

# Painterly Image Harmonization using Diffusion Model

Lingxiao Lu  
MoE Key Lab of Artificial Intelligence,  
Shanghai Jiao Tong University  
China  
lulingxiao@sjtu.edu.cn

Jiangtong Li  
MoE Key Lab of Artificial Intelligence,  
Shanghai Jiao Tong University  
China  
keep\_moving-lee@sjtu.edu.cn

Junyan Cao  
MoE Key Lab of Artificial Intelligence,  
Shanghai Jiao Tong University  
China  
joy\_c1@sjtu.edu.cn

Li Niu\*  
MoE Key Lab of Artificial Intelligence,  
Shanghai Jiao Tong University  
China  
ustcnewly@sjtu.edu.cn

Liqing Zhang\*  
MoE Key Lab of Artificial Intelligence,  
Shanghai Jiao Tong University  
China  
zhang-lq@cs.sjtu.edu.cn

## ABSTRACT

Painterly image harmonization aims to insert photographic objects into paintings and obtain artistically coherent composite images. Previous methods for this task mainly rely on inference optimization or generative adversarial network, but they are either very time-consuming or struggling at fine control of the foreground objects (e.g., texture and content details). To address these issues, we propose a novel Painterly Harmonization stable Diffusion model (PHDiffusion), which includes a lightweight adaptive encoder and a Dual Encoder Fusion (DEF) module. Specifically, the adaptive encoder and the DEF module first stylize foreground features within each encoder. Then, the stylized foreground features from both encoders are combined to guide the harmonization process. During training, besides the noise loss in diffusion model, we additionally employ content loss and two style losses, *i.e.*, AdaIN style loss and contrastive style loss, aiming to balance the trade-off between style migration and content preservation. Compared with the state-of-the-art models from related fields, our PHDiffusion can stylize the foreground more sufficiently and simultaneously retain finer content. Our code and model are available at <https://github.com/bcmi/PHDiffusion-Painterly-Image-Harmonization>.

## CCS CONCEPTS

• **Computing methodologies** → **Appearance and texture representations; Image manipulation; Computer vision.**

## KEYWORDS

painterly image harmonization, diffusion model, style transfer

### ACM Reference Format:

Lingxiao Lu, Jiangtong Li, Junyan Cao, Li Niu, and Liqing Zhang. 2023. Painterly Image Harmonization using Diffusion Model. In *Proceedings of*

\*Corresponding authors

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).

MM '23, October 29–November 3, 2023, Ottawa, ON, Canada

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0108-5/23/10...\$15.00

<https://doi.org/10.1145/3581783.3612451>



**Figure 1: Painterly image harmonization aims to harmonize the inserted photographic foreground according to the background painting. From left to right, we present the background image, the composite image via cut-and-paste, the harmonized results of PHDNet [3] and our method.**

the 31st ACM International Conference on Multimedia (MM '23), October 29–November 3, 2023, Ottawa, ON, Canada. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3581783.3612451>

## 1 INTRODUCTION

The goal of painterly image harmonization is to integrate photographic objects into background paintings and achieve visual coherence. While standard image harmonization [7, 8] focuses on adapting low-level statistics (e.g., color, brightness), painterly image harmonization [3, 44, 58] is more challenging as it requires transferring high-level styles in addition to low-level statistics.

Existing works for this task can be roughly divided into two categories: optimization-based [36, 58] and feed-forward [3, 44, 56] approaches. For the optimization-based approaches [36, 58], they optimize over the composite image by minimizing the designed losses, which makes them very time-consuming and unsuitable for real-time applications. Besides, the feed-forward [3, 44, 56] methods mainly rely on Generative Adversarial Network (GAN) [15] and the trained model can directly generate the harmonized images. However, one limitation of GAN-based approaches is limited control over complex foregrounds [49], resulting in unsatisfactory harmonized foregrounds (e.g., loss of content and style details).

In recent years, diffusion models [21] have demonstrated comparable or better performance compared with state-of-the-art image generation models, by formulating image generation as sequential stochastic transitions from a simple distribution to data distribution. Diffusion methods can be divided into unconditional diffusion methods [21, 51] and conditional diffusion methods [39, 46]. The

unconditional diffusion methods aim to generate realistic images by modeling the distribution of natural images without conditioning on any specific input. Whereas, the conditional diffusion methods aim to generate images with the guidance of conditional information (e.g., text, semantic mask, etc.). Among them, Stable Diffusion (SD) [46] is one of the most popular models, which successfully integrates the text CLIP [45] into latent diffusion. Further, some more recent models (e.g., T2I-adapter [39], ControlNet [57]) freeze the SD model and introduce trainable adapters to encode different types of conditions into SD through multi-step guidance. Therefore, conditional diffusion models offer a promising and flexible approach to improve painterly image harmonization by enabling multi-step guidance for the photographic foreground.

Several existing works have introduced diffusion methods into similar tasks like cross-domain image composition [18] and image editing [38]. For example, CDC [18] proposed an inference-time conditioning method that uses high-frequency details from the background and low-frequency style from the foreground object for image composition. However, CDC [18] assumes that high-frequency (*resp.*, low-frequency) feature in the image represents style (*resp.*, content) information, which does not always hold. Another work SDEdit [38] synthesizes images by first adding noise to the input image and then iteratively denoising through a stochastic differential equation. However, this approach lacks proper and sufficient guidance during the denoising process, leading to the final image lacking sufficient styles and contents.

In this paper, we introduce Painterly Harmonization stable Diffusion model (PHDiffusion), which exploits two extra modules based on the Stable Diffusion (SD) model. Inspired by the conditional diffusion model [39], we equip SD with a lightweight adaptive encoder, which aims to extract the required condition information (*i.e.*, background style, image content) from the composite image. As part of denoising U-Net, the denoising encoder in SD takes composite image as input. The adaptive encoder takes in the concatenation of composite image and foreground mask, producing residuals added to the feature maps in the denoising encoder. Based on the denoising encoder and the adaptive encoder, we introduce a Dual Encoder Fusion (DEF) module to fuse the information from two encoders. Specifically, given the image features extracted by two encoders, our DEF module incorporates the background style into foreground content and generates the stylized foreground features. Then, the stylized foreground features from two encoders are combined to provide multi-step guidance in the denoising steps.

To utilize the rich prior knowledge in pretrained SD and relieve the training burden, following [39], we freeze the model parameters of SD, and only update the adaptive encoder and DEF module during training. The standard noise loss used in diffusion models [46] could maintain the image content, but cannot migrate background style to the foreground. Therefore, we further introduce two additional style losses, *i.e.*, AdaIN loss and contrastive style loss, to balance the style and content for foreground object. The AdaIN loss [22] aligns the multi-scale statistics (e.g., mean, variance) of the foreground object with the background painting, while the contrastive style loss [5] aims to push the foreground style towards background style. In addition, we also incorporate a content loss to address the issue of excessive content preservation by using noise loss alone. With noise loss, style losses, and content loss, our PHDiffusion is able

to comprehend the background style and preserve the foreground content. During testing, our PHDiffusion could be directly used to produce harmonized image, preventing additional time-consuming inference optimization [18, 26].

To verify the effectiveness of our PHDiffusion, we compare our methods with the state-of-the-art methods, and conduct experiments on the benchmark datasets COCO [32] and WikiArt [41]. **The experimental results show that our PHDiffusion can achieve certain visually pleasant results that previous methods cannot achieve, especially when the background has dense textures or abstract style.** Our contributions can be summarized as follows: 1) We are the first work focusing on painterly image harmonization using diffusion model. 2) We propose a Painterly Harmonization stable Diffusion model (PHDiffusion) by using dual encoder fusion to provide effective guidance and reasonable loss designs to achieve sufficient stylization. 3) The experimental results show that our PHDiffusion strikes a good balance between adapting styles and maintaining structures.

## 2 RELATED WORK

### 2.1 Image Harmonization

As a subtask of image composition [42], image harmonization aims to adjust the color and illumination statistics of foreground to be compatible with background in a composite image. In recent years, deep learning methods [4, 6, 17, 54, 62] play an important role in this field. Especially after the first large-scale image harmonization dataset iHarmony4 [8] was released, supervised image harmonization [2, 9, 16, 19, 20, 33] methods have received more and more attention. For example, DoveNet [8] approached image harmonization as a domain translation task. Hao *et al.* [20] utilized attention block to calculate non-local information for foreground adjustment. SSAM [9] focused on relation between the spliced region and non-spliced region by exploiting a dual path attention model to fuse them together. CDTNet [7] combined pixel-to-pixel transformation and RGB-to-RGB transformation for high-resolution image harmonization. Recently, DCCF [55] and  $S^2CRNet$  [31] are applied in this field for high resolution image harmonization. Note that the abovementioned methods require ground-truth image to supervise, which is not suitable for our task.

### 2.2 Painterly Image Harmonization

As a similar task to image harmonization, painterly image harmonization aims to blend a photographic foreground into an artistic background painting, resulting in a visually coherent painting. Compared with image harmonization, painterly image harmonization is more challenging as it needs to adapt high-level styles beyond low-level statistics. Deep Painterly Harmonization [36] introduced a two-pass algorithm to ensure both spatial and inter-scale statistical consistency. Meanwhile, Deep Image Blending [58] utilized a two-stage blending algorithm and proposed a Poisson blending loss to guide blending together with content and style loss. However, both Deep Painterly Harmonization [36] and Deep Image Blending [58] are optimization-based method, which optimizes the input image during inference, making them unusable from real-time harmonization. On the other hand, E2STN [44], PHDNet [3], and Yan *et al.* [56] exploited the feed-forward scheme by first training

the generator and then directly producing the harmonized image during inference. Specifically, E2STN [44] took advantages of both global and local discriminators to harmonize the embedded element with the background image. PHDNet[3] exploited spatial and frequency domains to capture different types of background type, and then adjusted the foreground in both domains. Yan *et al.* [56] integrated GP-GAN [53], WCT [29], and StyleTr<sup>2</sup> [10] together to fuse the global and local information together. Note that, these feed-forward approaches [3, 44, 56] are mainly based on adversarial learning by playing a minimax game between generator and discriminator, which have limited control over complex photographic foreground [49] and have difficulty in leveraging prior knowledge across different image domains [15]. Different from them, our method is built upon the diffusion model, with stylized foreground features as guidance and a combo of style losses to produce the harmonized result.

