

ArtCrafter: Text-Image Aligning Artistic Attribute Transfer via Embedding Reframing

Nisha Huang, Kaer Huang, Yifan Pu, Jiangshan Wang, Jie Guo, Yiqiang Yan,
Xiu Li, *Member, IEEE*, Tong-Yee Lee, *Senior Member, IEEE*

Abstract—Recent years have witnessed significant advancements in text-guided style transfer, primarily attributed to innovations in diffusion models. These models excel in conditional guidance, utilizing text or images to direct the sampling process. Traditional style transfer focuses on low-level visual features, such as brushstroke textures and color distributions, and appears more like applying an artistic filter to an image. Artistic attribute transfer, however, transcends the limitations of traditional style transfer by achieving the transfer of visual concepts from color and brushstrokes to high level aesthetic attributes such as composition, pose, and key semantic elements, resulting in more natural outcomes. Therefore, we propose an innovative text-to-image artistic attribute transfer framework named ArtCrafter. Specifically, we introduce an attention-based style extraction module, meticulously engineered to capture the subtle artistic attribute elements within an image. This module features a multi-layer architecture that leverages the capabilities of perceiver attention mechanisms to integrate fine-grained information. Additionally, we present a novel text-image aligning augmentation component that adeptly balances control over both modalities, enabling the model to efficiently map image and text embeddings into a shared feature space. We achieve this through attention operations that enable smooth information flow between modalities. Lastly, we incorporate an explicit modulation that seamlessly blends multimodal enhanced embeddings with original embeddings through an embedding reframing design, empowering the model to generate diverse outputs. Extensive experiments demonstrate that ArtCrafter yields impressive results in visual stylization, exhibiting exceptional levels of artistic attribute intensity, controllability, and diversity. The code is available at <https://github.com/haha-lisa/ArtCrafter>.

Index Terms—Diffusion models, text-image alignment, artistic attribute transfer

1 INTRODUCTION

Diffusion-based text-to-image generation models [1], [2] have made significant strides in the areas of personalization and customization, particularly in consistent synthesis tasks such as identity protection [3], [4], object customization [5], [6], and style transfer [7]–[13]. Among these applications, text-guided style transfer has emerged as a powerful tool, focusing on fine-grained style representation that captures abstract concepts like texture, color, composition, brushstroke, and genre. This approach enables the creation of a diverse range of personalized outputs that are deeply rooted in the semantic essence of the input text.

Current methods for stylization tasks often leverage pre-trained diffusion models [14]–[22], enhancing model features by adding a trainable adapter module without full retraining. In text-to-image style transfer applications, adapter-based methods shape the style and content of the output by adjusting the condition guidance scales over the

input image and text prompts. However, we have identified three main issues with previous research: **1) Inadequate representation of artistic attributes.** Traditional style transfer settings and image encoder architectures limit their ability to understand artistic attributes beyond textures and colors, such as composition, pose, and key semantic elements. **2) Suboptimal text-guided conditions.** As shown in Fig. 2, the usual adapter-based method [16] fails to deliver the expected results when using the text condition “Fashion shoes”. This discrepancy arises because the amount of information in image and text embeddings is not equivalent, yet adapter-based methods [14]–[20] often directly concatenate the two without addressing the imbalance and disparity, leading to image data overshadowing text prompts during the sampling process. **3) Lack of output diversity.** The constrained liberation of textual guidance results in generated outputs that closely resemble the reference images, thereby limiting the diversity of the results.

To address the aforementioned challenges, we introduce ArtCrafter, a novel embedding reframing solution based on diffusion models, specifically tailored for text-guided stylization tasks. ArtCrafter comprises three key components: **1) Attention-based Artistic Attribute Extraction** (Sec. 3.2): This component leverages multi-level features to capture aesthetic attributes at various levels, ensuring a more coherent and comprehensive encoding of artistic features. It employs a non-hierarchical refinement module and a multi-layer attention architecture to capture both local and global stylistic elements. Additionally, we introduce the ArtMarket dataset, which pairs art images with descriptive texts, enabling us to fine-tune an encoder initially trained

The work is supported in part by the Shenzhen Key Laboratory of Next Generation Interactive Media Innovative Technology, China (No. ZDSYS20210623092001004), and the National Science and Technology Council, Taiwan (No. 114-2221-E-006-114-MY3). (Co-corresponding authors: Xiu Li and Tong-Yee Lee.)

- N. Huang, Y. Pu, J. Wang, and X. Li are with Tsinghua International Graduate School, Tsinghua University, Shenzhen 518071, China. E-mail: {hns24, puylf23, wjs23}@mails.tsinghua.edu.cn and li.xiu@sz.tsinghua.edu.cn.
- N. Huang and G. Jie are also with Pengcheng Laboratory, Shenzhen 518055, China. E-mail: {huangnsh, guoj01}@pcl.ac.cn
- K. Huang and Y. Yan are with Lenovo, Inc. E-mail: {huangke1, yangq}@lenovo.com.
- T.-Y. Lee is with National Cheng Kung University, Tainan 701, Taiwan. E-mail: tonylee@mail.ncku.edu.tw.

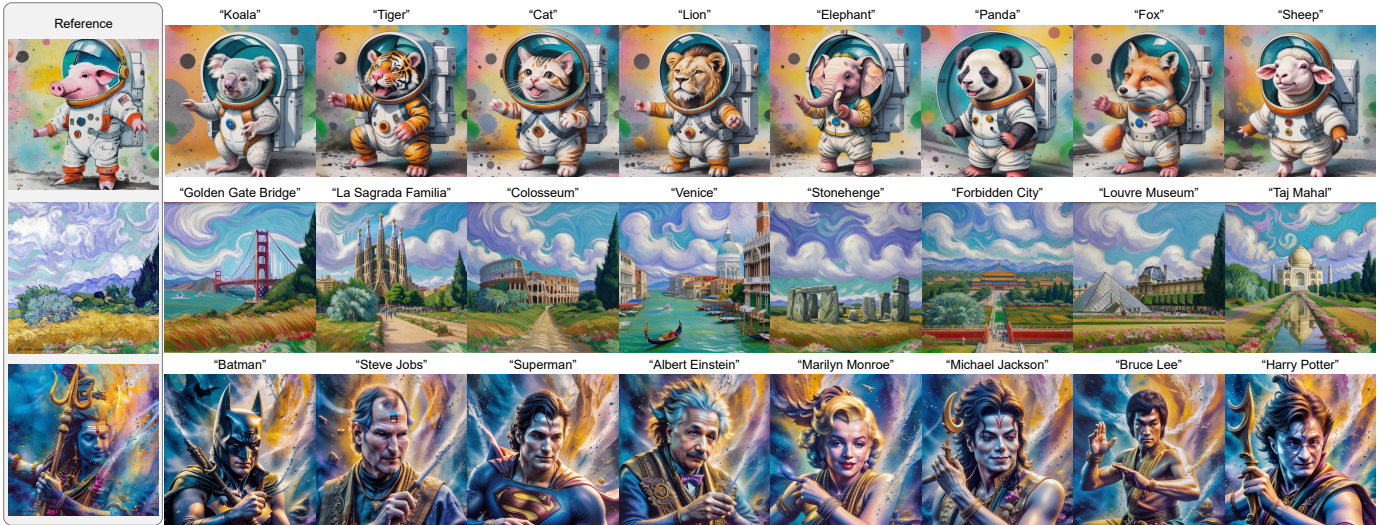


Fig. 1. **ArtCrafter generation results.** By injecting the features of the reference images and text prompts during the diffusion process, our method is capable of capturing and generating a faithful style representation.

on natural images. This approach successfully maintains strong generalization capabilities while being sensitive to the unique visual elements of artistic paintings. **2) Text-Image Aligning Augmentation** (Sec. 3.3): This module enhances the alignment of image and text embeddings through crafted attentional interactions, mapping them into a shared feature space. This alignment ensures that the generated images reflect both the style of the reference image and the content of the textual conditions, thereby enhancing controllability. **3) Explicit Modulation** (Sec. 3.4): This component enhances the adaptability of conventional fusion techniques through the implementation of linear interpolation and concatenation schemes. This method facilitates the merging of original image and text embeddings with multimodal embeddings, yielding enhanced embeddings. It allows ArtCrafter to generate images that are relevant to the text prompts and exhibit diverse visual representations. Our main contributions are summarized as follows:

- We introduce ArtCrafter, a lightweight adapter designed to enhance the capabilities of pre-trained diffusion models in text-guided stylized image generation. Our approach focuses on attention-based style feature extraction, effectively capturing both local and global features.
- We propose an innovative text-image aligning augmentation module that enables robust interaction between reference images and textual descriptions within a shared feature space, significantly enhancing the influence of text prompts on the generative process.
- The explicit modulation within ArtCrafter optimizes the utilization of multimodal embeddings, offering greater flexibility and diversity than conventional methods. Additionally, ArtCrafter is compatible with additional control conditions and achieves superior performance across various experimental benchmarks compared to state-of-the-art approaches.



Fig. 2. **Generic adapter-based vs. ArtCrafter generation results.** Given a content description of “Fashion Shoes”, generic adapter-based generation (above) results in unaligned results and limited result diversity. In contrast, our approach (below) generates text-aligned content as well as multiple shoe types.

