|x_t)p(\mathbf{v_c}|\mathbf{v_s}, x_t). \quad (7)$$

From Eq. (6), we derive the following expression:

$$p(\mathbf{v}|x_t) \propto \frac{p(x_t|\mathbf{v_c})}{p(x_t)} \frac{p(x_t|\mathbf{v_s}, \mathbf{v_c})}{p(x_t|\mathbf{v_c})} \quad (8)$$

As in Ho and Salimans (2022), we deduce the above equation to the style and context guidance as in below.

$$\hat{\epsilon}_\theta^1(x_t, \mathbf{v}) = \epsilon_\theta(x_t, \emptyset)$$
$$+ \lambda_c[\epsilon_\theta(x_t, \mathbf{v_c}) - \epsilon_\theta(x_t, \emptyset)] \quad (9)$$
$$+ \lambda_s[\epsilon_\theta(x_t, \mathbf{v}) - \epsilon_\theta(x_t, \mathbf{v_c})]$$

Similarly, using Eq. (7), we derive another guidance term:

$$\hat{\epsilon}_\theta^2(x_t, \mathbf{v}) = \epsilon_\theta(x_t, \emptyset)$$
$$+ \lambda_s[\epsilon_\theta(x_t, \mathbf{v_s}) - \epsilon_\theta(x_t, \emptyset)] \quad (10)$$
$$+ \lambda_c[\epsilon_\theta(x_t, \mathbf{v}) - \epsilon_\theta(x_t, \mathbf{v_s})]$$

Rather than relying solely on either Eq. (9) or Eq. (10), we employ a balanced guidance approach by integrating both forms of guidance as shown in Eq. (3) in the main text. We will elaborate on the benefits of utilizing both terms in the following section. In practice, alongside Eq. (3), we also leverage the conventional classifier-free guidance. Therefore, the final style and context guidance is as follows:

$$\hat{\epsilon}_\theta(x_t, \mathbf{v}) = \epsilon_\theta(x_t, \emptyset)$$
$$+ \lambda_n[\epsilon_\theta(x_t, \mathbf{v}) - \epsilon_\theta(x_t, \emptyset)]$$
$$+ \lambda_c[\epsilon_\theta(x_t, \mathbf{v_c}) - \epsilon_\theta(x_t, \emptyset)]$$
$$+ \lambda_s[\epsilon_\theta(x_t, \mathbf{v}) - \epsilon_\theta(x_t, \mathbf{v_c})] \quad (11)$$
$$+ \lambda_s[\epsilon_\theta(x_t, \mathbf{v_s}) - \epsilon_\theta(x_t, \emptyset)]$$
$$+ \lambda_c[\epsilon_\theta(x_t, \mathbf{v}) - \epsilon_\theta(x_t, \mathbf{v_s})]$$

| Phase | Hyperparameter | Value |
|---|---|---|
| Optimization | Optimization steps | 500 |
| | Learning rate | 0.002 |
| | Batch size | 8 |
| | $T$ (multi-stage TI) | 6 |
| Inference | Inference steps | 25 |
| | Scheduler | DPM |
| | $\lambda_n$ (null guidance) | 5.0 |
| | $\lambda_s$ (style guidance) | 0.5~3.0 |
| | $\lambda_c$ (context guidance) | 3.0 |

Table 4: **Detailed hyperparameters of DreamStyler.**

By utilizing both the proposed and the classifier-free guidances, we can subtly adjust the scaling of style and context guidance ($\lambda_s, \lambda_c$), starting from a zero value. For instance, to minimize the influence of the style guidance on the output images, one merely needs to set $\lambda_s$ to zero without changing other hyperparameters, and similarly, this applies to the context guidance as well. Without the incorporation of classifier-free guidance, users would be exhausted from initiating the guidance scaling search from scratch. We describe the guidance scaling parameters in Table 4.

## Experimental Settings

**Implementation details.** Table 4 demonstrates hyperparameters used for our method. We optimize DreamStyler with Eq. (2) (in the main text), which is a similar training procedure to textual inversion (Gal et al. 2022). We utilize a single A100 GPU for both optimization and inference.

For the style transfer task, we extract a depth map condition from a content image and then encode structure features by adopting ControlNet (Zhang and Agrawala 2023). Although various conditional modalities can be considered, we observed that a depth map is the most suitable choice for style transfer. An edge map tends to strongly preserve the structure of content images so that the structural styles are not effectively conveyed to output images. While one might consider employing a segmentation map to better preserve the structure, this approach can introduce artifacts, especially if the segmentation map is not flawlessly generated.

**Datasets.** To build a dataset used in our model evaluation, we collect most of the artistic paintings from the WikiArt dataset (Tan et al. 2019), while some of the modern art and illustration are from Unsplash[1], licensed free images only. In Table 5, we show inference prompts used in the evaluation.

**Evaluation.** For model evaluation, we employ CLIP-based text and image scores as well as Gram-based style score. The CLIP text score evaluates the alignment of a generated image $I$ with a given text prompt $C$ as in below.

$$\text{TextScore}(I, C) = max(100 * cos(E_I, E_C), 0),$$

---

[1]https://unsplash.com

where, $E_I$ denotes visual CLIP embedding of image $I$ and $E_C$ is textual CLIP embedding of prompt $C$. The CLIP image score measures the alignment of a generated image $I$ and a given style reference image $S$ with the following equation.

$$\text{ImageScore}(I, S) = max(100 * cos(E_I, E_S), 0).$$

The Gram-based style score is a widely used metric in style transfer. It quantifies the textural alignment between two images using Gram features, as below equation.

$$\text{StyleScore}(I, S) = 50 - \frac{1}{L} \sum_{l=1}^{L} \frac{1}{\mathcal{B}} \sum_{\forall (i,j) \in \mathcal{B}} cos(P_I^i, P_S^j)$$

Here, $L$ is the number of layers of VGG16 (Simonyan and Zisserman 2014), which is used to calculate Gram features. $\mathcal{B}$ represents the selected patch pair set and $P^i$ is $i$-th image patch of an image. Notably, we assess the style alignment patch-wise; this technique is more adept at capturing local style features than when analyzing the entire image. Additionally, to align the magnitude of this score with others, we subtract Gram similarity to 50. For an efficient implementation, we choose five patches, each of size 224×224, by cropping from the four corners and the center of the image.