### 2.3 Artistic Style Transfer

Artistic style transfer aims to stylize a content image given a style image. Previous optimization-based methods [13, 14, 25] optimize the content image to match its style with the style image. In contrast, feed-forward [10, 22, 24, 28, 35, 43] methods generate stylized images by training a generator to combine the content image and the style image. For example, style-relevant statistics [30, 60] (*e.g.*, mean and standard deviation of feature map) between the style image and fused image should be similar, and content-relevant information [11] (*e.g.*, the categories of objects) within the fused image should also be kept from the content image. Moreover, to enhance the visual quality in artistic style transfer, contrastive learning is also introduced [5, 61] to capture sufficient style information. Artistic style transfer methods stylize the entire content image, while painterly image harmonization focuses on the inserted object in the background painting.

### 2.4 Diffusion Models

Recently, diffusion models have shown remarkable performance in image generation [21, 34, 51], text-to-image generation [40, 46], image translation [26], image inpainting [37, 47], and image editing [18, 23, 38, 59]. Image editing and cross-domain image composition are the most relevant fields to our painterly image harmonization. Hence, we focus on these two fields with diffusion models in this section. Specifically, SDEdit [38] employed the image synthesis approach that commences with the addition of noise to the input image, followed by iterative denoising process with stochastic differential equation. CDC [18] proposed to harmonize the image in frequency domain by exploiting high-frequency details from the background and low-frequency style from the foreground object. Besides, there are a few diffusion models designed for artistic style transfer. To name a few, DiffStyle [23] disentangled representations for content and style, and fused them in h-space, which lies in the bottleneck of U-Net. InST [59] was motivated by the belief that an unique artwork can not be directly explained by words, so it designed an encoding module that maps style image into text domain through a CLIP image encoder.

However, previous diffusion-based methods can not provide powerful style guidance and hold adequate content in painterly

image harmonization. In this work, we endow stable diffusion with stylized feature guidance and well-designed losses for adjusting sufficient styles and maintaining content details, leading to better harmonization performance.

## 3 METHOD

The overall framework of our PHDiffusion is depicted in Figure 2, which consists of a Stable Diffusion (SD) model, an adaptive encoder, and a Dual Encoder Fusion (DEF) module. Given a composite image  $I_c$  and the corresponding foreground mask  $M$ , we first exploit the adaptive encoder to extract the multi-scale composite feature maps  $F_c^i$ ,  $i \in \{1, 2, 3, 4\}$ . In order to guide the generation of SD, we fuse the composite feature map  $F_c^i$  with the corresponding denoising feature map  $F_{z_t}^i$  from the U-Net encoder of SD according to their resolutions through the DEF module to generate new denoising feature map  $\hat{F}_{z_t}^i$ . Note that, we freeze the SD, and supervise the adaptive encoder and the DEF module using noise loss, content loss, and two style losses (AdaIN loss and contrastive style loss).

Next, we will first briefly review SD model, introduce the adaptive encoder, and then elaborate on our DEF Module. Finally, we will introduce the objectives for training adaptive encoder and DEF. As for the notation in the remainder of this section,  $z'_0$  is the initial latent feature extracted by encoder from image  $I$ .  $z'_t$ ;  $t = 1, 2, \dots, T$  represents the latent feature that is deduced from  $z'_0$  in the forward process of diffusion  $q(z'_t|z'_0)$ . While  $z_t$ ;  $t \in \{T-1, \dots, 0\}$  is predicted in the backward denoising process.  $\hat{I}_0$  is the decoded harmonized image from  $z_0$ .  $\hat{M}$  is the mask that is down-sampled to size of  $z'_0$ .

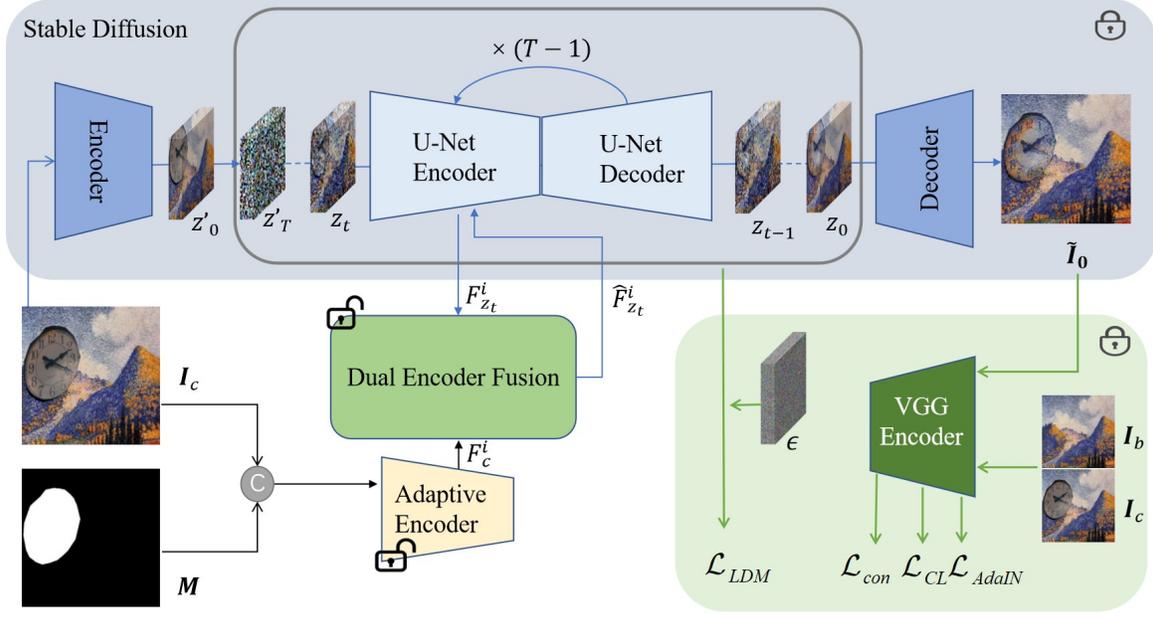
### 3.1 Preliminaries

Our method is built upon the Stable Diffusion (SD) [46] model, where SD is a latent diffusion model pretrained in two stages comprised of an auto-encoder and a denoising U-Net. In the first stage, the SD model trains the auto-encoder, in which the encoder  $\mathcal{E}$  first encodes images  $I$  into latent space  $z'_0 = \mathcal{E}(I)$  and then the decoder  $\mathcal{D}$  reconstructs them into original images  $\hat{I} = \mathcal{D}(z'_0)$ . In the second stage, the auto-encoder is frozen and the SD constructs the denoising U-Net  $\epsilon_\theta$  [21] by first adding  $T$ -step noise to latent space feature  $z'_0$  to generate  $z'_t$ ;  $t = 1, 2, \dots, T$ , and then training the denoising U-Net with latent denoising loss, which is formulated as

$$\mathcal{L}_{LDM} := \mathbb{E}_{z'_0, y, \epsilon \sim \mathcal{N}(0,1), t} \left[ \left\| \epsilon - \epsilon_{\theta_1}(z'_t, t, \tau_{\theta_2}(y)) \right\|_2^2 \right], \quad (1)$$

where  $\epsilon$  is the noise that is added in latent space feature  $z'_0$  in each noising step,  $\epsilon_{\theta_1}$  is the denoising U-Net that predicts the noise  $\epsilon$  in current step  $t$ ,  $y$  stands for extra condition (*e.g.*, text, mask, *etc.*), and  $\tau_{\theta_2}$  is a domain specific encoder that projects  $y$  to intermediate representation. In this work, we add condition information, *i.e.*, the composite image with foreground mask, using an adaptive encoder similar to [39].

During inference, noise is first added to  $z'_0$  to generate  $z'_T$ , and then  $z'_T$  is used as  $z_T$ , the initial input for  $\epsilon_{\theta_1}$ .  $\epsilon_{\theta_1}$  is then iteratively used to estimate the noise at each denoising step  $t$ , thus the latent map  $z_T$  is gradually refined and ultimately becomes clean latent feature  $z_0$ . Finally, the clean latent feature  $z_0$  is fed into the decoder  $\mathcal{D}$  to generate the image. For more details about the training and inference of Stable Diffusion, please refer to [46].



**Figure 2: The architecture of our PHDiffusion.** Given an composite image  $I_c$  and its foreground mask  $M$ ,  $I_c$  is sent to a pretrained Stable Diffusion [7] model for painterly image harmonization. The input  $I_c$  is first encoded to the latent space  $z'_0 = \mathcal{E}(I_c)$ , which is followed by the forward process of diffusion to deduce  $z'_t; t = 1, 2, \dots, T$  with noise  $\epsilon$ . During inference,  $z'_T$  is used as the initial input  $z_T$  for the backward process to predict  $z_t; t = T - 1, T - 2, \dots, 0$  through the denoising U-Net. Finally, the harmonized image  $\tilde{I}_0$  is generated through the decoder by  $\tilde{I}_0 = \mathcal{D}(z_0)$ . In the meanwhile, the composite image  $I_c$  concatenated with foreground mask  $M$  is sent to the adaptive encoder, which is followed by the Dual Encoder Fusion to provide guidance to the denoising process in U-Net. The denoising process is supervised by the noise loss  $\mathcal{L}_{LDM}$ . Besides, we exploit two style losses ( $\mathcal{L}_{AdaIN}$  and  $\mathcal{L}_{CL}$ ) for foreground stylization and content loss ( $\mathcal{L}_{con}$ ) for content preservation.