2 RELATED WORK

2.1 Attention Control in Diffusion Models

Following the remarkable progress made in the field of pre-training text-to-image diffusion models [1], [2], [23], [24], a series of image editing efforts [25]–[28] have emerged. Hertz *et al.* [29] propose the Prompt-to-Prompt method, which achieves text-based partial image editing and generates edited images that conform to textual conditions by replacing original vocabulary and cross-attention maps. Plug-and-play [30] utilizes the spatial features and self-attention mapping of the original image to guide the diffusion model for text-guided image-to-image translation while preserving the spatial layout of the original image. Later, MasaCtrl [31] proposes a mutual self-attention control technique for coherent image editing. Consistent image editing is achieved by preserving the key and value of the self-attention layer of the source image while conditioning the model with desired text prompts. Recently, StyleID [7] proposes a style migration method without training by injecting styles in the self-attention layer and introduces query preservation, attention temperature scaling, and initial latent AdaIN techniques to minimize the impact of style injection on the original

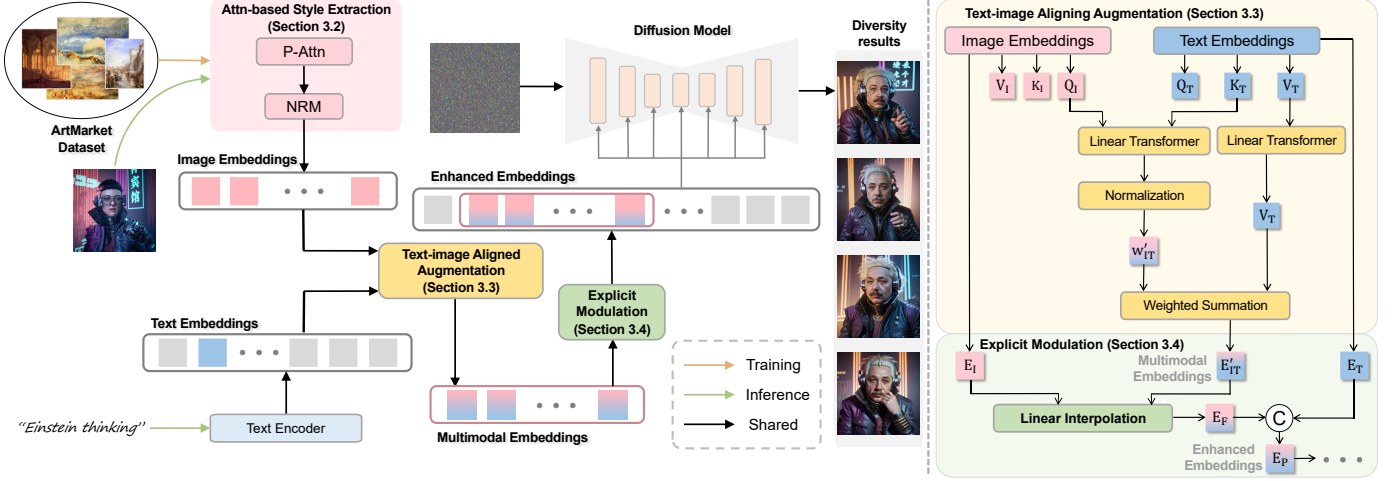


Fig. 3. Our proposed ArtCrafter’s overview pipeline consists of three modules: 1) Attention-based style extraction captures fine-grained style features from images using multi-head perceiver attention and feed-forward networks. 2) Text-image aligning augmentation transfers the style embedding to the textual space to improve the image-text consistency in the generated images. 3) Explicit modulation seamlessly combines multimodal embeddings with original image embeddings and text embeddings in different ways, increasing the versatility and diversity of the method.

content. \mathcal{Z}^* [10] shows how to extract style information directly from style images and fuse it with content images without training by means of an attentional reweighting strategy. Unlike the methods mentioned above, our approach focuses on the interaction of image and text information, as well as a balance for guiding the generation process.

2.2 Text-Guided Style Transfer

Stylization [32]–[35] is to control the content through text and make the generated image consistent while keeping the style of the reference image. StyleDrop [36] is a method of achieving arbitrary style synthesis from a small number of stylized images and textual descriptions by efficiently fine-tuning a small number of parameters of a pre-trained model and combining iterative training with feedback to improve the quality of the generated images. The well-received work IP-Adapter [16] improves stylization image generation by introducing the embedding of input images in an additional layer of cross-attention, which enhances the model’s ability to capture features from the input images. Building on the IP-Adapter, InstantStyle [18] manually selects specific attention layers to control the style of the output. However, for the IP-Adapter conditioned on natural images, the expected conditioning of the input artistic image does not always work. Visual Style Prompting (VSP) [37] captures stylistic details by fusing key features of the reference image in a later self-attention layer while preserving the content consistency of the original image. However, compared to cross-attention, self-attention provides a weaker ability to control the semantic content of the generated image. Style Aligned [38] attempts to align style and content through a shared attention mechanism. However, it generates results with content information leaked from the style image, and there are challenges in disentangling content and style. StyleShot [17] is trained by a two-stage style control method. However, detailed information within the control is easily lost due to sparse rows and columns. Furthermore, [17], [14], and [27] only trained on the art-text datasets, making it difficult to broadly adapt to arbitrary stylistic features.

3 METHOD

The overall architecture of ArtCrafter is shown in Fig. 3. As reviewed in Sec. 3.1, ArtCrafter is built upon diffusion model [1], [39]. In Section 3.2, we introduce attention-based style extraction, which captures multi-level style information by non-layer refinement module and multi-layer design. Text-image aligning augmentation (Section 3.3) allows the model to dynamically weigh the importance of different parts of the text prompt. This enables a more nuanced and context-aware generation process, resulting in images that are more closely related to the text prompt. Explicit modulation (Section 3.4) effectively combines textual and visual information, enabling the model to generate images that are relevant to the text prompt and have diverse visual representations. In Section 3.5, we provide a detailed description of the training and inference processes. During training, we optimize the adapter while keeping the parameters of the pre-trained diffusion model frozen. The adapter is trained on the ArtMarket dataset. During inference, we enhance controllability through classifier-free guidance with dual conditioning, which enables the model to generate images more aligned with the text prompt while retaining the stylistic features of the art image.

3.1 Preliminary

The forward diffusion process incrementally adds Gaussian noise ϵ to the data x_0 through a Markov chain. Specifically, the data x_0 is gradually corrupted by noise over T timesteps, resulting in a sequence of progressively noisier data points x_0, x_1, \dots, x_T . Each step t in the Markov chain is defined by:

$$x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon, \quad (1)$$

where β_t is a predefined noise schedule that controls the amount of noise added at each step, and $\epsilon \sim \mathcal{N}(0, 1)$ is Gaussian noise. By the final step T , the data x_T is typically indistinguishable from pure Gaussian noise.

The reverse denoising process aims to reconstruct the original data from the noisy data x_T . This process is driven

by a learnable denoising model $\epsilon_\theta(x_t, t, c)$ parameterized by θ . The denoising model is implemented using a U-Net architecture [40], which is capable of capturing complex patterns and dependencies in the data. The denoising process starts from $x_T \sim \mathcal{N}(0, 1)$ and iteratively refines the data by predicting and removing the noise at each timestep t :

$$\hat{x}_{t-1} = \sqrt{1 - \beta_t} \hat{x}_t - \sqrt{\beta_t} \epsilon_\theta(\hat{x}_t, t, c). \quad (2)$$

The denoising model $\epsilon_\theta(\cdot)$ is trained to predict the noise ϵ added at each step, allowing the model to progressively reconstruct the original data x_0 .

The denoising model $\epsilon_\theta(\cdot)$ is trained with a mean-squared loss derived from a simplified variant of the variational bound:

$$\mathcal{L} = \mathbb{E}_{t, x_0, \epsilon} [\|\epsilon - \hat{\epsilon}_\theta(x_t, t, c)\|^2], \quad (3)$$

where c denotes an optional condition. In the diffusion model, c is generally represented by the text embeddings E_T encoded from a text prompt using CLIP [41]. These text embeddings are integrated into the diffusion model, allowing the model to generate samples conditioned on the provided text prompt. This conditioning mechanism enables the diffusion model to produce outputs that are semantically aligned with the text description, making it a powerful tool for text-to-image synthesis and other conditional generation tasks.

3.2 Attention-based Style Extraction (ASE)

Our ASE method innovates in the extraction of artistic features. Traditional encoders (such as CLIP) are pre-trained on natural images, and their representation spaces are suboptimal for extracting artistic styles. Aesthetic attributes are distributed both in the global composition and local brushstrokes of an image, necessitating a mechanism that can distill discriminative artistic factors from redundant pixel information. ASE enhances the encoding capability of artistic attributes by integrating fine-grained features through a multi-layer architecture. Unlike the methods in [42], [43], which focus on general feature extraction, our ASE method specifically targets intricate artistic attributes from images by leveraging Perceiver Attention [42] (P-Attn) and a Non-Linear Refinement Module (NRM). Our technical contributions lie in the novel combination and adaptation of these components to artistic attribute feature extraction, which significantly enhances the representation of attribute features. The method involves several key steps: initializing latent variables, expanding them to match the input batch size, applying P-Attn, and using the NRM to refine the latent variables. The method involves several key steps: initializing latent variables, expanding them to match the input batch size, applying P-Attn, and using the NRM to refine the latent variables.

Given the reference image, we obtain the input image embeddings through CLIP, denoted as x . The latent variables z are initialized as a tensor with a shape of $(1, N, D)$, where N is the number of queries and D is the dimension of the latent space. To stabilize the training process, the latent variables are normalized by dividing by the square root of D :

$$z = \frac{\mathcal{N}(0, 1)}{D^{0.5}}. \quad (4)$$

To match the batch size of the input x , we expand the latent variable z by repeating it along the batch dimension. This process can be represented as:

$$z = z \otimes \mathbf{1}_B, \quad (5)$$

where $\mathbf{1}_B$ is a tensor of ones with shape $(B, 1, 1)$, and B is the batch size of x . This operation results in a tensor z with shape (B, N, D) , where N is the number of queries and D is the dimension of the latent space.