**User study.** We asked 37 participants where they were asked to select the best results based on how closely the outputs adhered to the style of reference images and the context of inference prompts (or content images). Each participant was tasked to vote on 9 questions for text-to-image synthesis (yielding 333 responses in total) and 7 questions for style transfer (yielding 259 responses in total). They were presented with style images, inference prompts, and the results from DreamStyler to other methods. For the style transfer evaluations, content images were displayed in place of the inference prompts. The questionnaires used in the user study are detailed in Table 6. We determined the User score based on the ratio of instances voted as the best.

## Additional Model Analysis

**Training/inference time.** Textual inversion methods typically demand increased training time due to the forward/backward passes in the CLIP text encoder. Nevertheless, the inference time difference with the model optimization-based approach is minimal. In our measurement, with a batch size of one, DreamBooth, TI, and DreamStyler requires 300s, 620s, and 580s, respectively. With 8 batch size, DreamBooth, TI, and DreamStyler takes 60s, 500s, and 480s. Despite the additional time required, DreamStyler proves its worth by delivering superior stylization outcomes as shown in a series of experiments.

**Few-shot text-to-image synthesis.** Figure 13 depicts the change in performance trade-offs for diffusion-based personalization methods when transitioning from a one-shot to a few-shot training regime. For this analysis, all methods are trained using five artistic-style images. Importantly, model optimization-based frameworks, such as DreamBooth and CustomDiffusion, exhibit marked improvements in text scores (as compared to Figure 6, in the main text) due to
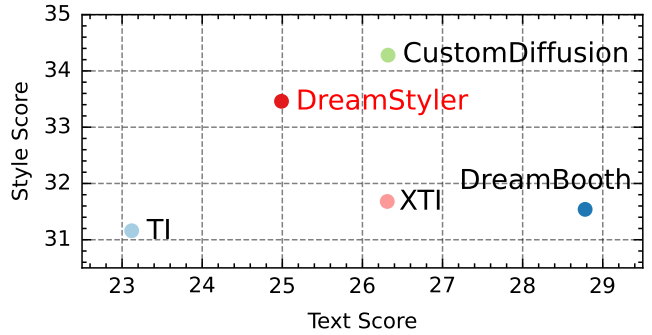


Figure 13: **Performance of few-shot style-guided** text-to-image synthesis. All methods are trained using five artistic style images. In a few-shot training regime, the model optimization-based approaches (DreamBooth, CustomDiffusion) significantly enhance text score. This is because they now can distinguish context from style by referring to multiple style images, leveraging their powerful capacity inherited from learning the model parameters. However, previous textual inversion-based methods (TI, XTI) cannot enjoy using multiple images, showing lower style and text scores. DreamStyler strikes a good balance, exhibiting significant improvement over other textual inversion-based methods.

their capability derived from directly optimizing parameters of denoising UNet. Conversely, textual inversion-based methods, such as TI and XTI, fail to leverage the benefits of multiple images due to the inherent capacity constraint of the embedding optimization nature. Despite DreamStyler employing a textual inversion-based approach, our method not only distinguishes from other textual inversion frameworks but also achieves a trade-off performance comparable to CustomDiffusion, attributed to the timestep-aware textual embedding and context-aware training prompt strategy.

**Style and context guidance.** In Figure 14, we compare three different guidance settings: 1) $\hat{\epsilon}_c(\mathbf{v})$, shown in Eq. (3) in the main text, 2) $\hat{\epsilon}_\theta^1(\mathbf{v})$, shown in Eq. (9), and 3) $\hat{\epsilon}_\theta^2(\mathbf{v})$, shown in Eq. (10). Note that the difference in the outputs ($\epsilon(\mathbf{v})$) relative to the one with null condition ($\epsilon(\emptyset)$) exceeds those of with style-only ($\epsilon(\mathbf{v_s})$) or context-only ($\epsilon(\mathbf{v_c})$). This is because the null conditioned outputs are presumed to produce arbitrary imagery. Hence, when we only rely on $\hat{\epsilon}_\theta^1(\mathbf{v})$ or $\hat{\epsilon}_\theta^2(\mathbf{v})$, the guidance influence between style and context might be imbalanced even with the same scaling parameters, $\lambda_s, \lambda_c$. For instance, with the $\hat{\epsilon}_\theta^1(\mathbf{v})$ guidance, the computation involves differences between $\epsilon(\mathbf{v_c})$ and $\epsilon(\emptyset)$, which amplifies contextual details at the expense of style attributes. On the other hand, $\hat{\epsilon}_\theta^2(\mathbf{v})$ emphasizes style, often overlooks the contextual facets. To achieve a harmonious blend of style and context guidance, we incorporate both guidance forms.

**Ablation study on multi-stage TI.** To demonstrate the necessity of multi-stage TI, we compare our proposed TI with a naive multi-token approach in Figure 15. Naive multi-token TI ($V = 2, 6$) is the method in which only the number of embedding tokens is increased. Our multi-stage TI approach outperforms both StyleScore and TextScore.
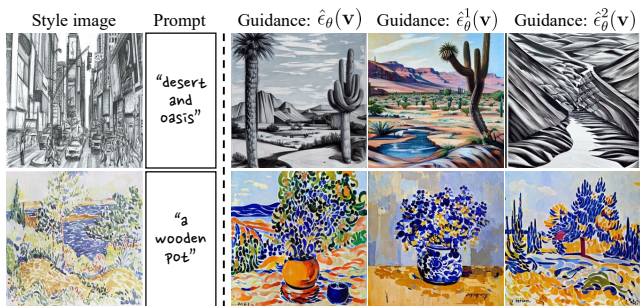
Figure 14: **Comparison of different guidance forms.** Applying $\hat{\epsilon}^1_\theta(\mathbf{v})$ guidance effectively captures the context of a given inference prompt. However, it sometimes fails to adopt the style of a reference image. Conversely, $\hat{\epsilon}^2_\theta(\mathbf{v})$ guidance adeptly aligns the stylistic elements with the reference image but often overlooks certain contexts in the inference prompt. By merging these guidance terms, as depicted in Eq. (3) in the main text, the model achieves a well-balanced expression of both style and context.