### 3.2 Adaptive Encoder

As introduced before, the adaptive encoder accounts for encoding additional condition and providing multi-step guidance in the denoising steps. Previous implementations of the adaptive encoder [39] focus more on rough structures (e.g., sketch, pose, semantic mask), and exploit the text condition [46] to indicate extra demands (e.g., styles or environments). Different from previous works, we discard the text CLIP model and adopt the lightweight adaptive encoder [39] to encode the concatenation of composite image and foreground mask, simultaneously preserving content details and extracting background styles. In detail, the architecture of adaptive encoder comprises of four feature extraction blocks and three DownSample (DS for short) blocks. The input with resolution  $512 \times 512$  is first downsampled to  $64 \times 64$  (named as  $F_c^0$ ) through pixel unshuffle [48]. By combining one convolutional layer and two residual blocks as an extraction module (EM for short), the generation of  $F_c^i, i \in \{1, 2, 3, 4\}$  can be formulated as:

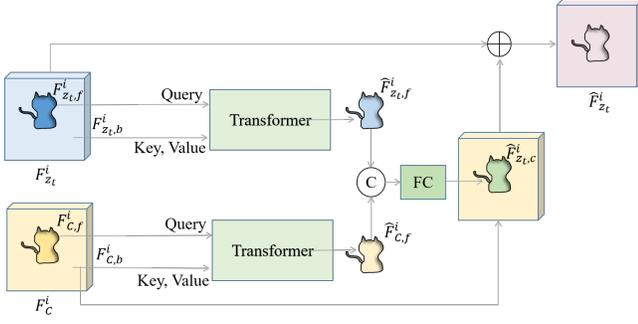
$$\begin{aligned} F_c^1 &= \text{EM}_1(F_c^0), & F_c^2 &= \text{EM}_2(\text{DS}(F_c^1)), \\ F_c^3 &= \text{EM}_3(\text{DS}(F_c^2)), & F_c^4 &= \text{EM}_4(\text{DS}(F_c^3)), \end{aligned} \quad (2)$$

where the resolutions of  $F_c^1, F_c^2, F_c^3$ , and  $F_c^4$  are  $64 \times 64, 32 \times 32, 16 \times 16$ , and  $8 \times 8$ , respectively. Similar structures also exist in U-Net encoder to generate  $F_{z_t}^i$  with the same resolution as  $F_c^i, i \in \{1, 2, 3, 4\}$ .

### 3.3 Dual Encoder Fusion Module

In the section, we introduce our Dual Encoder Fusion (DEF) module to preserve content of foreground object and extract reasonable style from background painting. As shown in Figure 2, the adaptive encoder first takes composite image  $I_c$  and foreground mask  $M$  as input, then generates the composite feature maps  $F_c^i, i \in \{1, 2, 3, 4\}$  with different resolutions. Each of the composite feature maps  $F_c^i$  is then fused with the corresponding denoising feature maps  $F_{z_t}^i$ . However, we find that if we directly add or concatenate the feature maps from these two encoders, the final generation result cannot stylize the foreground well, leading to notable style discrepancy between foreground and background.

Considering that CNN can only expand the receptive field of the foreground features within a certain range and struggles to capture long-range dependency [12], we propose to endow the foreground features with global receptive field to some extent. To balance the global-local receptive field of the foreground features, we design different fusion strategies for the features with different resolutions. For shallow features with high resolutions, where  $i \in \{1, 2\}$ , composite feature maps  $F_c^i$  are simply added to  $F_{z_t}^i$  to maintain the local structures. While for deeper features with low resolutions, where  $i \in \{3, 4\}$ ,  $F_c^i$  and  $F_{z_t}^i$  are fused through our DEF module to



**Figure 3: The architecture of Dual Encoder Fusion (DEF) module.** Given the composite feature map  $F_c^i$  (*resp.*, the denoising feature map  $F_{z_t}^i$ ), we stylize the foreground of composite feature map (*resp.*, the denoising feature map) by regarding the foreground feature map  $F_{c,f}^i$  (*resp.*,  $F_{z_t,f}^i$ ) as query and the background feature map  $F_{c,b}^i$  (*resp.*,  $F_{z_t,b}^i$ ) as key/value through a transformer layer to obtain the stylized foreground feature map  $\hat{F}_{c,f}^i$  (*resp.*,  $\hat{F}_{z_t,f}^i$ ). Then,  $\hat{F}_{c,f}^i$  and  $\hat{F}_{z_t,f}^i$  are fused through concatenation and a fully-connected layer, which is then combined with the composite background feature map  $F_{c,b}^i$  and denoising feature map  $F_{z_t}^i$  to further guide the denoising process.

capture the global styles. The above process can be formulated as

$$\hat{F}_{z_t}^i = \begin{cases} F_c^i + F_{z_t}^i, & i = 1, 2, \\ \text{DEF}(F_c^i, F_{z_t}^i), & i = 3, 4. \end{cases} \quad (3)$$

The structure of our DEF module is illustrated in Figure 3, which consists of stylized feature extraction and stylized feature fusion.

**3.3.1 Stylized Feature Extraction.** Before we fuse composite feature maps  $F_c^i$  and denoising feature maps  $F_{z_t}^i$ , we first need to expand the receptive field of the foreground features and extract the desired style from background. We utilize the foreground features to search for relevant background styles through a transformer layer [52]. In detail, by taking composite feature map  $F_c^i$  as example, we first extract its foreground features  $F_{c,f}^i$  and background features  $F_{c,b}^i$  by masking and flattening, which can be formulated as

$$\begin{aligned} F_{c,f}^i &= \text{Flatten}(F_c^i \circ \hat{M}), \\ F_{c,b}^i &= \text{Flatten}(F_c^i \circ (1 - \hat{M})), \end{aligned} \quad (4)$$

where  $\hat{M}$  denotes the foreground mask that is down-sampled to the corresponding size,  $\circ$  represents the element-wise product, and  $\text{Flatten}(\cdot)$  means the conversion from 2D feature map to 1D feature sequence. To search for the relevant background styles, we enhance the foreground features through a transformer layer, where  $F_{c,f}^i$  serve as queries and  $F_{c,b}^i$  serve as keys/values. The stylized composite foreground features  $\hat{F}_{c,f}^i$  can be represented by

$$\hat{F}_{c,f}^i = \text{Transformer}(F_{c,f}^i, F_{c,b}^i, F_{c,b}^i), \quad (5)$$

where Transformer is a transformer encoder layer [52].

Similar to  $F_c^i$ , the denoising feature map  $F_{z_t}^i$  can also be used to get the stylized denoising foreground features  $\hat{F}_{z_t,f}^i$ , so that the foreground features are stylized by relevant background styles.

**3.3.2 Stylized Feature Fusion.** After extracting  $\hat{F}_{c,f}^i$  and  $\hat{F}_{z_t,f}^i$  for  $i \in \{3, 4\}$ , we need to leverage both  $\hat{F}_{c,f}^i$  and  $\hat{F}_{z_t,f}^i$  from dual encoders to help the denoising process. In particular, we first concatenate  $\hat{F}_{c,f}^i$  and  $\hat{F}_{z_t,f}^i$ , and then pass them through a fully-connected layer to acquire the stylized foreground features  $\hat{F}_{z_t,c}^i$ . The stylized foreground features  $\hat{F}_{z_t,c}^i$  are then combined with the background of composite feature map  $F_c^i$  and the denoising feature map  $F_{z_t}^i$  to guide the denoising steps. The above steps could be formulated as

$$\hat{F}_{z_t,c}^i = \text{FC}(\hat{F}_{c,f}^i \oplus \hat{F}_{z_t,f}^i), \quad (6)$$

$$\hat{F}_{z_t}^i = F_{z_t}^i + \text{Fold}(\hat{F}_{z_t,c}^i) + F_c^i \circ (1 - \hat{M}), \quad (7)$$

where  $\text{FC}(\cdot)$  means the fully-connected layer,  $\oplus$  means the concatenation between two vectors,  $\text{Fold}(\cdot)$  means folding the  $\hat{F}_{z_t,c}^i$  into 2D foreground map.

## 3.4 Objective Function

First, we employ the standard noise loss from diffusion models [46], which aims to reconstruct the image feature within the latent space. However, merely using noise loss can only reconstruct the composite image without changing the foreground style. Therefore, we employ a combination of noise loss, AdaIN loss, and contrastive style loss, with the goal of attaining a balance between reasonable style and preservation of image structures/details. Besides, content loss is employed to assist in balancing noise loss and style losses. In the following, we will detail four losses one by one.

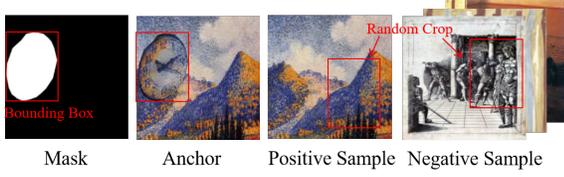
**3.4.1 Noise Loss.** The denoising step of DDPM [21] is to remove noise from  $z_t'$  step by step and finally reconstruct the original composite input  $z_0'$ . Therefore, the goal of the noise loss is to predict the noise in step  $t$ , which can be formulated as

$$\mathcal{L}_{LDM} := \mathbb{E}_{z_0', y, \epsilon \sim \mathcal{N}(0,1), t} \left[ \left\| \epsilon - \epsilon_{\theta_1} \left( z_t', t, \tau_{\hat{\theta}_2}(y) \right) \right\|_2^2 \right], \quad (8)$$

in which  $t$  is sampled from  $\{1, \dots, T\}$ ,  $y$  is the condition information (*i.e.*, composite image and foreground mask) in our problem.  $\epsilon_{\theta_1}$  includes the model parameters of denoising U-Net, while  $\tau_{\hat{\theta}_2}$  includes the model parameters of adaptive encoder and dual encoder fusion module.

**3.4.2 AdaIN Loss.** Note that the noise loss in Eqn. (8) is calculated in the latent space, whereas the style losses cannot be directly calculated in the latent space. Thus, we calculate two style losses based on the decoded image through decoder  $\mathcal{D}$ , giving  $\hat{I}_0 = \mathcal{D}(z_0')$ .  $\hat{z}_0^t = \left( z_t' - \sqrt{1 - \bar{\alpha}_t} \epsilon_{\theta_1} \left( z_t', t, \tau_{\hat{\theta}_2}(y) \right) \right) / \sqrt{\bar{\alpha}_t}$  [21], in which  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ ,  $\alpha_s = 1 - \beta_s$  and  $\beta_s$  represents forward process variances. Based on  $\hat{I}_0$ , we can easily calculate style losses.