The P-Attn mechanism denoted as P-Attn, is then applied to update the latent variables by attending to the input x and the repeated latent variables:

$$\text{P-Attn}(x, z) = \text{softmax}\left(\frac{zx^T}{\sqrt{d_k}}\right) \cdot x, \quad (6)$$

$$z' = \text{P-Attn}(x, z) + z, \quad (7)$$

where d_k is the dimension of the key tensor, typically equal to D . This operation allows the model to selectively focus on different parts of the input data based on the learnable latent variables, thereby capturing intricate style details from the reference image. Compared to [42], [43], our use of P-Attn in the context of style extraction is novel and specifically tailored to enhance the representation of stylistic nuances.

The Non-Linear Refinement Module (NRM) consists of two linear transformations with a GELU activation function in between:

$$\text{NRM}(z') = W_2 \cdot \text{GELU}(W_1 \cdot z' + b_1) + b_2, \quad (8)$$

where W_1 and W_2 are weight matrices, and b_1 and b_2 are bias terms. The NRM further refines the latent variables by applying non-linear transformations, enhancing the representation of the style features. In contrast to [42], [43], our NRM is designed to specifically refine style-related latent variables, resulting in a more detailed and robust style embedding.

The output E_I represents the style embeddings extracted from the input image, obtained by combining the NRM output and the updated latent variables z' :

$$E_I = \text{NRM}(z') + z'. \quad (9)$$

This final step integrates the refined latent variables with the non-linear transformations from the NRM, resulting in a robust and detailed style embedding that captures the intricate style details from the reference image. This style embedding can be used to guide the generation process in downstream applications, ensuring that the generated outputs inherit the desired stylistic features.

The output E_I represents the artistic attribute feature embedding extracted from the input image, obtained by combining the output of the NRM and the updated latent variable z' :

$$E_I = \text{NRM}(z') + z'. \quad (10)$$

This final step integrates the refined latent variables with the non-linear transformations from the NRM, resulting in a robust and detailed feature embedding that captures the complex artistic attributes from the reference image. This embedding can be used to guide the generation process in downstream applications, ensuring that the generated outputs inherit the desired artistic features. Experimental evaluation in Fig. 4 and Table 1 and 2 verify our success in guiding artistic attribute transfer.

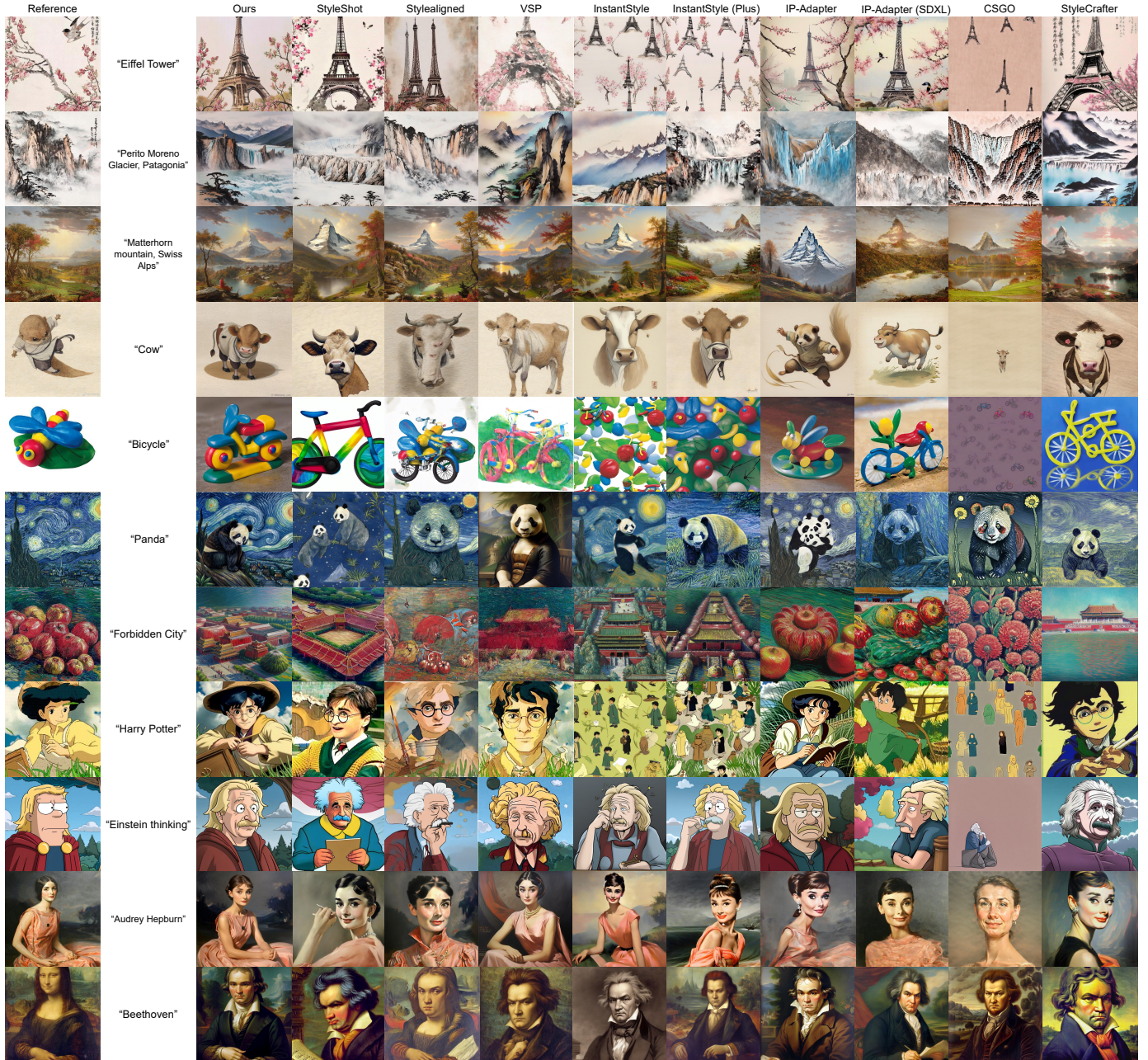


Fig. 4. **Qualitative Comparison with Other State-of-the-Art Text-Guided Stylization Methods.** We conducted a comprehensive qualitative evaluation by comparing our method with various state-of-the-art text-guided stylization methods, including StyleShot [17], Style Aligned [38], VSP [37], InstantStyle [18], InstantStyle (Plus) [18], IP-Adapter [16], IP-Adapter (SDXL) [16], CSGO [27], and StyleCrafter [14].

3.3 Text-Image Aligning Augmentation (TIAA)

The TIAA method focuses on addressing the issue that textual conditions are not overwhelmed by image information during the generation process. In traditional methods, directly concatenating image embeddings and text embeddings leads to modality imbalance, usually manifested as image information dominating the generation process while textual guidance is neglected. TIAA allows the model to integrate image and text embeddings more effectively, projecting them into a shared feature space where their interactions can be more nuanced. Our novel process involves three main steps: linear transformation, attention-weight calculation, and attention-weighted value matrix computation.

We start by transforming the image prompt embeddings E_I and the text prompt embeddings E_T into query, key, and value matrices through linear layers. These transformations are represented as:

$$Q_I = W_{Q_I} E_I, \quad K_T = W_{K_T} E_T, \quad V_T = W_{V_T} E_T. \quad (11)$$

Here, W_{Q_I} , W_{K_T} , and W_{V_T} are the weight matrices associated with the query for images, key for text, and value for text, respectively. These linear transformations map the embeddings into a shared feature space, enabling the subsequent attention mechanism to effectively capture the interactions between the image and text representations.

The attention weights w_{IT} are calculated by taking the dot product of the query matrix Q_I and the key matrix K_T ,

scaled by the square root of the key dimension d_{k_T} to prevent gradient disappearance:

$$\mathbf{w}_{IT} = \frac{Q_I K_T^T}{\sqrt{d_{k_T}}}. \quad (12)$$

The softmax function is then applied to these raw attention scores to obtain the normalized attention weights \mathbf{w}'_{IT} :

$$\mathbf{w}'_{IT} = \text{softmax}(\mathbf{w}_{IT}). \quad (13)$$

The normalized attention weights \mathbf{w}'_{IT} represent the importance of each text feature with respect to the image features, allowing the model to dynamically focus on the most relevant parts of the text prompt.

Using the normalized attention weights \mathbf{w}'_{IT} , we compute the weighted sum of the value matrix V_T to generate the multimodal embeddings E'_{IT} :

$$E'_{IT} = \mathbf{w}'_{IT} \cdot V_T. \quad (14)$$

The resulting multimodal embeddings E'_{IT} capture the multimodal context more effectively, allowing the model to generate images that are more closely aligned with the semantic content of the text prompt. From our experiments, this approach is particularly beneficial in scenarios where the textual and visual information needs to be tightly integrated to produce coherent outputs. By leveraging cross-attention mechanisms, the model can dynamically prioritize different aspects of the text prompt, leading to more accurate and contextually relevant image generation.

3.4 Explicit Modulation (EM)

EM provides an explicit control knob for output diversity. Traditional fusion methods are rigid and lack the ability to finely adjust the strength of style control and output diversity. Traditional stylization methods lack flexibility. For example, the adapter-based method in Fig. 2 generates similar results, failing to produce diverse outputs. In contrast, the EM method integrates image embeddings with multimodal embeddings through linear interpolation, flexibly balancing both influences. This significantly enhances the diversity of the output.

Specifically, we fuse the image embeddings E_I with the multimodal embeddings E'_{IT} through linear interpolation:

$$E_F = \alpha E_I + (1 - \alpha) E'_{IT}, \quad (15)$$

where α is a predefined constant controlling the fusion ratio between the original and enhanced embeddings. The value of α can be adjusted to balance the contribution of the image embeddings and the multimodal embeddings, ensuring that the fused embeddings E_F retain the essential features of both modalities.