Figure 15: **Multi-stage TI ($T$) vs. naive multi-token ($V$).**

**The effectiveness of # T.** We observed a saturation of Style-Score when $T \geq 6$ (Figure 8). This indicates that more embeddings beyond this threshold do not contribute new stylistic nuances and may instead introduce redundancy in artistic expression. For TextScore, we observed an unexpected trend at $T \geq 6$; it starts to generate subjects not specified in the prompt, but presented in style image. We hypothesize that as $T$ increases to a high value, the embeddings try to learn context of style image, beyond stylistic focus; a potential overextension of the embedding capacity. Comparatively, DreamBooth's approach, which involves modifying a large set of parameters, manifests a substantially lower TextScore. Thus, we speculate that increasing $T$ to an extremely high value may yield trends similar to DreamBooth. Nevertheless, regarding TextScore, we believe further investigation is required, and we are grateful for the highlight of this aspect, which had previously eluded our consideration.

**When the style guidance is effective?** In Figure 18, we investigate when the style guidance significantly influences the resulting synthesized image. Our observations indicate that style guidance plays a crucial role in ensuring proper style expression, especially when the output, with zero style guidance ($\lambda_s = 0$), fails to capture the stylistic nuances present in a reference image. Conversely, when the output already effectively expresses the stylistic elements, the influence of style guidance becomes marginal, even with a substantial increase in the scaling parameters $\lambda_s$.

Given these observations, it is now essential to distinguish when the model can accurately replicate style images without the need for style guidance. We hypothesize that style pattern repeatability (Hong et al. 2023) is linked to this ability to some extent. Specifically, when a stylistic feature is complex, the model struggles to capture all its stylistic details. In such cases, the style guidance serves as an effective booster for the model's style representation. Another influencing factor, we conjecture, is the global colortone of the
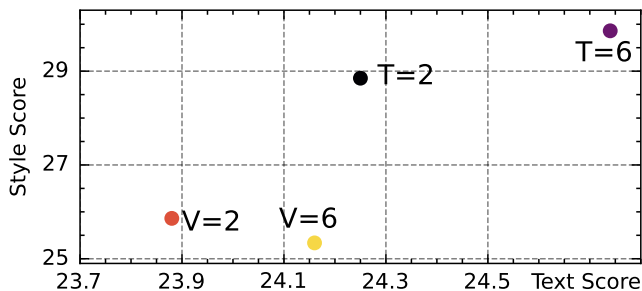
style image. We observed that when there is a pronounced disparity in the colortone of the style image compared to a natural photo, the model often struggles to mimic it, even if the style is not particularly complicated. Despite our observations, determining whether a model can easily adapt to a particular style remains a non-trivial question, and we leave this investigation for future research.

**When the context guidance is effective?** Figure 19 demonstrates scenarios in which the context guidance proves beneficial for artistic image synthesis. We hypothesize that the context guidance is particularly effective when the style image is expressed in a very abstract manner, especially when the structures of the desired subjects to draw deviate significantly from the stylistic expression of the reference image. Conversely, when the style images aptly capture the details of some subjects (as realism paintings do usually), the model can render the desired subject in that style without needing context guidance. However, determining when the personalized diffusion model can adeptly convey the subject with a given style image remains an open question, and we leave this exploration for future work.

**Style transfer** In Figure 16, we illustrate the role of condition modality in the style transfer task. The figure indicates that when the additional condition modality is incorporated via ControlNet (Zhang and Agrawala 2023), the resulting style-transferred images retain a significant degree of fidelity to the original contents' structure. On the other hand, methods that bypass this step and rely solely on the image inversion technique (Meng et al. 2021) often introduce considerable alterations to the initial structure.

**Training progress.** In Figure 20,21, and 22, we compare model optimization- and textual inversion-based methods throughout their training progress. To conduct this, we double the training steps and plot the intermediate results of each method. This allows us to examine the training tendencies and to determine whether the methods fall into overfitting. Through this inspection, we highlight the strengths and weaknesses of these approaches. The model optimization-based approach, DreamBooth (Ruiz et al. 2023), CustomDiffusion (Kumari et al. 2023), exhibits superior style adaptation capability owing to its rich capacity as we train the model directly. However, this approach is prone to overfitting. When overfitting occurs, the model tends to generate images that are very similar to the style images, disregarding all the context from the inference prompt. A critical