We utilize the AdaIN style loss [14] to achieve consistency in multi-scale feature statistics (*i.e.*, mean, standard deviation) between the background painting and the foreground in the harmonized



**Figure 4: Construction of a triplet of anchor, positive sample, and negative sample for contrastive style loss.**

image  $\hat{I}_0$ . The style loss can be written as

$$\mathcal{L}_{AdaIN} = \sum_{l=1}^L \left\| \mu \left( \phi^l \left( \hat{I}_0 \right) \circ \bar{M}^l \right) - \mu \left( \phi^l \left( I_b \right) \right) \right\|_2^2 + \sum_{l=1}^L \left\| \sigma \left( \phi^l \left( \hat{I}_0 \right) \circ \bar{M}^l \right) - \sigma \left( \phi^l \left( I_b \right) \right) \right\|_2^2, \quad (9)$$

where  $\phi^l, l \in \{1, 2, 3, 4\}$  represents the  $l$ -th ReLU <sub>$l-1$</sub>  layer in a pre-trained VGG-19 [50] network.  $I_b$  is the complete background painting and  $\bar{M}^l$  denotes the foreground mask that is down-sampled to the corresponding size.  $\mu(\cdot)$  means the calculation of mean value and  $\sigma(\cdot)$  means the calculation of standard deviation.

**3.4.3 Contrastive Style Loss.** To better migrate background style to foreground, we introduce another contrastive style loss, which is complementary with AdaIn loss. Contrastive style loss was first introduced into style transfer task by [5], which distinguishes the image rendered by the reference style from other styles. Here, we adapt contrastive style loss to our painterly image harmonization task. Specifically, we construct a triplet of three elements: anchor, positive sample, and negative sample, as shown in Figure 4. To acquire anchor, we first feed the harmonized output  $\hat{I}_0$  into pre-trained VGG-19 network and extract the output feature map of ReLU<sub>3\_1</sub> layer. Then, we crop the feature map with the downsampled mask followed by average pooling to obtain the foreground feature, which is projected to the anchor vector  $f_q$ . Through the same procedure, positive sample  $f_b^+$  is extracted from the background painting  $I_b$ , so that  $f_q$  and  $f_b^+$  share the same style. And negative samples  $f_b^-$  are extracted from other style images. Given the triplet, we tend to pull close the anchor and positive example while separating the anchor from negative examples, which can be represented as

$$\mathcal{L}_{CL} = -\log \left( \frac{\exp \left( \left( f_q \right)^T \left( f_b^+ \right) / \eta \right)}{\exp \left( \left( f_q \right)^T \left( f_b^+ \right) / \eta \right) + \sum f_b^- \exp \left( \left( f_q \right)^T \left( f_b^- \right) / \eta \right)} \right), \quad (10)$$

where the temperature  $\eta$  regulates the push and pull forces. We set  $\eta$  as 0.2 following [5].

**3.4.4 Content Loss.** In addition, when balancing between noise loss and style losses, chances are that the content details are excessively preserved, leading to insufficient style transfer. Therefore, we reduce the weight of noise loss and incorporate content loss [14], which is commonly used in style transfer tasks. The content loss

can help preserve the high-level content information without the sacrifice of styles. The content loss can be written as

$$\mathcal{L}_{con} = \left\| \phi^4 \left( \hat{I}_0 \right) - \phi^4 \left( I_c \right) \right\|_2^2, \quad (11)$$

where  $\phi^4$  has been defined below Eqn. (9).

**3.4.5 Total Loss.** By summarizing the noise loss, two style losses, and content loss, the total loss can be written as

$$\mathcal{L}_{total} = \lambda_1 \mathcal{L}_{LDM} + \mathcal{L}_{AdaIN} + \lambda_2 \mathcal{L}_{CL} + \mathcal{L}_{con}, \quad (12)$$

in which  $\lambda_1$  and  $\lambda_2$  are hyper-parameters. We empirically set them as 60 and 5 respectively.

## 4 EXPERIMENTS

### 4.1 Experiment Settings

We train the adaptive encoder and the fusion module for 10 epochs with a batch size of 2. We utilize Adam as the optimizer with the learning rate of  $2 \times 10^{-4}$ . During training, we resize the input images and the mask to  $512 \times 512$  and use the pretrained Stable Diffusion model [7] with the version of sd-v1-4. We utilize the training data of COCO [32] and WikiArt [41]. For more implementation details, please refer to the supplementary.

### 4.2 Baselines

Based on the target task, existing baselines can be categorized into three groups: painterly image harmonization [3, 36, 58], cross-domain composition methods [18, 38], and artistic style transfer methods [10, 22, 35, 43, 59]. The painterly image harmonization methods include DIB [58], DPH [36], and PHDNet [3]. The cross-domain composition methods include CDC [18] and SDEdit [38]. The artistic style transfer methods include AdaIN [22], AdaAttN [35], SANet [43], StyTr2 [10], and InST [59]. Among them, CDC [18], SDEdit [38], and InST [59] are diffusion-based methods.

For the first and the second groups, these works can stylize a certain region, so we directly compare them with our results. However, for the third group, these works stylize the entire photographic image. To adapt artistic style transfer methods to our task, we stylize the content image according to the background image, followed by cutting and pasting the stylized foreground object onto the background image. We set *Strength* to 0.7 by default to control total inference steps for our model. More details of implementations including hyper-parameters of baselines are in the supplementary.

### 4.3 Comparisons with Baselines

**4.3.1 Visualization Analysis.** Compared with the first group of baselines, painterly harmonization methods, we can refer to Figure 5 for the visualization results. It is shown that our PHDiffusion can endow the foreground with more abundant and coherent styles (row 1, 2, 3, 4). For example, as illustrated in row 1, our PHDiffusion can not only learn the local dotted textures, but also learn the global stripe arrangement. Besides, our method can strike great balance between content and style. For content preservation, our method holds more semantic edge information (row 1, 2, 3), and also maintains more refined details (row 1). In row 3, our method preserves the clear umbrella frame while capturing the textures

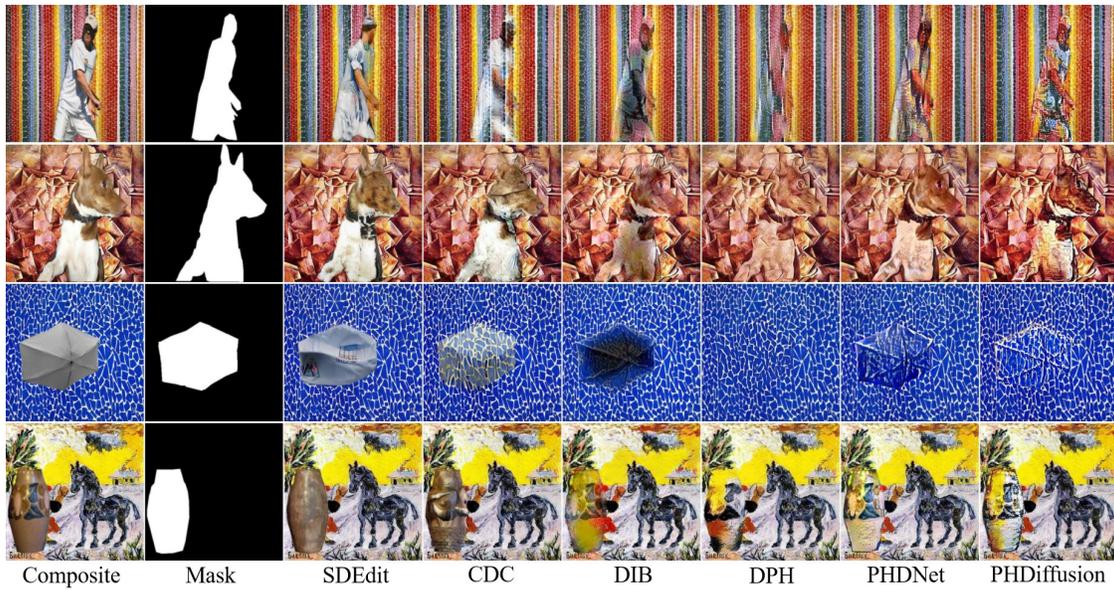


Figure 5: From left to right, we show the composite image, mask, harmonized results of SDEdit, CDC, DIB, DPH, PHDNet and our PHDiffusion. Best viewed in color and zoom in.

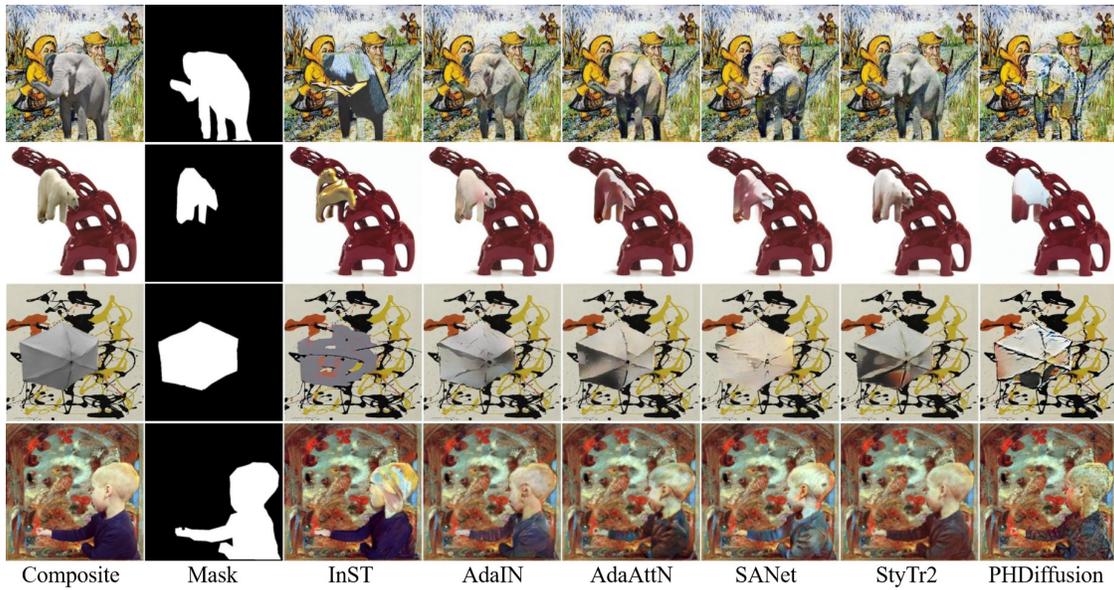


Figure 6: From left to right, we show the composite image, mask, example results for InST, AdaIN, AdaAttN, SANet, StyTr2 and our PHDiffusion. Best viewed in color and zoom in.