Ultimately, we concatenate the fused image embeddings E_F with the text prompt embeddings E_T to form the complete prompt embeddings for image generation:

$$E_P = E_T \oplus E_F, \quad (16)$$

where \oplus denotes the concatenation operation. The resulting prompt embeddings E_P integrate the enhanced multimodal information into the diffusion model. By carefully balancing the contributions of the text and image embeddings, the model gains a robust and controlled representation that effectively captures multimodal conditions, thereby improving the overall generation performance.

3.5 Training and Inference

During training, we optimize the adapter while keeping the parameters of the pre-trained diffusion model frozen. The adapter is trained on the ArtMarket dataset, using the same loss function as the original Stable Diffusion (SD):

$$\mathcal{L} = \mathbb{E}_{z_0, \epsilon, c_t, c_{art}, t} \|\epsilon - \epsilon_\theta(z_t, c_t, c_{art}, t)\|^2. \quad (17)$$

Here, ϵ is the randomly sampled Gaussian noise, ϵ_θ is the noise prediction model, and t is the timestep. Note that during training, the latent variable z is constructed with the art image c_{art} as follows:

$$z_t = \sqrt{\bar{\alpha}_t} \psi(c_{art}) + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad (18)$$

where $\psi(\cdot)$ is the function mapping the original input to the latent space, and $\bar{\alpha}_t$ represents the cumulative variance at timestep t .

Classifier-free guidance [44] is a well-established technique in the field of diffusion models. During the inference phase, we enhance controllability through classifier-free guidance with dual conditioning. The output at timestep t is given by:

$$\hat{\epsilon}_\theta(x_t, c_t, c_{art}, t) = w \epsilon_\theta(x_t, c_t, c_{art}, t) + (1 - w) \epsilon_\theta(x_t, t), \quad (19)$$

where the classifier-free guidance factor w is set to 0.6.

4 EXPERIMENTS

4.1 Experimental Setup

We train ArtCrafter on a training data consisting of about 500,000 real image-text pairs from LAION-Aesthetic [45] and 50k art-text pairs from our proposed ArtMarket dataset. The images are paired with text descriptions generated by BLIP-2 [46], forming image-text data pairs. During both training and inference, we resize the input images to a spatial resolution of 512×512 . We have implemented our method over the stable diffusion 1.5 version [1]. Our training processes are conducted using 8 NVIDIA A100 GPUs, each with 80GB of memory, and a batch size of 8 per GPU. The inference phase, which consumes 5185 MiB of memory, takes about one second of sampling time on a single A100 at denoising steps of 50.

4.2 Qualitative Evaluations

We evaluate our proposed method by comparing it with various existing methods, including but not limited to StyleShot [17], Style Aligned [38], VSP [37], InstantStyle [18], InstantStyle (Plus) [2], [18], IP-Adapter [16], IP-Adapter (SDXL) [2], [16], CSGO [27], and StyleCrafter [14]. We utilized the publicly available implementations of these methods and followed their recommended configurations for testing.

The qualitative comparison in Fig. 4 provides a visual assessment of the results achieved by various related methods. StyleShot generally increases the saturation and brightness of the hue. Style Aligned often produces double images in semantic expression, which may overlap with the semantic information of the input image. VSP has significant differences from the reference image in terms of color and texture. InstantStyle and its enhanced version repeatedly generate the same object within a single image, possibly misinterpreting semantic information as patterns and thus

TABLE 1

Quantitative Comparison. To provide a comprehensive evaluation of our method against other state-of-the-art text-guided stylization techniques, we conducted extensive quantitative experiments using a diverse set of metrics. These metrics cover key aspects of image quality, including image consistency (measured by CLIP-Image and DINO-v2), text consistency (assessed via CLIP-Text), and diversity (quantified using LPIPS). The best results are highlighted in **bold**, and the second-best results are marked with underline.

Metric	Ours	Styleshot	Style Aligned	VSP	InstantSt.	InstantSt. (Plus)	IP-Ada.	IP-Ada. (SDXL)	CSGO	StyleCrafter
CLIP-Text \uparrow	22.57	<u>22.46</u>	18.26	22.29	19.78	19.53	15.01	19.14	17.82	19.37
CLIP-Image \uparrow	69.48	59.47	63.82	66.31	62.59	64.65	<u>68.90</u>	67.43	55.16	58.68
DINO-v2 \uparrow	40.92	23.45	27.89	35.34	29.76	30.62	<u>38.12</u>	32.98	24.53	31.60
LPIPS \uparrow	0.4908	<u>0.4561</u>	0.3892	0.2653	0.1478	0.1234	0.3456	0.2987	0.3655	0.3194

TABLE 2

User Study. To further validate the effectiveness of our method from a human perception perspective, we conducted a comprehensive user study. Participants were asked to evaluate the generated images based on three key criteria: content consistency (Human-Content), style consistency (Human-Style), and overall quality (Human-Overall). The best results are in **bold** while the second-best results are marked with underline.

Metric	Ours	Styleshot	Style Aligned	VSP	InstantSt.	InstantSt. (Plus)	IP-Ada.	IP-Ada. (SDXL)	CSGO	StyleCrafter
Human-Content \uparrow	4.27	<u>3.96</u>	3.19	3.85	3.08	2.72	2.92	3.27	2.38	3.42
Human-Style \uparrow	4.08	<u>3.08</u>	3.04	3.58	2.92	3.15	<u>3.81</u>	3.46	2.65	2.81
Human-Overall \uparrow	4.19	<u>3.73</u>	3.23	3.65	2.81	2.96	3.50	3.15	2.54	2.88

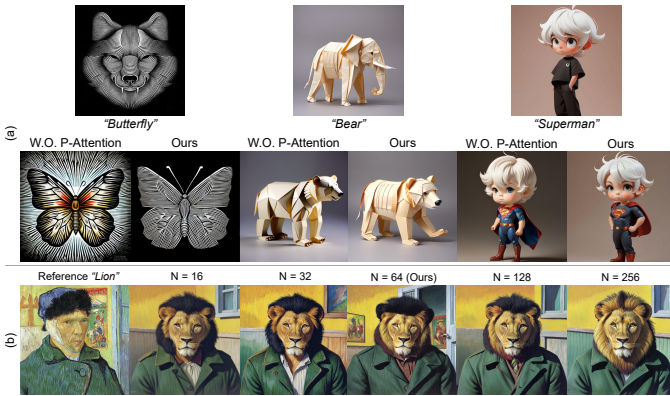


Fig. 5. **Ablation studies for attention-based Style Extraction.** (a) Visual differences between ArtCrafter generated samples without preceiver attention and our method. (b) Discussion on the number of learnable queries N.

causing their repeated appearance. The content information in the style image also affects the IP-Adapter and its SDXL version; for example, the “apple” content is introduced in the “Forbidden City, Beijing” case. The fidelity of the results generated by CSGO is constrained by the number of art datasets. StyleCrafter produces outputs that are more consistent with the text, but the styles differ significantly from the reference image. As shown in Fig. 4, the text-guided semantic content is skillfully integrated into the painting, and our method achieves a more harmonious and aesthetically pleasing transfer result globally.

4.3 Quantitative Evaluations

Metric Description. To comprehensively assess the quality of the models, we computed the following metrics: **CLIP-Text** [41] measures the textual consistency of the generated image with the target description by calculating the cosine similarity between the target caption and the generated image. **CLIP-Image** [41] evaluates the similarity of the generated image to the target style by calculating the cosine

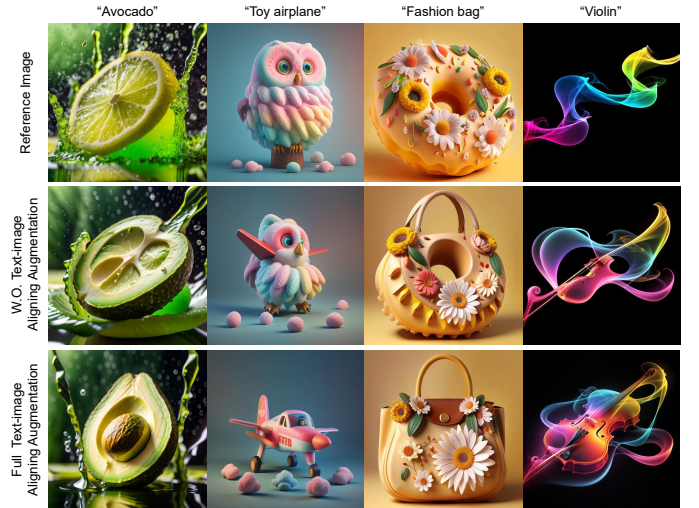


Fig. 6. **Ablation studies for text-image aligning augmentation.** All results use the same seed as well as setup factors. From the results at the bottom of the figure, it can be observed that this component plays a key role in the function of text conditioning.

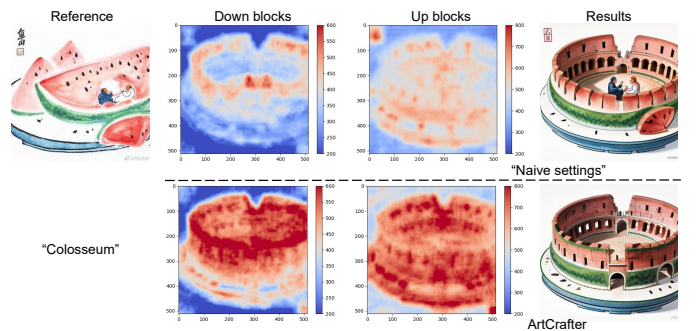


Fig. 7. **Ablation studies for text-image aligning augmentation.** Results of the “Naive setting” show a tendency to favor the reference image content over the textual content and fail to adequately reconstruct the target content. With the addition of the text-image aligning augmentation module (bottom), the rearranged attention enables more favorable text-based content reconstruction.