Figure 16: **Study on the role of additional conditions in style transfer.** When using additional conditions encoded through ControlNet (Zhang and Agrawala 2023), the outputs faithfully reflect the structure of content images. In contrast, outputs without additional condition and relying solely on image inversion (Meng et al. 2021) considerably changes the structure. This issue is also observed in other inversion-based methods (Zhang et al. 2023; Ahn et al. 2023).
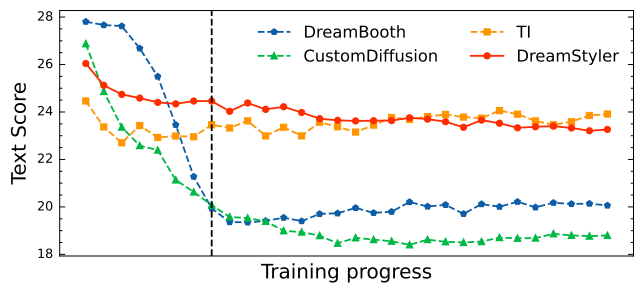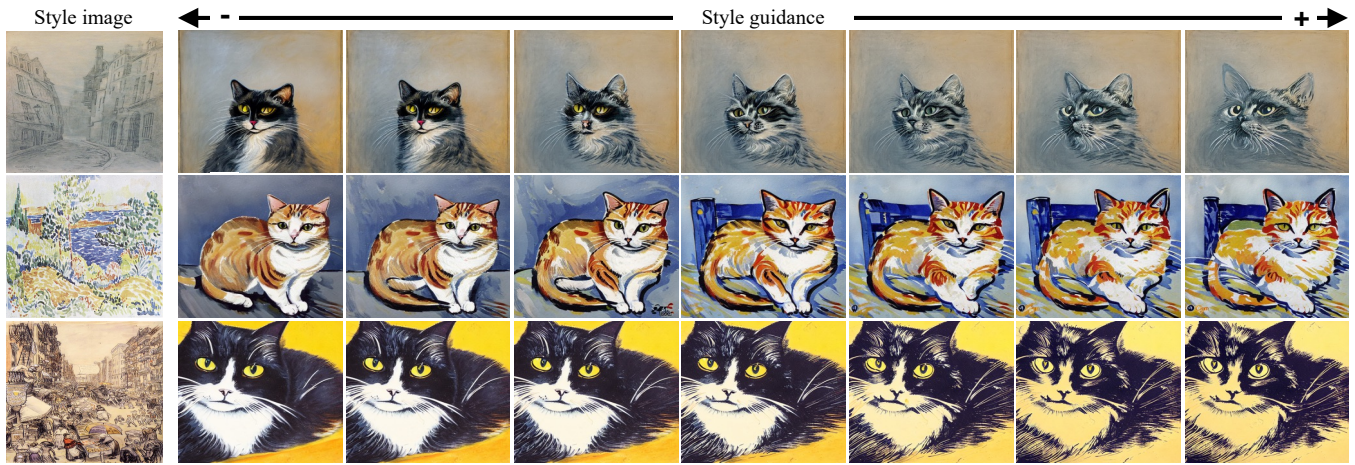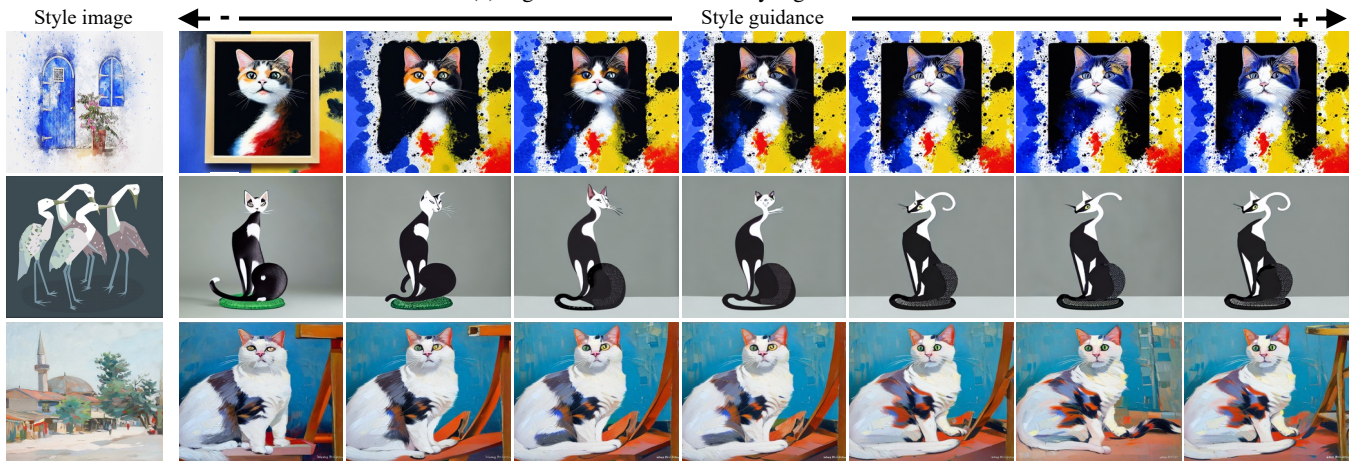


Figure 17: **Training progress.**

additional comparison results on the style transfer task in Figure 24. DreamStyler transfer the style impeccably while preserving the semantic information of content images.

point to note is that the training steps at which overfitting occurs vary significantly to style images. In practice, this is problematic because users cannot select a fixed number of training steps. Instead, they need to inspect all intermediate samples to identify the best checkpoint, which is highly time-consuming and unsustainable in real-world production. Conversely, the textual inversion-based approach (Gal et al. 2022) might not replicate the style image as effectively as model optimization does, but it is less prone to overfitting. This tendency might appear to be symptomatic of underfitting. DreamStyler takes the strengths of both approaches while mitigating their weaknesses. It avoids overfitting in the same way as the textual inversion-based method, yet it also adeptly captures the stylistic nuances found in the style image, akin to the model optimization-based approach. This emphasizes the superiority and practicality of DreamStyler.

In Figure 17, we demonstrate quantitative results of such overfitting phenomena. In our task, overfitting is characterized by the model's neglect of the text prompt and verbatim reproduction of the style image and this can be identifiable through changes in the TextScore. Notably, CustomDiffusion and DreamBooth exhibit a sudden and substantial drop in TextScore at the point when the results are merely copied from style images. On the contrary, DreamStyler performs better in mitigating overfitting than DreamBooth.