	AdaIN	AdaATTN	SANet	StyTr2	InST	DIB	DPH	PHDNet	SDEdit	CDC	PHDiffusion
BT	-0.224	0.180	0.596	0.912	-1.779	-1.123	1.376	1.811	-2.476	-1.863	2.590

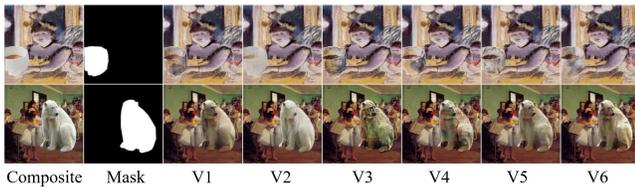
Table 1: Comparisons with baselines. "BT" stands for B-T score.

closest to the background. However, DPH loses details and semantic edge, while DIB and PHDNet fail to transfer coherent textures.

Compared with the second group, cross-domain composition methods, as illustrated in Figure 5 (left), we can see that SDEdit

Version	Method	BT
V1	w/o $\mathcal{L}_{CL}$	0.636
V2	w/o $\mathcal{L}_{AdaIN}$	-1.584
V3	w/o TA,TU	-0.421
V4	w/o TA	0.182
V5	w/o TU	-0.112
V6	full	1.300

**Table 2: The results of the ablation experiments. "TA" and "TU" stand for the transformer layer in adaptive encoder and U-Net encoder respectively. "BT" stands for B-T score.**



**Figure 7: Examples of ablation experiments. Best viewed in color and zoom in.**

and CDC struggle to transfer sufficient style while keeping original content (row 1, 3, 4). In row 1, the content has already changed sharply while the style fails to be adapted to the target. Though CDC has learnt some dotted textures, the style and content exhibit a lack of cohesion, resembling distinct layers. And in row 3 and row 4, the content is lost to some extent with details unrecognized. The poor performance of baselines is probably caused by the stochastic nature of diffusion model, which can be a double-edged sword for the tasks that require handling delicate details, since it is very hard to adjust hyper-parameters, (e.g., strength) to balance between style and content without proper guidance. In contrast, our method provides more effective guidance for the denoising process.

Compared with the third group, artistic transfer methods, as shown in Figure 6, it can be seen that InST loses content and exhibits style incompatible with the background (row 1, 4). Other style transfer methods can not produce adequate styles. Our method not only produces textures that highly match the background (row 1, 2, 3, 4), but also enables the overall color distribution to be strongly correlated with the background (row 1, 3), leading to better visual harmony.

The great performance of our PHDiffusion is attributed to two aspects. Firstly, for producing sufficient styles, our DEF module can query backgrounds for reasonable styles and utilize prior knowledge in pretrained stable diffusion model. Secondly, for balancing content and style, the combination of noise loss, content loss, and style losses enable the adaptive encoder and DEF module to store appropriate guiding information.

**4.3.2 User Study.** We also conduct a user study to compare the effectiveness of various methods, following [3]. Specifically, we randomly select 100 content images from COCO [32] and 100 style images from WikiArt [41] to generate 100 composite images. Given

each composite image, we can obtain 11 harmonized results including 10 baselines and our method. Then pairwise comparisons are conducted, resulting in 5,500 image pairs. We invite 50 users to identify the more harmonious one in each pair. Finally 275,000 comparison results are collected, followed by using the Bradley-Terry (B-T) model [1, 27] to calculate an overall ranking of all methods. As presented in Table 1, our PHDiffusion achieves the highest B-T score.

#### 4.4 Ablation Studies

As described in Section 3, our PHDiffusion exploits an adaptive encoder along with the dual encoder fusion module to guide the denoising process and two style losses to balance the content and style. Therefore, in this section, we demonstrate their effectiveness, and report B-T score in Table 2 and visual results in Figure 7. For the effectiveness of style losses, we conduct experiments without contrastive style loss (V1) or AdaIN loss (V2). For the effectiveness of transformer layer in DEF module, we conduct experiments in the following three settings: (1) remove transformer layers in both encoders (V3); (2) remove transformer layer in adaptive encoder (V4); (3) remove transformer layer in U-Net encoder (V5).

Comparing V1, V2, and V6 in Table 2, we find that the AdaIN loss is more important for style transfer, while the contrastive style loss assists in capturing more reasonable styles. Moreover, by comparing V3, V4, V5, and V6 in Table 2, we can find that the transformer layers in U-Net encoder and adaptive encoder are both helpful to generate reasonable styles, and exploiting transformer layer in both encoders can further boost the painterly image harmonization.

For visual results in Figure 7, comparing V1, V2, and V6, it is observed that AdaIN loss (V1) can learn styles that appear to be blended, while contrastive style loss (V2) tends to learn fine textures (more haziness for cup and more textures of the fur for bear) while maintaining the original color. So the combination of two style losses helps the transformer capture adequate local textures and fine global styles. Comparing V3 and V4 (V3 and V5) in Figure 7, we can find that the transformer can help learn more consistent and reasonable styles from the background. Comparing V4, V5, and V6 in Figure 7, it is observed that the adaptive encoder prefers more subdued color (V5) while the U-Net encoder tends to perform more exaggerated color (V4). Balancing them can achieve more reasonable and harmonized styles for our final result (V6).

## 5 CONCLUSION

In this work, we have introduced diffusion model into painterly image harmonization. We have proposed a novel Painterly Harmonization stable Diffusion model (PHDiffusion), in which the denoising process of diffusion is under the guidance of lightweight adaptive encoder and dual encoder fusion. Experiments have demonstrated that our approach can simultaneously preserve detailed content and produce sufficient styles, surpassing the state-of-the-art methods.

## ACKNOWLEDGMENTS

The work was supported by the Shanghai Municipal Science and Technology Major / Key Project, China (Grant No. 20511100300 / 2021SHZDZX0102) and the National Natural Science Foundation of China (Grant No. 62076162).