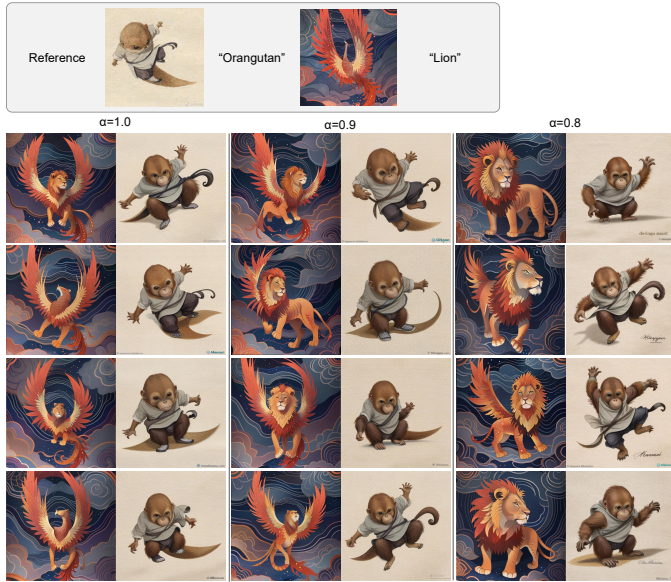


Fig. 8. **The effect of varying degrees of explicit modulation.** As α is increased, more varied results can be obtained (right columns), with some styles of misalignment of details.

similarity between the target styling image and the generated image. **DINO-v2** [47] further verifies the style consistency of the generated image with the target style by calculating the feature similarity between the target style image and the generated image. **LPIPS** [48] measures the perceived similarity between two generated images, with higher values indicating less image similarity and better diversity.

As shown in Table 1, our method outperformed other text-guided stylization methods in the CLIP-Text, CLIP-Image, DINO-v2, and LPIPS metrics. This indicates that our method has a significant advantage in generating images that are highly consistent with the target description and style. Moreover, our method effectively retains the details and semantics of the content while maintaining high stylization quality as well as diversity.

4.4 User Study

To obtain a more comprehensive assessment, we conducted a user study. We randomly selected 30 generated results for each method covering a wide range of styles and hired 26 professionals in the art field to evaluate these generated results. Specifically, they rated text and image consistency on a scale of 1 – 5, resulting in the Human-Text and Human-Image scores in Table 2. The images with the best overall results are then pulled out to obtain the Human-Overall percentage results. The results of Table 2 indicate that the results generated by our proposed ArtCrafter are more favored in all three aspects. We notice the difference between objective evaluation metrics and subjective evaluation metrics because the former assesses each aspect in isolation. Whereas users may integrate information across various aspects, despite separate options provided. Human Preference suggests that our generated results have struck a better balance among text consistency, image consistency, and overall visual appeal.

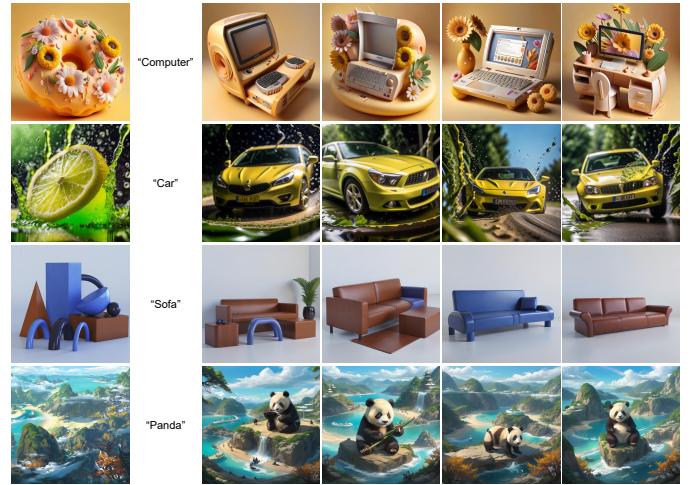


Fig. 9. **Diversity results with explicit modulation ablation study.** We conducted an ablation study to assess the impact of the explicit modulation component on the diversity of generated images. The results demonstrate that the explicit modulation component significantly enhances the diversity of the outputs.

4.5 Ablation Study

4.5.1 Visualization analysis

Attention-based style extraction. In our study, we observed that naive cross-attention tends to cause a shift in style, as shown in the results of Fig. 4 (Styleshot, Stylealigned, Stylecrafter, CSGO). The novel approach is by processing the embedded image’s latent features with perceptual attention (Sec. 3.2). This method not only effectively prevents unintended style shifts but also significantly enhances artistic expression, particularly in terms of color and texture. As demonstrated in Fig. 5 (a), the images processed with perceptual attention show a remarkable improvement in the richness of color and the detail of texture, making the generated images more artistically appealing and visually attractive.

To further optimize the performance of our model, we conducted an in-depth ablation study on the key parameter N . As shown in Figure 5 (b), we found that the value of N has a crucial impact on the model’s attention distribution and generation results. When N is low, the model’s attention becomes overly focused on narrow feature regions, leading to the loss of feature information and resulting in images that appear overly localized and lack a sense of coherence. Conversely, when N is high, the model’s attention becomes too dispersed, blurring the key attributes of the noise and causing the generated images to lack clear structure and detail. Through this series of experiments, we identified an optimal balance point, enabling the model to capture and reproduce key image features accurately while maintaining the overall style.

Text-image aligning augmentation. To evaluate the effectiveness of the TIAA (Sec. 3.3), we present visual representations of the stylization results in Fig. 6 and Fig. 7. Fig. 6 clearly illustrates that the TIAA component significantly enhances the alignment between text and image, thereby improving the overall efficacy of the text role. Moreover, the heatmap in Fig. 7 reveals the cosine similarity between the cross-attention of the generated results and the textual self-

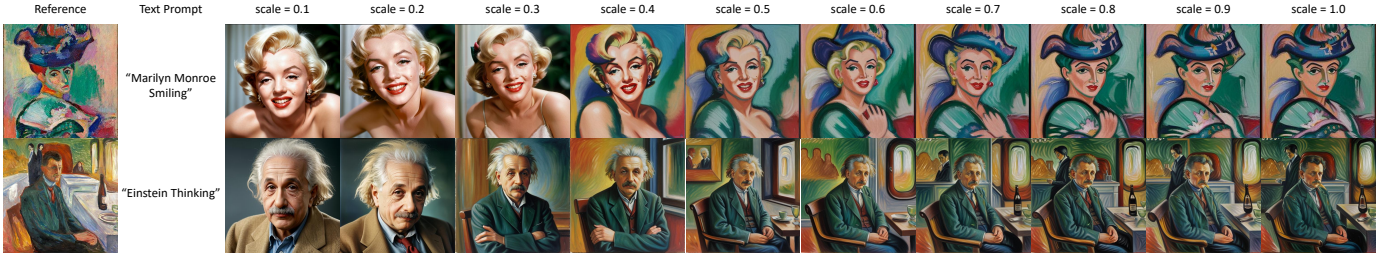


Fig. 10. **Visual demonstrations of ArtCrafter’s influence on pre-trained diffusion models.** The illustrations reveal how ArtCrafter’s modulation enhances the model’s ability to generate images that are stylistically coherent with varying degrees of influence. As the intensity of ArtCrafter’s application increases, there is a noticeable evolution in the image’s stylistic attributes, demonstrating the adaptability and controllability of our method.

TABLE 3

Quantitative ablation study of proposed components. The results provide insights into how each component affects the performance of key metrics such as image consistency, text consistency, and diversity.

Configuration	CLIP-Text \uparrow	CLIP-Image \uparrow	DINO-v2 \uparrow	LPIPS \uparrow
W.O. ASE 3.2	23.49	66.81	37.25	0.4971
W.O. TIAA 3.3	19.13	70.52	41.77	0.4316
W.O. EM 3.4	22.19	71.39	42.63	0.3572

erating images that not only capture the essence of the input text but also exhibit a wide range of artistic styles, thereby validating its importance in achieving diverse and compelling visual outcomes.

4.5.2 quantitative analysis

To fully evaluate the performance of the proposed components, we conduct an ablation study using both quantitative and qualitative methods. Table 3 summarizes the quantitative results, in which we test ArtCrafter’s performance after removing each essential technique component individually. There are inherent trade-offs between image consistency, text consistency, and diversity. Each module enhances a specific aspect: ASE improves image performance (CLIP-Image, DINO-v2), TAA boosts text performance (CLIP-Text), and EM parametrically balances and diversifies the results (LPIPS). Our goal is to achieve an optimal balance among these aspects to maximize overall effectiveness. Collectively, these findings underscore the significance and efficacy of each component within the overall strategy.

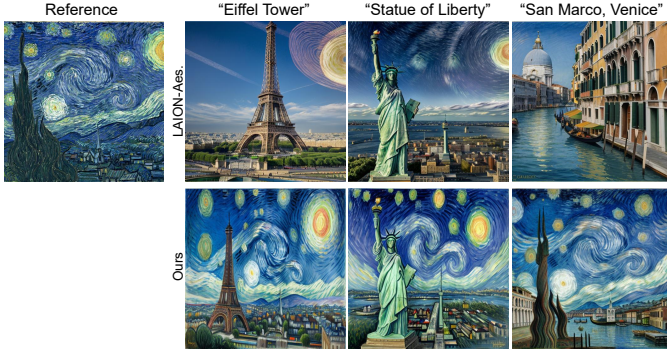


Fig. 11. **Visualization results of ArtCrafter architecture trained on different datasets.** Specifically, when trained on a dataset with rich artistic styles, ArtCrafter generates images with intricate and varied artistic features. In contrast, when trained on a dataset with more natural images, the model produces outputs that are highly consistent with natural visual elements. These results underscore the robustness and versatility of ArtCrafter in generating high-quality images.

attention results. This indicates that ArtCrafter, through its fine-grained attention mechanism, effectively aligns the generated images with the textual descriptions. This visual evidence strongly supports the success of our strategy in achieving tighter compatibility and refinement between text and images via the TIAA module.