## Additional Qualitative Results

In this section, we provide additional qualitative results for better comparisons with previous studies. As shown in Figure 23, DreamStyler successfully generates results with given prompts while depicting the specific style of reference images such as pointillism or watercolor. We also provide
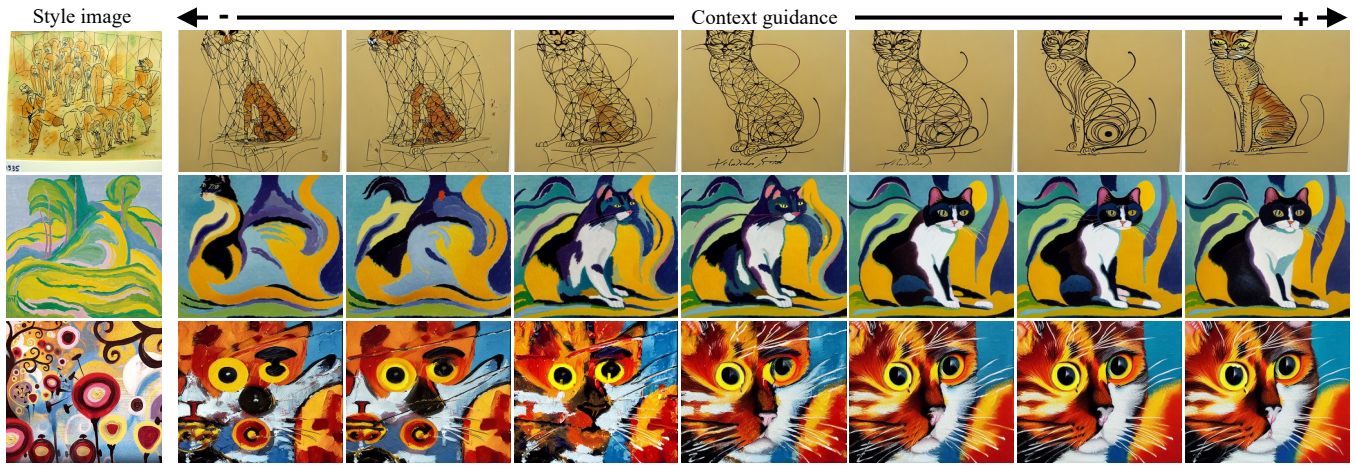
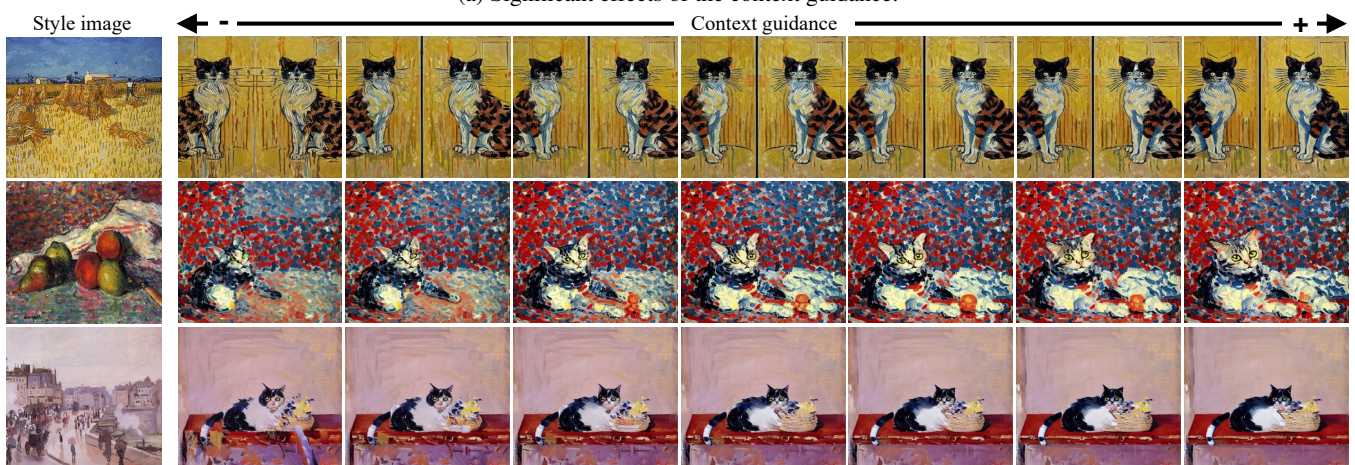(a) Significant effects of the style guidance.



(b) Moderate effects of the style guidance.

Figure 18: **When is the style guidance beneficial?** Inference prompt: "A cat". We investigate the role of the style guidance in stylistic representation. **(a)** Our findings show that the style guidance significantly impacts style expression when generated samples with zero-guidance (leftmost ones) fails to adequately capture the style image. **(b)** In contrast, when the samples without the style guidance effectively embody stylistic nuances, the introduction of the style guidance has minimal effect. We hypothesize that the style guidance is particularly beneficial for images with high pattern complexity (as shown in (a)), so that the model cannot easily adapt to. In cases of low pattern complexity (as shown in (b)), its influence appears to be marginal.

(a) Significant effects of the context guidance.



(b) Moderate effects of the context guidance.