## REFERENCES

- [1] Ralph Allan Bradley and Milton E Terry. 1952. Rank analysis of incomplete block designs: I. The method of paired comparisons. *Biometrika* 39, 3/4 (1952), 324–345.
- [2] Junyan Cao, Wenyan Cong, Li Niu, Jianfu Zhang, and Liqing Zhang. 2022. Deep Image Harmonization by Bridging the Reality Gap. *BMVC* (2022).
- [3] Junyan Cao, Yan Hong, and Li Niu. 2023. Painterly Image Harmonization in Dual Domains. *AAAI* (2023).
- [4] Bor-Chun Chen and Andrew Kae. 2019. Toward realistic image compositing with adversarial learning. In *CVPR*.
- [5] Haibo Chen, Zhizhong Wang, Huiming Zhang, Zhiwen Zuo, Ailin Li, Wei Xing, Dongming Lu, et al. 2021. Artistic style transfer with internal-external learning and contrastive learning. In *NeurIPS*.
- [6] Wenyan Cong, Li Niu, Jianfu Zhang, Jing Liang, and Liqing Zhang. 2021. Bargainnet: Background-guided domain translation for image harmonization. In *ICME*.
- [7] Wenyan Cong, Xinhao Tao, Li Niu, Jing Liang, Xuesong Gao, Qihao Sun, and Liqing Zhang. 2022. High-resolution image harmonization via collaborative dual transformations. In *CVPR*.
- [8] Wenyan Cong, Jianfu Zhang, Li Niu, Liu Liu, Zhixin Ling, Weiyuan Li, and Liqing Zhang. 2020. Dovenet: Deep image harmonization via domain verification. In *CVPR*.
- [9] Xiaodong Cun and Chi-Man Pun. 2020. Improving the harmony of the composite image by spatial-separated attention module. *IEEE Trans. Image Process.* 29 (2020), 4759–4771.
- [10] Yingying Deng, Fan Tang, Weiming Dong, Chongyang Ma, Xingjia Pan, Lei Wang, and Changsheng Xu. 2022. Stytr2: Image style transfer with transformers. In *CVPR*.
- [11] Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Darrell. 2014. Decaf: A deep convolutional activation feature for generic visual recognition. In *ICML*.
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiuhua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
- [13] Len Du. 2020. How much deep learning does neural style transfer really need? an ablation study. In *WACV*.
- [14] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. 2016. Image style transfer using convolutional neural networks. In *CVPR*.
- [15] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative adversarial networks. *Commun. ACM* 63, 11 (2020), 139–144.
- [16] Zonghui Guo, Dongsheng Guo, Haiyong Zheng, Zhaorui Gu, Bing Zheng, and Junyu Dong. 2021. Image harmonization with transformer. In *ICCV*.
- [17] Zonghui Guo, Haiyong Zheng, Yufeng Jiang, Zhaorui Gu, and Bing Zheng. 2021. Intrinsic image harmonization. In *CVPR*.
- [18] Roy Hachnochi, Mingrui Zhao, Nadav Orzech, Rinon Gal, Ali Mahdavi-Amiri, Daniel Cohen-Or, and Amit Haim Bermanto. 2023. Cross-domain Compositing with Pretrained Diffusion Models. *arXiv preprint arXiv:2302.10167* (2023).
- [19] Yucheng Hang, Bin Xia, Wenming Yang, and Qingmin Liao. 2022. Scs-co: Self-consistent style contrastive learning for image harmonization. In *CVPR*.
- [20] Guoqing Hao, Satoshi Iizuka, and Kazuhiro Fukui. 2020. Image Harmonization with Attention-based Deep Feature Modulation. In *BMVC*.
- [21] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *NeurIPS*.
- [22] Xun Huang and Serge Belongie. 2017. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*.
- [23] Jaeseok Jeong, Mingi Kwon, and Youngjung Uh. 2023. Training-free Style Transfer Emerges from h-space in Diffusion models. *arXiv preprint arXiv:2303.15403* (2023).
- [24] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. 2016. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*.
- [25] Nicholas Kolkin, Jason Salavon, and Gregory Shakhnarovich. 2019. Style transfer by relaxed optimal transport and self-similarity. In *CVPR*.
- [26] Gihyun Kwon and Jong Chul Ye. 2022. Diffusion-based image translation using disentangled style and content representation. *arXiv preprint arXiv:2209.15264* (2022).
- [27] Wei-Sheng Lai, Jia-Bin Huang, Zhe Hu, Narendra Ahuja, and Ming-Hsuan Yang. 2016. A comparative study for single image blind deblurring. In *CVPR*.
- [28] Chuan Li and Michael Wand. 2016. Precomputed real-time texture synthesis with markovian generative adversarial networks. In *ECCV*.
- [29] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. 2017. Universal style transfer via feature transforms. In *NeurIPS*.
- [30] Yanghao Li, Naiyan Wang, Jiaying Liu, and Xiao Di Hou. 2017. Demystifying neural style transfer. *arXiv preprint arXiv:1701.01036* (2017).
- [31] Jingtang Liang, Xiaodong Cun, Chi-Man Pun, and Jue Wang. 2022. Spatial-separated curve rendering network for efficient and high-resolution image harmonization. In *ECCV*.
- [32] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *ECCV*.
- [33] Jun Ling, Han Xue, Li Song, Rong Xie, and Xiao Gu. 2021. Region-aware adaptive instance normalization for image harmonization. In *CVPR*.
- [34] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. 2022. Pseudo numerical methods for diffusion models on manifolds. *arXiv preprint arXiv:2202.09778* (2022).
- [35] Songhua Liu, Tianwei Lin, Dongliang He, Fu Li, Meiling Wang, Xin Li, Zhengxing Sun, Qian Li, and Errui Ding. 2021. Adaattn: Revisit attention mechanism in arbitrary neural style transfer. In *ICCV*.
- [36] Fujun Luan, Sylvain Paris, Eli Shechtman, and Kavita Bala. 2018. Deep painterly harmonization. In *Comput Graph Forum*.
- [37] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. 2022. Repaint: Inpainting using denoising diffusion probabilistic models. In *CVPR*.
- [38] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. 2021. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *ICLR*.
- [39] Chong Mou, Xintao Wang, Liangbin Xie, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaoou Qie. 2023. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453* (2023).
- [40] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2021. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741* (2021).
- [41] Kiri Nichol. 2016. *Painter by numbers*. <https://www.kaggle.com/c/painter-by-numbers/overview>
- [42] Li Niu, Wenyan Cong, Liu Liu, Yan Hong, Bo Zhang, Jing Liang, and Liqing Zhang. 2021. Making Images Real Again: A Comprehensive Survey on Deep Image Composition. *arXiv preprint arXiv:2106.14490* (2021).
- [43] Dae Young Park and Kwang Hee Lee. 2019. Arbitrary style transfer with style-attentional networks. In *CVPR*.
- [44] Hwai-Jin Peng, Chia-Ming Wang, and Yu-Chiang Frank Wang. 2019. Element-Embedded Style Transfer Networks for Style Harmonization. In *BMVC*.
- [45] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *ICML*.
- [46] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *CVPR*.
- [47] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. 2022. Palette: Image-to-image diffusion models. In *SIGGRAPH*.
- [48] Wenze Shi, Jose Caballero, Ferenc Huszar, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. 2016. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *CVPR*.
- [49] Alon Shoshan, Nadav Bthonker, Igor Kviatkovsky, and Gerard Medioni. 2021. Gan-control: Explicitly controllable gans. In *ICCV*.
- [50] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [51] Jiaming Song, Chenlin Meng, and Stefano Ermon. 2020. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502* (2020).
- [52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *NeurIPS* (2017).
- [53] Huikai Wu, Shuai Zheng, Junge Zhang, and Kaiqi Huang. 2019. Gp-gan: Towards realistic high-resolution image blending. In *ACM MM*.
- [54] Yazhou Xing, Yu Li, Xintao Wang, Ye Zhu, and Qifeng Chen. 2022. Composite photograph harmonization with complete background cues. In *ACM MM*.
- [55] Ben Xue, Shenghui Ran, Quan Chen, Rongfei Jia, Binqiang Zhao, and Xing Tang. 2022. Dccf: Deep comprehensible color filter learning framework for high-resolution image harmonization. In *ECCV*.
- [56] Xiao Yan, Yang Lu, Juncheng Shuai, and Sanyuan Zhang. 2022. Style Image Harmonization via Global-Local Style Mutual Guided. In *ACCV*.
- [57] Lvmin Zhang and Maneesh Agrawala. 2023. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543* (2023).
- [58] Lingzhi Zhang, Tarmily Wen, and Jianbo Shi. 2020. Deep image blending. In *CVPR*.
- [59] Yuxin Zhang, Nisha Huang, Fan Tang, Haibin Huang, Chongyang Ma, Weiming Dong, and Changsheng Xu. 2022. Inversion-Based Creativity Transfer with Diffusion Models. *arXiv preprint arXiv:2211.13203* (2022).
- [60] Yabin Zhang, Minghan Li, Ruihuang Li, Kui Jia, and Lei Zhang. 2022. Exact feature distribution matching for arbitrary style transfer and domain generalization. In *CVPR*.
- [61] Yuxin Zhang, Fan Tang, Weiming Dong, Haibin Huang, Chongyang Ma, Tong-Yee Lee, and Changsheng Xu. 2022. Domain enhanced arbitrary image style transfer

via contrastive learning. In *SIGGRAPH*.  
[62] Jun-Yan Zhu, Philipp Krahenbuhl, Eli Shechtman, and Alexei A Efros. 2015. Learning a discriminative model for the perception of realism in composite

images. In *ICCV*.

# Supplementary Material for Painterly Image Harmonization using Diffusion Model

Lingxiao Lu  
MoE Key Lab of Artificial Intelligence,  
Shanghai Jiao Tong University  
China  
lulingxiao@sjtu.edu.cn

Jiangtong Li  
MoE Key Lab of Artificial Intelligence,  
Shanghai Jiao Tong University  
China  
keep\_moving-lee@sjtu.edu.cn

Junyan Cao  
MoE Key Lab of Artificial Intelligence,  
Shanghai Jiao Tong University  
China  
joy\_c1@sjtu.edu.cn

Li Niu\*  
MoE Key Lab of Artificial Intelligence,  
Shanghai Jiao Tong University  
China  
ustcnewly@sjtu.edu.cn

Liqing Zhang\*  
MoE Key Lab of Artificial Intelligence,  
Shanghai Jiao Tong University  
China  
zhang-lq@cs.sjtu.edu.cn

## ACM Reference Format:

Lingxiao Lu, Jiangtong Li, Junyan Cao, Li Niu, and Liqing Zhang. 2023. Supplementary Material for Painterly Image Harmonization using Diffusion Model. In *Proceedings of the 31st ACM International Conference on Multimedia (MM '23)*, October 29–November 3, 2023, Ottawa, ON, Canada. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3581783.3612451>

In the supplementary, we will first introduce the dataset and implementation details in Appendix A. Then the hyper-parameter *Strength* will be studied for style strength control in Appendix B. The visualization of attention maps in the transformers of our dual encoder fusion module will be explained in Appendix C. We will also provide more details of implementing baselines and offer more visual comparison results in Appendix D. Finally, we will discuss the limitations of our method in Appendix E.

## A DATASET AND IMPLEMENTATION DETAILS

We conduct experiments on two benchmark datasets, *i.e.*, COCO [5] and WikiArt [9], where COCO is a large-scale photograph dataset with the instance segmentation annotation for 80 different object categories and WikiArt is a large-scale digital art dataset consisting of 27 distinct styles. These two datasets are used to produce composite images by inserting photographic foreground objects from COCO into painterly backgrounds from WikiArt.

In detail, to obtain the foreground object with proper size and resolution, we select 9,100 foreground images from the COCO dataset, whose foreground ratio is between 0.05 and 0.3, and width and height are  $\geq 480$ . Moreover, we select 37,931 background images from the WikiArt dataset, whose width and height are  $\geq 512$ . During training, we use instance annotation to extract the foreground

objects from the foreground images and then place it onto a randomly chosen painterly background from the background images, leading to 37,931 composite images in each epoch. Finally, all the composite image are resized to  $512 \times 512$  for training. This process can produce composite images with discordant visual elements.

Our network is implemented using Pytorch 1.11.0. And the training process is executed on an Ubuntu 20.04 LTS operating system, utilizing a computing environment comprising of 32GB memory, Intel Xeon Silver 4116 CPU, and two GeForce RTX 3090 GPUs.

## B STRENGTH CONTROL

*Strength* is a hyper-parameter that decides the total step in the inference process. For example, the total step is equal to 35 when *Strength* is 0.7 and default total step is 50. The larger total step means more noise to be added and removed, leading to larger variability, so that the guidance from condition information could have greater impact on the harmonized results.