Explicit modulation. We investigate the impact of the dynamic scale α in EM (Sec. 3.4) on the diversity of generated results by adjusting its value. As illustrated in Fig. 8, increasing α leads to a more diverse set of images while maintaining common stylistic attributes with the reference image. We typically set α to 0.8.

Fig. 9 provides detailed ablation results for our EM component, further demonstrating its effectiveness in enhancing content diversity. These results highlight the component’s ability to significantly enrich the variety of outputs, ensuring a broader spectrum of creative expressions. The experiments detailed in Fig. 9 underscore how EM contributes to gen-

4.5.3 ArtCrafter control

We test the role of ArtCrafter in the pre-training diffusion model. The result of adjusting the scale of ArtCrafter in the pre-trained diffusion model is shown in Fig. 10. The figure illustrates the outcomes at intervals of 0.1 for the guidance coefficient, thereby providing a granular view of how varying levels of input influence the final generation quality. In this work, we usually set the scale to 0.6. This verifies that ArtCrafter can effectively enhance the applicability of pre-trained diffusion models in the art generation domain at a lower training cost. With this adjustment, we can more flexibly control the degree of text consistency and stylization of the generated images to meet different creative needs.

4.5.4 Dataset

We constructed the ArtMarket dataset using art images from WikiArt [49] and LAION-Aesthetics [50]. This dataset is specifically designed to optimize our model’s image encoder, while the text encoder remains frozen throughout the training process and is not fine-tuned on this dataset. During inference, all text conditions are processed by a frozen CLIP text encoder, ensuring the universality of text semantic understanding while allowing the image encoder to focus on artistic style extraction.

We utilized BLIP-2 [46] as the textual style descriptive model to create art-description data pairs. The text



Fig. 12. **Integration of ArtCrafter with additional conditions.** The outcomes of this integration indicate that ArtCrafter not only effectively incorporates images as content guidance for style transfer but also seamlessly integrates with 3D content information. This compatibility showcases ArtCrafter’s potential to enhance the richness of generated images by leveraging various types of input data, thereby expanding its applicability in diverse creative scenarios.

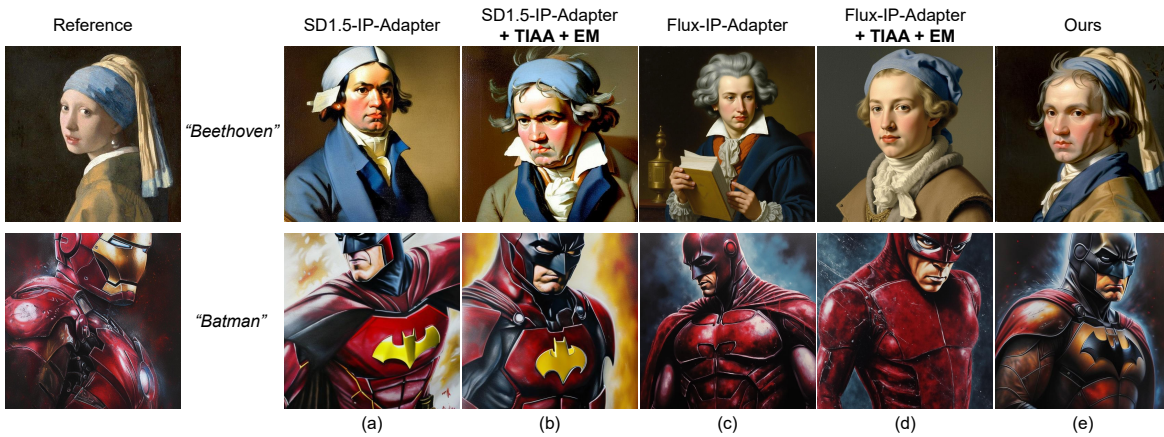


Fig. 13. **Enhanced comparisons with IP-Adapter variants.** To further validate the superiority of our approach, we conducted additional comparative analyses with the SD1.5 and Flux versions of the IP-Adapters, which were enhanced with our TIAA and EM algorithms. These comparisons provide a comprehensive view of how our algorithms contribute to the performance of these models.

descriptions in our dataset serve as high-level, semantic supervisory signals to guide the image encoder to learn discriminative artistic style representations that align with human perception, rather than focusing on low-level textures. This ensures that the extracted style embedding E_I captures the semantic essence of “style” rather than irrelevant visual noise.

To evaluate the impact of our dataset, we conducted an ablation study by training the model solely on the LAION-Aesthetics dataset. As shown in Fig. 11, the model trained only on LAION-Aesthetics tends to generate realistic land-

scape images. In contrast, our model trained on the combined dataset effectively captures the distinctive style of Van Gogh’s *Starry Night*. This demonstrates that incorporating diverse artistic data significantly enhances the model’s ability to generate a wide range of artistic styles.

5 DISCUSSION

5.1 Application

5.1.1 ArtCrafter with additional conditions applied

Benefiting from our design without any changes to the network structure of the original diffusion model, ArtCrafter

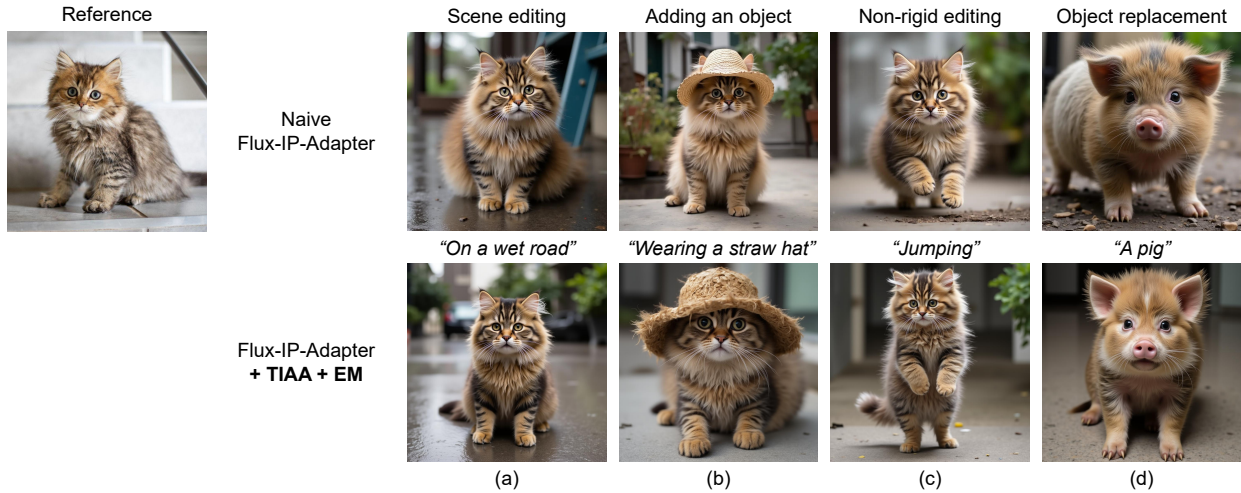


Fig. 14. **Versatility of our algorithms in natural image editing.** Our algorithms are designed to be versatile and applicable across a broad spectrum of natural image editing applications. They are effective regardless of the artistic dataset used, demonstrating their robustness and adaptability. This capability allows for a seamless integration into various workflows, enhancing the potential for creative expression.

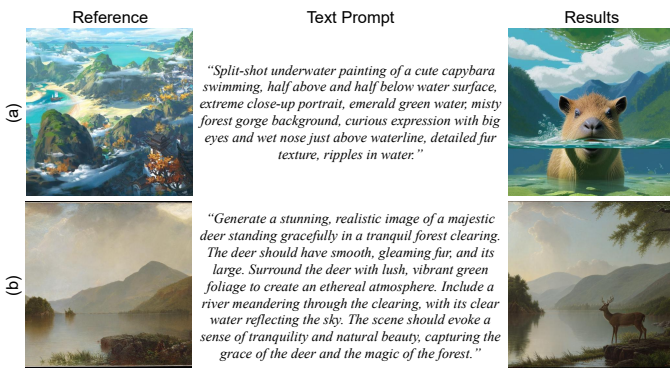


Fig. 15. **Enhanced comprehension of extended text.** Our model demonstrates an advanced ability to interpret and understand lengthy and intricate textual descriptions.



Fig. 16. **Dataset samples.** The dataset samples cover a diverse range of artistic styles and periods, ensuring a comprehensive representation of the artistic spectrum. The textual descriptions are crafted to capture the essence of each artwork, providing detailed insights into the visual elements, thematic content, and stylistic nuances.

is seamlessly compatible with existing controllable tools [25]. Fig. 12 shows the diverse examples generated by applying different structural controls, including canny edge detection [51], normal map [52], HED edge detection [53], lineart edge, and MiDaS depth map [54]. These examples not only highlight ArtCrafter’s flexibility in adapting to a variety of conditions but also foretell its great potential for application in the 3D field, opening up new possibilities for artistic creation and visual design.