Figure 19: **When is the context guidance beneficial?** Inference prompt: "A cat". We investigate the role of the context guidance in the context (subject) production aspect. **(a)** Our findings show that the context guidance substantially impacts to appearance of subjects when generated samples with zero-guidance (leftmost ones) exhibit abstract subject structures. **(b)** In contrast, when subjects already appear detailed even without the context guidance, the introduction of the context guidance has minimal effect. We hypothesize that the context guidance is particularly beneficial for images with a high degree of structural abstraction (as shown in (a)). In contrast, for the images that present detailed textures, its influence appears to be marginal.
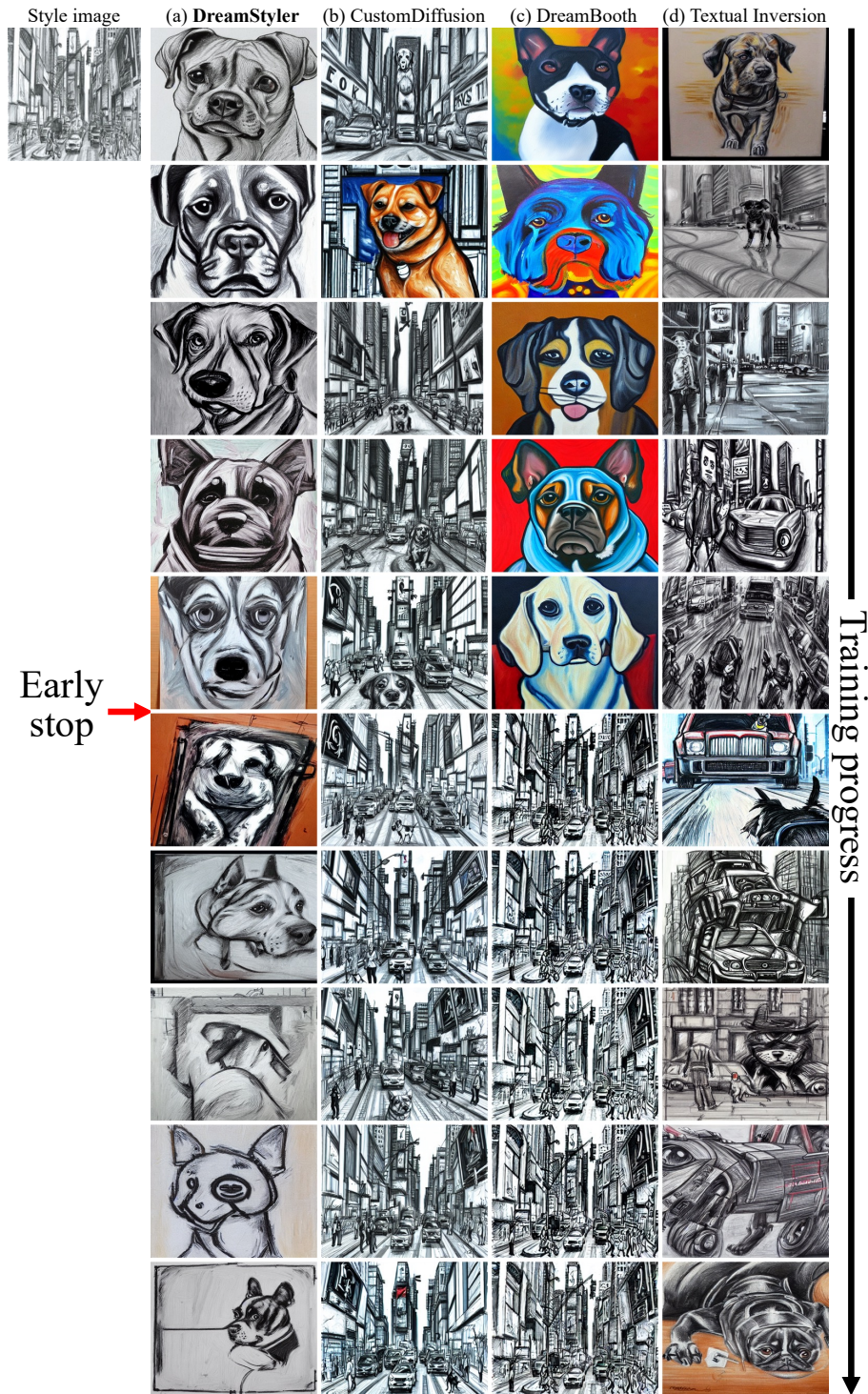
Figure 20: **Comparison of training progress.** Inference prompt: "A painting of a dog". We double the training steps for each model and visualize their training tendencies. "Early stop" refers to the original (official) training step which we also used in the paper. Model optimization-based methods tend to suffer from overfitting to the style image, and we observed that the occurrence of overfitting varies depending on the style image. Notably, DreamBooth often fails to reflect the style image initially and then abruptly overfits to it. Since it's challenging to predict when overfitting will begin, users would need to monitor all intermediate training samples, which is impractical for production-level tasks. While textual inversion-based methods (including ours) are immune to overfitting, the vanilla TI often fails to accurately capture both the style images and the inference prompt. In contrast, DreamStyler avoids overfitting and consistently demonstrates high-quality stylization in terms of both style and context.
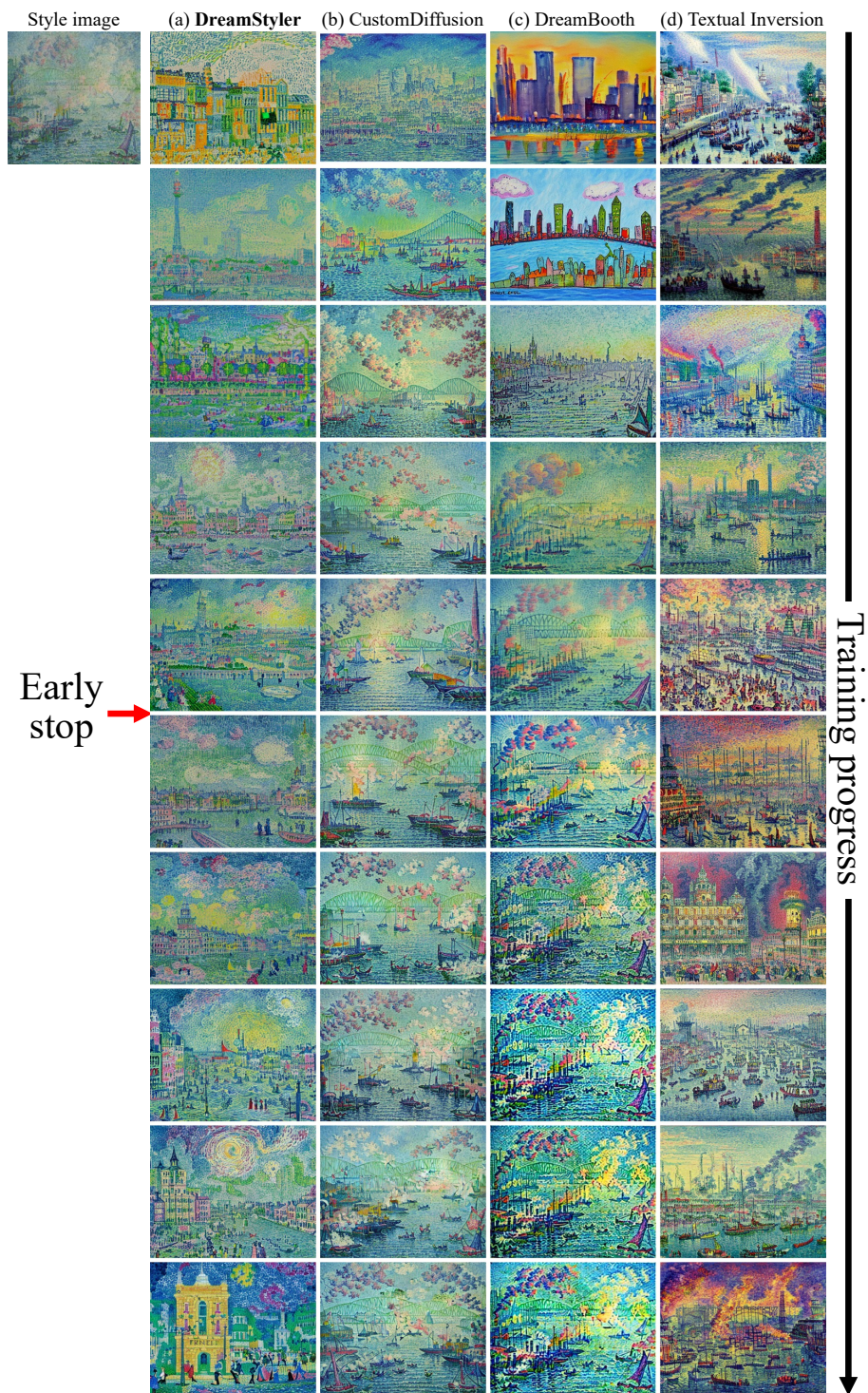
Figure 21: **Comparison of training progress.** Inference prompt: "A painting of a house".
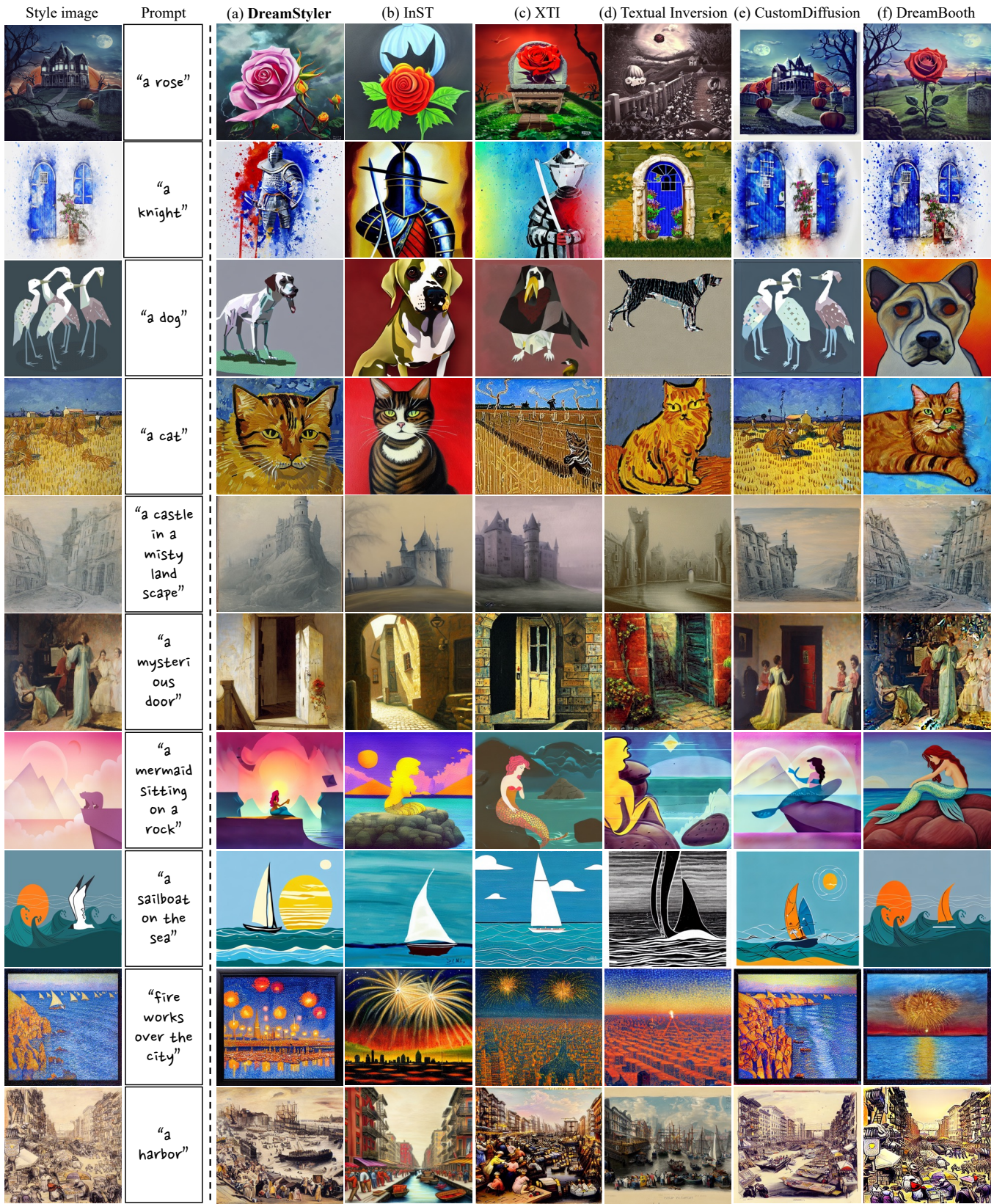
Figure 22: **Comparison of training progress.** Inference prompt: A painting of a city".