We observe that when *Strength* grows larger in a proper range (*i.e.*, 0.1 - 0.7), the style of our harmonized result gets transferred gradually while the content details are barely changed. However, as the *Strength* changes, other diffusion model baselines cannot balance the style and content. In Figure 1, we visualize how the harmonization results changes as the strength changes. In detail, we compare the harmonization results of SDEdit [8], CDC [3], InST [12], and our PHDiffusion with the *Strength* ranging from 0.1 to 0.9. Recall that the inference process of diffusion model is under control of the *Strength*, where the smaller the *Strength* is, the smaller the denoising step is. If the denoising process is guided improperly or without guidance, the style and the content cannot be balanced as the *Strength* changes. In detail, for SDEdit [8] and InST [12] in Figure 1, as the *Strength* becomes larger, the style is more sufficient while the content is destroyed. For CDC [3], since the denoising process is guided by the composite image, as the *Strength* gets larger, the content details are more preserved while the style is ignored; however, as the *Strength* gets smaller, the style is more sufficient while the content is destroyed. Besides, the balanced point for CDC [3] is also quite vulnerable.

However, if this process is guided by our DEF module, it can be tailored to painterly image harmonization with smooth transition between original and target styles without destroying content.

\*Corresponding authors

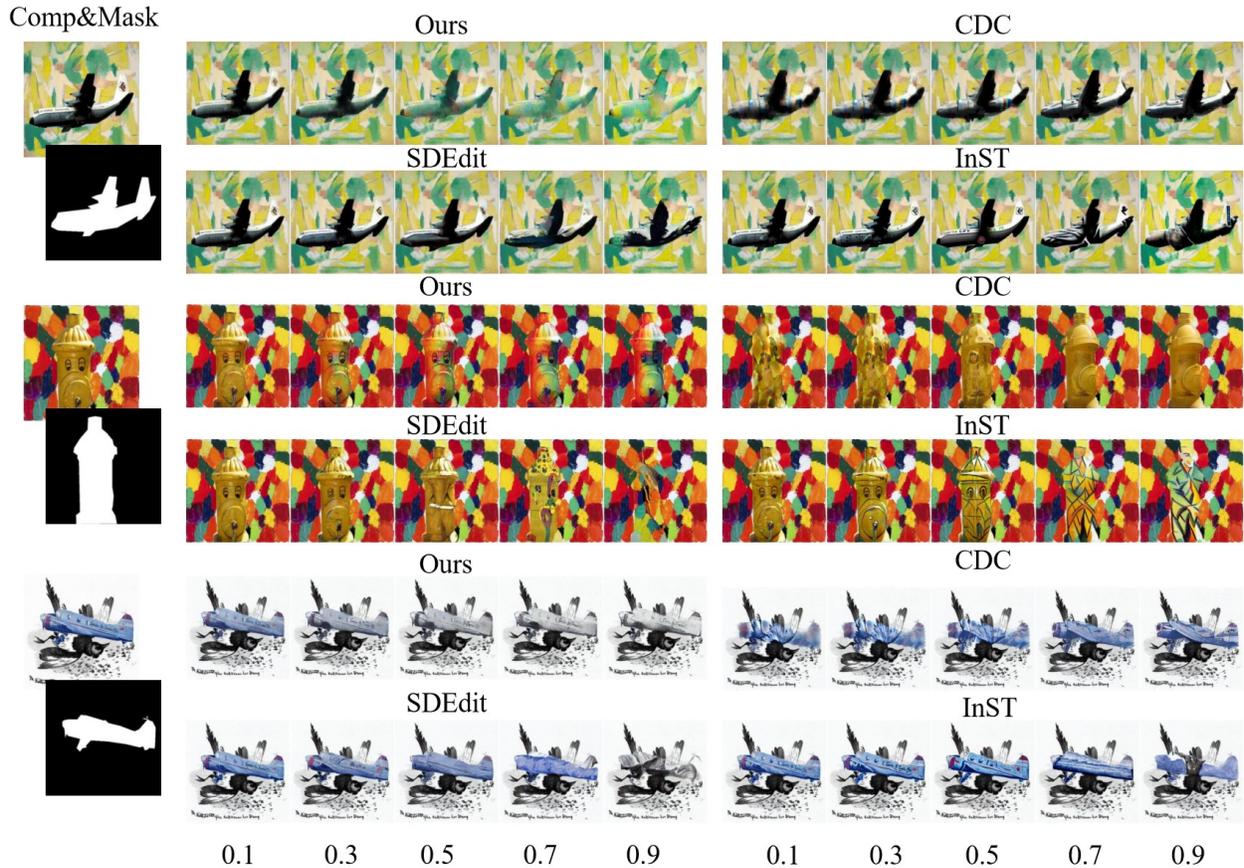
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

MM '23, October 29–November 3, 2023, Ottawa, ON, Canada

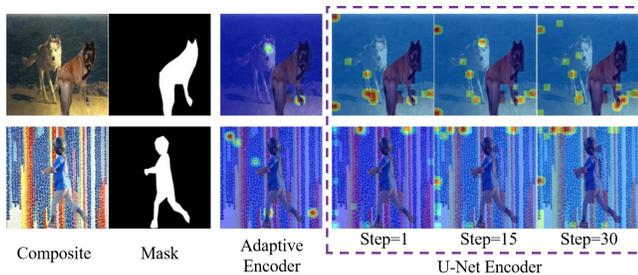
© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0108-5/23/10...\$15.00

<https://doi.org/10.1145/3581783.3612451>



**Figure 1: The results of adjusting *Strength* for controlling the degree of style transfer. *Strength* is set to 0.1, 0.3, 0.5, 0.7, 0.9 from left to right for each method. We also present the results of Ours, CDC, SDEdit, and InST for comparison.**



**Figure 2: The attention maps in DEF module for Adaptive Encoder and U-Net Encoder. For U-Net Encoder, we present its attention maps in different timesteps.**

Specifically, from our harmonization results in Figure 1, we observe that the styles are progressively strengthened while content is well-preserved (details such as eyes in the second example and words on the airplane in the third example remain visually clear). This proves that diffusion process is highly controllable and our mechanism can provide powerful guidance for the diffusion process.

## C VISUALIZATION OF ATTENTION MAPS FOR DUAL ENCODER FUSION

To better understand the transformers in DEF module, we visualize the attention maps of transformers during inference in Figure 2. Since the feature maps in adaptive encoder remain the same during multi-step inference, so its attention maps also remain the same. In contrast, the feature maps in U-Net encoder are updated during multi-step inference, so the attention maps vary in different steps. Therefore, in Figure 2, we show how the attention maps in U-Net encoder changes during denoising steps. Specifically, each attention map is obtained by averaging the attention maps from all attention heads, and then resized to the original resolution. In Figure 2, it can be observed that, for composite image whose background has similar objects with foreground, these objects can be detected and attended by the transformer layer, so that these objects in the background usually have larger weights. For example, in row 1 of Figure 2, the wolf in the background, which is similar to the foreground dog, gains attention in both adaptive encoder and U-Net encoder. For U-Net encoder in different steps, the wolves are all attended. Moreover, the ground surrounding the dog, which has similar color to the dog, is also attended by our DEF module.



Figure 3: Example failure cases of our PHDiffusion.

Besides, for composite image whose background has pure textures (row 2), the DEF module seems to pay attention to the background randomly for both adaptive and U-Net encoders to capture the overall pattern. These visualization results again prove that our DEF module can focus on meaningful background regions and provide rational guidance during the denoising steps.

## D COMPARISON WITH BASELINES

### D.1 Details of Diffusion-based Baselines

For all the diffusion-based baselines, since the *Strength* has great influence on the harmonization results, we adjust the *Strength* and select the optimal outcome for comparison. Therefore, we set *Strength* to 0.5, 0.5 and 0.7 for SDEdit [8], CDC [3] and InST [12], respectively, which are the best results for balancing content and style. For our PHDiffusion, we choose *Strength* = 0.7 by default.

Besides, to adapt style transfer method InST [12] for painterly image harmonization, for training process, we train the proposed small conditional network on WikiArt while freezing the diffusion model, after which the model is able to learn styles in WikiArt better. Moreover, during inference, we first add noise for  $T$  steps to the composite image in forward process. During each backward step, we exploit the corresponding background in  $i$ -th forward step to replace the background of predicted images in  $(T - i)$ -th background step, which aims to preserve the background and only adapt the foreground to satisfy the demand of painterly image harmonization.

### D.2 More Visual Results

We provide more visual results to compare with other baselines. As we have introduced in the main submission, we have three groups of baselines. The first group contains painterly image harmonization methods, DIB [11], DPH [7], and PHDNet [1]. The second group includes cross-domain composition methods, CDC [3] and SDEdit [8]. And the third group, artistic style transfer, consists of AdaIN [4], AdaAttN [6], SANet [10], StyTr2 [2], and InST [12]. The results for the first and second group are shown in Figure 4, while the results for the third group are shown in Figure 5.

In Figure 4, it can be seen that DIB [11], DPH [7] and PHDNet [1] can also achieve harmonization to some extent, but the learnt textures are not as accurate as ours (row 1, 3, 4, 5, 10). Besides, our PHDiffusion can capture more global styles, thus tending to be more naturally blended with the background (row 2, 6, 7, 8).

Specifically, our PHDiffusion is able to maintain more semantic information and content details (e.g., the stripes on the body of the cat in row 8 and the pattern on the shirt of the man in row 7). And for cross-domain methods in Figure 4, it is obvious that SDEdit [8] tends to directly copy the content of foreground objects (row 2, 3, 6, 9, 10), and the style is not compatible with the background images. Besides, CDC [3] fails to preserve the content details (row 1, 3, 5, 9, 10). The comparison between our PHDiffusion and other cross-domain methods (i.e., SDEdit [8] and CDC [3]) turns out that our method outperforms them in both style and content.

For the third group in Figure 5, it can be observed that InST [12] has lost too much content details. Though its learnt style is quite consistent with the background, the harmonization results from InST [12] are still less realistic (row 1, 4, 6). Moreover, for AdaIN [4], AdaAttN [6], SANet [10], and StyTr2 [2], they can also partially migrate the background style, however, the styles in our results are obviously more harmonized (row 3 to row 9). For various types of backgrounds, our harmonization results always behave well in holding semantic information (row 3, 7, 9). Overall, for our harmonization results, the inserted foreground objects can be better integrated into the background, making the whole harmonized images appear to be intact artistic paintings.