5.1.2 More applications on other T2I models

To further verify the transferability and universality of the proposed method in multi-functional application scenarios, we conducted a series of in-depth comparative experiments on various text-to-image generation models with different weight configurations and architectural designs. As shown in Fig. 13, we present the comparative results of two representative model versions—Stable Diffusion 1.5 (SD1.5) and Flux—after integrating two key modules of our method: the TIAA module and the EM module. The experimental results demonstrate that by embedding these two modules into the models in a lightweight integration manner, significant performance improvements can be achieved without the need for additional training or fine-tuning. This fully proves the transferability and efficiency of our method.

5.1.3 Natural image editing

The strengths of ArtCrafter are not confined to its training on the ArtMarket dataset; its performance in various natural image editing scenarios is also highly impressive. Fig. 14 demonstrates the remarkable capabilities of TIAA and EM in jointly enhancing image-text guidance. Our method can be widely applied to multiple natural image editing scenarios, including scene editing (such as adjusting the atmosphere or weather of a scene), adding objects (naturally integrating new objects into an image), non-rigid editing (deforming objects), and object replacement (replacing one object in an image with another). In Fig. 14 (a), the improved results more accurately reflect the “road” information from the text, such as more clearly presenting the texture, direction, or surrounding environment of the road. In Fig. 14 (b),

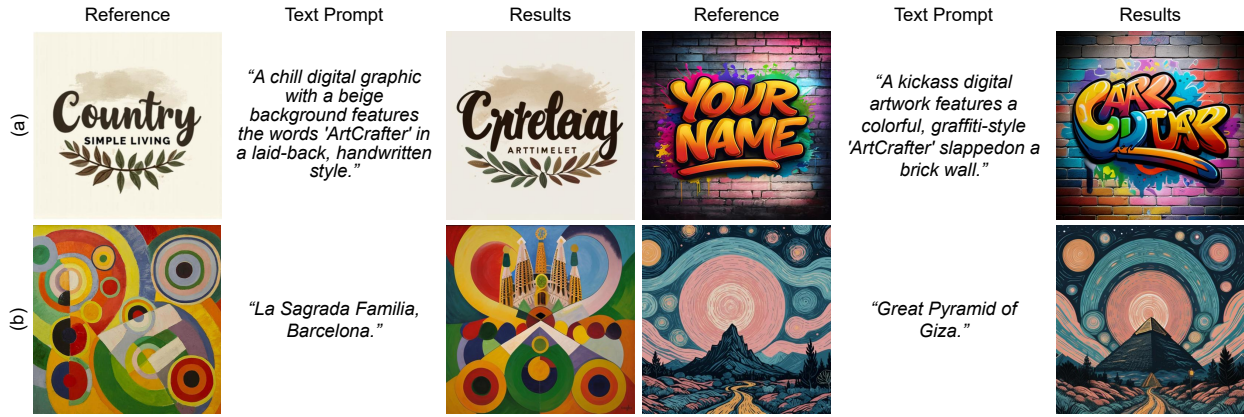


Fig. 17. **Analysis of failure cases.** (a) Text generation failed. (b) Instances of inconsistent line texture.

(c), and (d), not only are the requirements of the editing text better enhanced, but the information of non-edited parts, such as the details of the indoor environment, is also preserved, ensuring the naturalness and coherence of the overall image. This capability makes ArtCrafter valuable and broadly applicable in the field of natural image editing.

5.1.4 Long text prompt

In Fig. 15, we present a series of long and complex textual requirements that our model successfully translates into detailed visual outputs. For instance, in Fig. 15 (a), the model accurately generates an image with a deer having a "curious expression with big eyes and a wet nose just above the waterline", "detailed fur texture", "and ripples in the water", capturing the "half above and half below water surface" scenario. In Fig. 15 (b), the model creates a serene forest scene with a deer "surrounded by lush foliage," "a river meandering," and "clear water reflecting the sky." These examples demonstrate the robust capability of our text-image alignment design in handling intricate textual instructions and generating images that precisely match the detailed descriptions.

5.2 Dataset Samples

Our dataset ArtMarket, as shown in Fig. 16, contains a rich collection of artworks along with their corresponding textual descriptions. The visual content of the dataset is primarily sourced from the extensive art repositories of WikiArt [49] and LAION-Aesthetics [50], which cover a wide range of artistic styles and periods. The descriptive texts are derived from BLIP-2 [46], an advanced language model that can generate detailed and contextually rich captions. The combination of high-quality images and accurately descriptive texts enables ArtMarket to effectively support the training of models for style transfer and content generation in the field of text-to-image synthesis.

5.3 Limitation

The limitations of ArtCrafter are mainly reflected in two aspects. First, it is the ability to accurately generate text. As shown in Fig. 17 (a), ArtCrafter has difficulty in accurately generating text in images and cannot ensure the clarity and accuracy of the text content "ArtCrafter". This limitation may seriously affect its performance in image generation tasks

that require text elements. For example, in designing posters, brochures, or other image editing tasks that require precise text layout, ArtCrafter may not be able to meet the high-precision requirements of users for text content and layout. Secondly, the model can only partially follow large-scale textures and brushstrokes. For example, in Fig. 17 (b), our method incorrectly transforms the large-scale curves from left to right in the image into inconsistent curve shapes. This error indicates that the model may fail to fully understand and accurately reproduce the original shapes and directions when dealing with complex textures and brushstrokes. Accurately replicating complex long-range texture patterns is one of the challenges for our model to fully assimilate and reproduce in the generated output.

6 CONCLUSIONS AND FUTURE WORK

In this paper, we introduce ArtCrafter, a novel text-image aligning style transfer framework achieved through an embedding reframing architecture. Our approach ensures superior text-guided style transfer quality by integrating three core components: attention-based style extraction, text-image aligning augmentation, and explicit modulation. Comprehensive evaluations demonstrate ArtCrafter strengths in adapting to diverse artistic styles, maintaining textual prompt consistency, enhancing output diversity, and improving overall visual quality.

Acknowledging the current limitations of our work, for future research, we intend to enhance our approach by incorporating pattern reproducibility and contextual elements within style images, including the relative positioning of style patches, to facilitate a more cohesive art style transfer. We expect that advancements in the extraction and fusion of style and content features, coupled with an investigation into the method's scalability and adaptability, will markedly enhance the quality of style transfer and provide more precise control over the shape and appearance similarity of the generated images.

REFERENCES

- [1] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10 684–10 695.