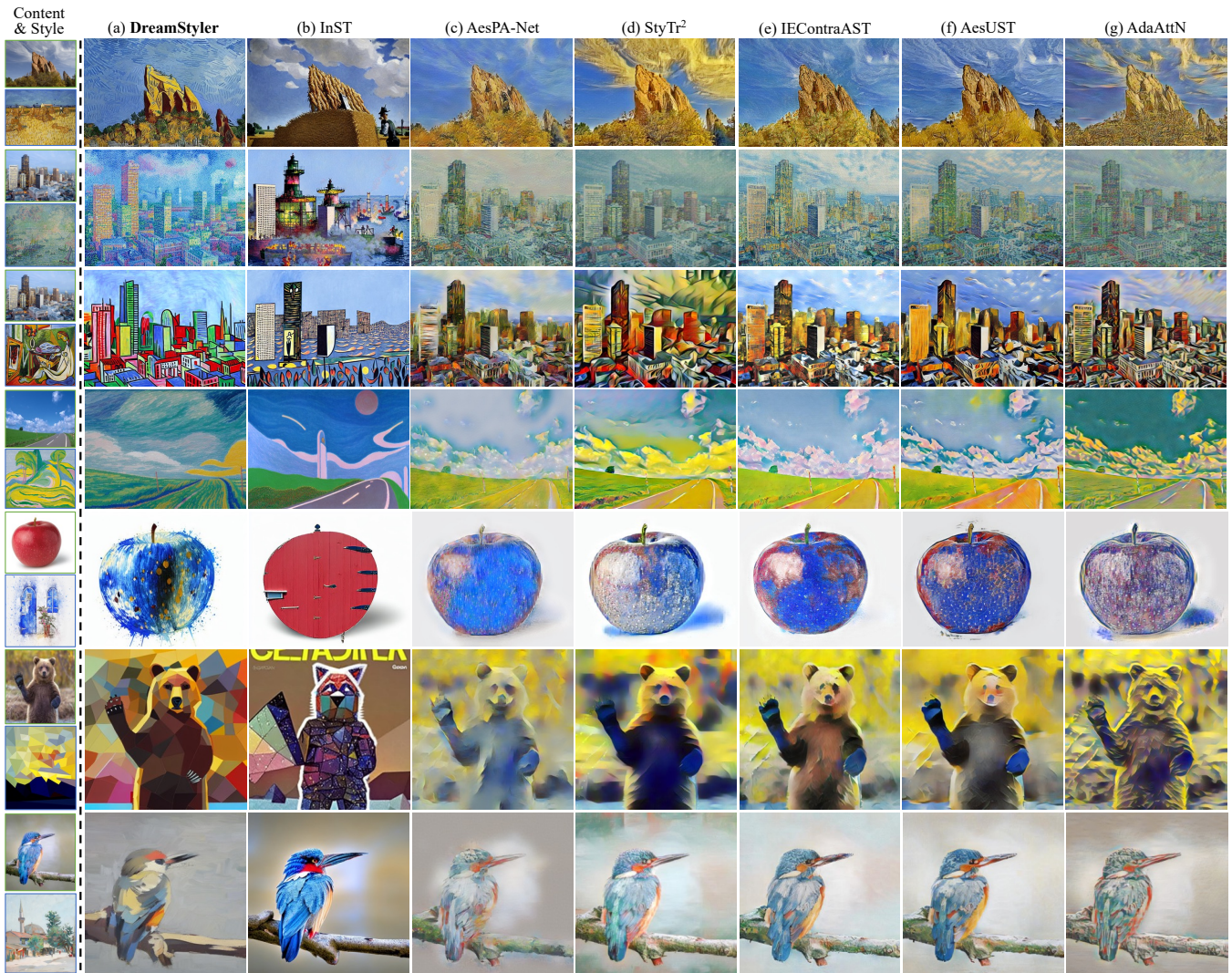Figure 23: **Additional qualitative comparison** on the style-guided text-to-image synthesis task.