- [2] D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, and R. Rombach, "Sdxl: Improving latent diffusion models for high-resolution image synthesis," *arXiv preprint arXiv:2307.01952*, 2023.
- [3] Z. Li, M. Cao, X. Wang, Z. Qi, M.-M. Cheng, and Y. Shan, "Photomaker: Customizing realistic human photos via stacked id embedding," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 8640–8650.
- [4] Q. Wang, X. Bai, H. Wang, Z. Qin, A. Chen, H. Li, X. Tang, and Y. Hu, "Instantid: Zero-shot identity-preserving generation in seconds," *arXiv preprint arXiv:2401.07519*, 2024.
- [5] M. Huang, Z. Mao, M. Liu, Q. He, and Y. Zhang, "Realcustom: Narrowing real text word for real-time open-domain text-to-image customization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 7476–7485.
- [6] J. Ma, J. Liang, C. Chen, and H. Lu, "Subject-diffusion: Open domain personalized text-to-image generation without test-time fine-tuning," in *ACM SIGGRAPH 2024 Conference Papers*, 2024, pp. 1–12.
- [7] J. Chung, S. Hyun, and J.-P. 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*, 2024, pp. 8795–8805.
- [8] Y. Zhang, N. Huang, F. Tang, H. Huang, C. Ma, W. Dong, and C. Xu, "Inversion-based style transfer with diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 10 146–10 156.
- [9] D.-Y. Chen, H. Tennent, and C.-W. Hsu, "Artadapter: Text-to-image style transfer using multi-level style encoder and explicit adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 8619–8628.
- [10] Y. Deng, X. He, F. Tang, and W. Dong, "Z*: Zero-shot style transfer via attention reweighting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 6934–6944.
- [11] Z. Zhang, Q. Zhang, W. Xing, G. Li, L. Zhao, J. Sun, Z. Lan, J. Luan, Y. Huang, and H. Lin, "Artbank: Artistic style transfer with pre-trained diffusion model and implicit style prompt bank," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 7, 2024, pp. 7396–7404.
- [12] R. Jiang, C. Wang, J. Zhang, M. Chai, M. He, D. Chen, and J. Liao, "Avatarcraft: Transforming text into neural human avatars with parameterized shape and pose control," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 14 371–14 382.
- [13] N. Huang, W. Dong, Y. Zhang, F. Tang, R. Li, C. Ma, X. Li, and C. Xu, "Creativesynth: Creative blending and synthesis of visual arts based on multimodal diffusion," *arXiv preprint arXiv:2401.14066*, 2024.
- [14] G. Liu, M. Xia, Y. Zhang, H. Chen, J. Xing, X. Wang, Y. Yang, and Y. Shan, "Stylecrafter: Enhancing stylized text-to-video generation with style adapter," *arXiv preprint arXiv:2312.00330*, 2023.
- [15] Y. Han, J. Zhu, K. He, X. Chen, Y. Ge, W. Li, X. Li, J. Zhang, C. Wang, and Y. Liu, "Face adapter for pre-trained diffusion models with fine-grained id and attribute control," *arXiv preprint arXiv:2405.12970*, 2024.
- [16] H. Ye, J. Zhang, S. Liu, X. Han, and W. Yang, "Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models," *arXiv preprint arXiv:2308.06721*, 2023.
- [17] J. Gao, Y. Liu, Y. Sun, Y. Tang, Y. Zeng, K. Chen, and C. Zhao, "Styleshot: A snapshot on any style," *arXiv preprint arXiv:2407.01414*, 2024.
- [18] H. Wang, Q. Wang, X. Bai, Z. Qin, and A. Chen, "Instantstyle: Free lunch towards style-preserving in text-to-image generation," *arXiv preprint arXiv:2404.02733*, 2024.
- [19] K. H. Lin, S. Mo, B. Klingher, F. Mu, and B. Zhou, "Ctrl-x: Controlling structure and appearance for text-to-image generation without guidance," *arXiv preprint arXiv:2406.07540*, 2024.
- [20] C. Rowles, S. Vainer, D. De Nigris, S. Elizarov, K. Kutsy, and S. Donné, "Ipadapter-instruct: Resolving ambiguity in image-based conditioning using instruct prompts," *arXiv preprint arXiv:2408.03209*, 2024.
- [21] J. Wang, Y. Pu, Y. Han, J. Guo, Y. Wang, X. Li, and G. Huang, "Gra: Detecting oriented objects through group-wise rotating and rotation," *arXiv preprint arXiv:2403.11127*, 2024.
- [22] N. Huang, Y. Zhang, F. Tang, C. Ma, H. Huang, W. Dong, and C. Xu, "Diffstyler: Controllable dual diffusion for text-driven image stylization," *IEEE Transactions on Neural Networks and Learning Systems*, 2024.
- [23] J. Singh, S. Gould, and L. Zheng, "High-fidelity guided image synthesis with latent diffusion models," in *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, 2023, pp. 5997–6006.
- [24] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.
- [25] L. Zhang, A. Rao, and M. Agrawala, "Adding conditional control to text-to-image diffusion models," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 3836–3847.
- [26] T. Brooks, A. Holynski, and A. A. Efros, "Instructpix2pix: Learning to follow image editing instructions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 18 392–18 402.
- [27] P. Xing, H. Wang, Y. Sun, Q. Wang, X. Bai, H. Ai, R. Huang, and Z. Li, "Csgo: Content-style composition in text-to-image generation," *arXiv preprint arXiv:2408.16766*, 2024.
- [28] O. Avrahami, D. Lischinski, and O. Fried, "Blended diffusion for text-driven editing of natural images," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18 208–18 218.
- [29] A. Hertz, R. Mokady, J. Tenenbaum, K. Aberman, Y. Pritch, and D. Cohen-Or, "Prompt-to-prompt image editing with cross attention control," 2022.
- [30] K. Zhang, Y. Li, W. Zuo, L. Zhang, L. Van Gool, and R. Timofte, "Plug-and-play image restoration with deep denoiser prior," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 10, pp. 6360–6376, 2021.
- [31] M. Cao, X. Wang, Z. Qi, Y. Shan, X. Qie, and Y. Zheng, "Masactrl: Tuning-free mutual self-attention control for consistent image synthesis and editing," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 22 560–22 570.
- [32] T. Qi, S. Fang, Y. Wu, H. Xie, J. Liu, L. Chen, Q. He, and Y. Zhang, "Deadiff: An efficient stylization diffusion model with disentangled representations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 8693–8702.
- [33] Z.-S. Liu, L.-W. Wang, W.-C. Siu, and V. Kalogeiton, "Name your style: text-guided artistic style transfer," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 3530–3534.
- [34] P. Ganugula, Y. S. S. S. Kumar, N. K. S. Reddy, P. Chellingi, A. Thakur, N. Kasera, and C. S. Anand, "Mosaic: Multi-object segmented arbitrary stylization using clip," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, October 2023, pp. 892–903.
- [35] H. Yang, Y. Chen, Y. Pan, T. Yao, Z. Chen, and T. Mei, "3dstyle-diffusion: Pursuing fine-grained text-driven 3d stylization with 2d diffusion models," in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 6860–6868.
- [36] K. Sohn, L. Jiang, J. Barber, K. Lee, N. Ruiz, D. Krishnan, H. Chang, Y. Li, I. Essa, M. Rubinstein *et al.*, "Styler: Text-to-image synthesis of any style," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [37] J. Jeong, J. Kim, Y. Choi, G. Lee, and Y. Uh, "Visual style prompting with swapping self-attention," *arXiv preprint arXiv:2402.12974*, 2024.
- [38] A. Hertz, A. Voynov, S. Fruchter, and D. Cohen-Or, "Style aligned image generation via shared attention," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 4775–4785.
- [39] SG161222, "Realistic-vision-v4.0-novae," 2024. [Online]. Available: https://huggingface.co/SG161222/Realistic_Vision_V4.0_noVAE
- [40] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention*. Springer, 2015, pp. 234–241.
- [41] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning*. PMLR, 2021, pp. 8748–8763.
- [42] A. Jaegle, F. Gimeno, A. Brock, O. Vinyals, A. Zisserman, and J. Carreira, "Perceiver: General perception with iterative attention," in *International Conference on Machine Learning*. PMLR, 2021, pp. 4651–4664.
- [43] A. Vaswani, "Attention is all you need," *Advances in Neural Information Processing Systems*, 2017.

- [44] J. Ho and T. Salimans, "Classifier-free diffusion guidance," in *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.
- [45] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta, C. Mullis, M. Wortsman *et al.*, "Laion-5b: An open large-scale dataset for training next generation image-text models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 25 278–25 294, 2022.
- [46] J. Li, D. Li, S. Savarese, and S. Hoi, "Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models," in *International Conference on Machine Learning*. PMLR, 2023, pp. 19 730–19 742.
- [47] M. Oquab, T. Darcet, T. Moutakanni, H. V. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby, R. Howes, P.-Y. Huang, H. Xu, V. Sharma, S.-W. Li, W. Galuba, M. Rabbat, M. Assran, N. Ballas, G. Synnaeve, I. Misra, H. Jegou, J. Mairal, P. Labatut, A. Joulin, and P. Bojanowski, "Dinov2: Learning robust visual features without supervision," 2023.
- [48] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 586–595.
- [49] W. Contributors, "Wikiart website," 2024. [Online]. Available: <https://www.wikiart.org>
- [50] LAION-AI, "LAION-Aesthetics V1," LAION-AI, Tech. Rep., 2024. [Online]. Available: https://laion-ai.github.io/laion-datasets/laion-aesthetic/laion_aesthetic.html
- [51] J. Canny, "A computational approach to edge detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, no. 6, pp. 679–698, 1986.
- [52] I. Vasiljevic, N. Kolkin, S. Zhang, R. Luo, H. Wang, F. Z. Dai, A. F. Daniele, M. Mostajabi, S. Basart, M. R. Walter *et al.*, "Diode: A dense indoor and outdoor depth dataset," *arXiv preprint arXiv:1908.00463*, 2019.
- [53] S. Xie and Z. Tu, "Holistically-nested edge detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1395–1403.
- [54] R. Ranftl, K. Lasinger, D. Hafner, K. Schindler, and V. Koltun, "Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 3, pp. 1623–1637, 2020.



Yifan Pu received the BS degree in Automation from Beihang University, Beijing, China, in 2020. After that, he was a Master's student with the Department of Automation, Tsinghua University, Beijing, China. He is currently pursuing the PhD degree after transferring into the doctoral program. His research interests include computer vision, machine learning, and deep learning.



Jiangshan Wang received his B.S degree in Beijing University of Posts and Telecommunications in 2023. He is currently pursuing the M.S. degree in Artificial Intelligence at Tsinghua University. His research interests include computer vision and generative models.



Jie Guo received her doctoral degree in School of Optics and Photonics from Beijing Institute of Technology in 2021. She is currently an assistant researcher in PengCheng Laboratory. Her main research areas include artificial intelligence, digital humans, virtual reality, and human-computer interaction.



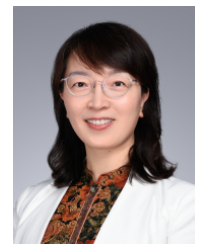
Nisha Huang received her B.S. degree in Aeronautical Information Engineering from Beihang University, Beijing, China, in 2021. After that, she obtained her M.S. degree in Artificial Intelligence from the University of Chinese Academy of Sciences, Beijing, China, in 2024. Currently, she is pursuing her PhD in Electronic Information at Tsinghua University. And she is conducting research at the PengCheng Laboratory. Her research interests include generative models, multimedia analysis, and computer vision.



Yiqiang Yan received the MS degree from Tsinghua University. He is currently the chief researcher of Lenovo, Beijing, China. His research interests include efficient large language models, reinforcement learning for GUI agents.



Kaer Huang received the MS degree from North China University, Taiyuan, China. He is currently working at Lenovo Research, Beijing, China. His research interests include perception algorithm, efficient large language models, reinforcement learning for GUI agents.



Xiu Li (Member, IEEE) received a Ph.D. degree in computer-integrated manufacturing from the Nanjing University of Aeronautics and Astronautics, Nanjing, China, in 2000. From then to 2002, she served as a Post-Doctoral Fellow at the Department of Automation, Tsinghua University, Beijing, China. From 2003 to 2010, she served as an Associate Professor at the Department of Automation, at Tsinghua University. Since 2016, she has been a Full Professor at the Tsinghua Shenzhen International Graduate School, Tsinghua University, Shenzhen, China. Her research interests are in the areas of data mining, deep learning, computer vision, and image processing.



Tong-Yee Lee (Senior Member, IEEE) received the Ph.D. degree in computer engineering from Washington State University, Pullman, in May 1995. He is currently a Chair Professor in the Department of Computer Science and Information Engineering, at National Cheng-Kung University, Tainan City, Taiwan. He leads the Computer Graphics Laboratory, National Cheng-Kung University (<http://graphics.csie.ncku.edu.tw>). His current research interests include computer graphics, non-photorealistic rendering, medical visualization, virtual reality, and media resizing. He is a Senior Member of the IEEE and a Member of the ACM. He also serves on the editorial boards of both the IEEE Transactions on Visualization and Computer Graphics, and IEEE Computer Graphics and Applications.