Figure 24: **Additional qualitative comparison** on the style transfer task.

| Evaluation prompt |
| --- |
| A painting of a house in the style of $S^*$ |
| A painting of a dog in the style of $S^*$ |
| A painting of a robot in the style of $S^*$ |
| A painting of a crying woman in the style of $S^*$ |
| A painting of a wooden pot in the style of $S^*$ |
| A painting of an ocean on a cloudy day in the style of $S^*$ |
| A painting of a city in the style of $S^*$ |
| A painting of a full moon in the mountains in the style of $S^*$ |
| A painting of a temple in the forest in the style of $S^*$ |
| A painting of a cat in the style of $S^*$ |
| A painting of a rose in the style of $S^*$ |
| A painting of a man with a bearded face in the style of $S^*$ |
| A painting of a solar system in the style of $S^*$ |
| A painting of a spaceship in the style of $S^*$ |
| A painting of a highway in the style of $S^*$ |
| A painting of a bridge in the style of $S^*$ |
| A painting of an apple and a banana in the style of $S^*$ |
| A painting of fireworks over the city in the style of $S^*$ |
| A painting of a harbor in the style of $S^*$ |
| A painting of a rock island and an ocean in the style of $S^*$ |
| A painting of a blossoming cherry tree in the style of $S^*$ |
| A painting of a snowy mountain peak in the style of $S^*$ |
| A painting of a knight in the style of $S^*$ |
| A painting of a marketplace in the style of $S^*$ |
| A painting of a desert and an oasis in the style of $S^*$ |
| A painting of elephants at sunset in the style of $S^*$ |
| A painting of a crowd under the stars in the style of $S^*$ |
| A painting of a lighthouse on a cliff in the style of $S^*$ |
| A painting of a mermaid sitting on a rock in the style of $S^*$ |
| A painting of a haunted house on a hill in the style of $S^*$ |
| A painting of a cafe in the morning in the style of $S^*$ |
| A painting of a castle in a misty landscape in the style of $S^*$ |
| A painting of birds taking flight in the style of $S^*$ |
| A painting of a family in the style of $S^*$ |
| A painting of a waterfall in a forest in the style of $S^*$ |
| A painting of a night of shooting stars in the style of $S^*$ |
| A painting of a sailboat on the sea in the style of $S^*$ |
| A painting of a girl under a tree in the style of $S^*$ |
| A painting of a fisherman casting his net at sunrise in the style of $S^*$ |

Table 5: **Inference prompts used in the model evaluation.**

| Task | Instruction and Question |
|---|---|
| Text-to-Image | **Instruction:**<br>Each question provides 1) a style image, 2) a text prompt, and 3) 6 generated samples.<br>Samples are results of synthesizing context of the text prompt to the style depicted in the style image.<br>Please select the sample that best represents both the text prompt and the style of the style image.<br><br>**Question:**<br>Which sample best captures the style of the style image and the content of the text? |
| Style Transfer | **Instruction:**<br>Each question provides 1) a style image, 2) a content image, and 3) 7 generated samples.<br>Samples are results of adapting the content image to the style depicted in the style image.<br>Please select the sample that best resembles with the artistic style of the artists' style image,<br>as if that artist used the content image as a reference in their painting.<br><br>**Question:**<br>If the artist of the style image had used the content image as a reference,<br>which sample best embodies their artistic style? |

Table 6: **Questionnaires used in the user study.**