## E LIMITATIONS

Generally speaking, our method is capable of producing visually appealing and harmonious results, however, some types of foreground objects such as human faces are still hard to be in great harmony with the backgrounds. Since human faces have delicate details and we are very sensitive to the subtle changes in human faces, it is very difficult to sufficiently stylize the human faces while preserving their delicate details.

## REFERENCES

- [1] Junyan Cao, Yan Hong, and Li Niu. 2023. Painterly Image Harmonization in Dual Domains. *AAAI* (2023).
- [2] Yingying Deng, Fan Tang, Weiming Dong, Chongyang Ma, Xingjia Pan, Lei Wang, and Changsheng Xu. 2022. Stytr2: Image style transfer with transformers. In *CVPR*.
- [3] Roy Hachnochi, Mingrui Zhao, Nadav Orzech, Rinon Gal, Ali Mahdavi-Amiri, Daniel Cohen-Or, and Amit Haim Bermano. 2023. Cross-domain Compositing with Pretrained Diffusion Models. *arXiv preprint arXiv:2302.10167* (2023).
- [4] Xun Huang and Serge Belongie. 2017. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*.
- [5] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *ECCV*.
- [6] Songhua Liu, Tianwei Lin, Dongliang He, Fu Li, Meiling Wang, Xin Li, Zhengxing Sun, Qian Li, and Errui Ding. 2021. Adaattn: Revisit attention mechanism in arbitrary neural style transfer. In *ICCV*.
- [7] Fujun Luan, Sylvain Paris, Eli Shechtman, and Kavita Bala. 2018. Deep painterly harmonization. In *Comput Graph Forum*.
- [8] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. 2021. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *ICLR*.
- [9] Kiri Nichol. 2016. *Painter by numbers*. <https://www.kaggle.com/c/painter-by-numbers/overview>
- [10] Dae Young Park and Kwang Hee Lee. 2019. Arbitrary style transfer with style-attentional networks. In *CVPR*.
- [11] Lingzhi Zhang, Tarmily Wen, and Jianbo Shi. 2020. Deep image blending. In *CVPR*.
- [12] Yuxin Zhang, Nisha Huang, Fan Tang, Haibin Huang, Chongyang Ma, Weiming Dong, and Changsheng Xu. 2022. Inversion-Based Creativity Transfer with Diffusion Models. *arXiv preprint arXiv:2211.13203* (2022).

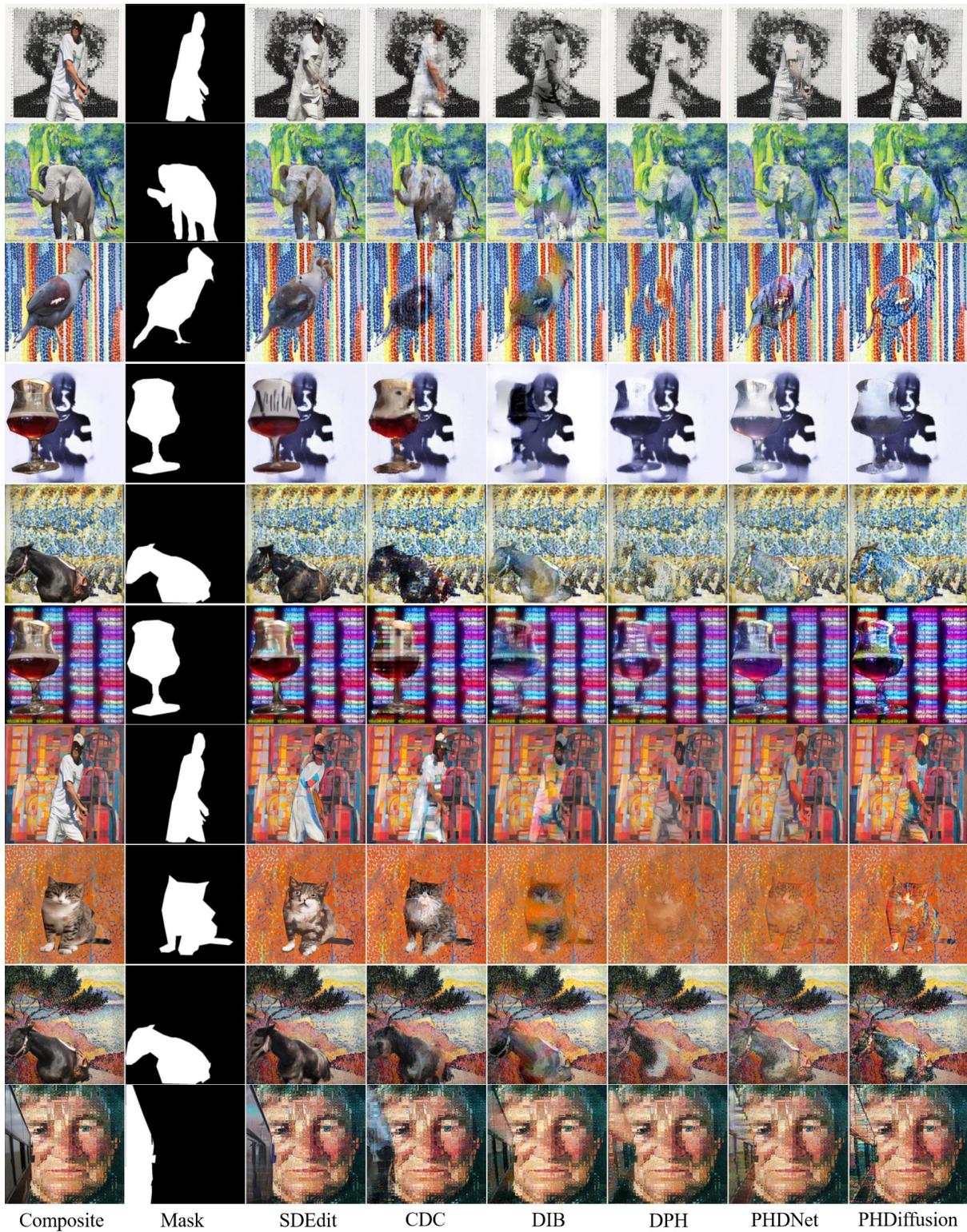


Figure 4: From left to right, we show the composite image, mask, harmonized results of SDEdit [8], CDC [3], DIB [11], DPH [7], PHDNet [1], and our PHDiffusion. Best viewed in color and zoom in.

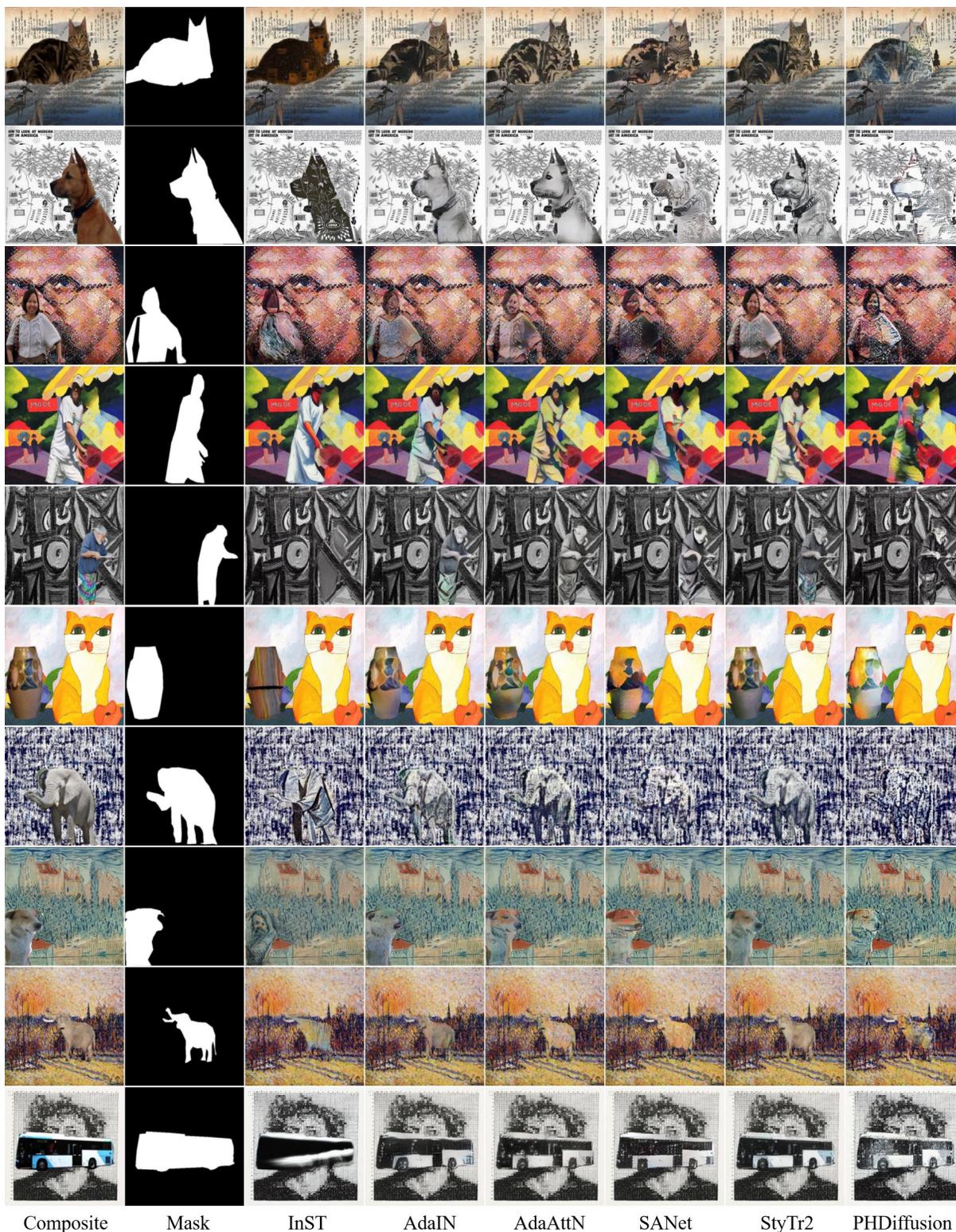


Figure 5: From left to right, we show the composite image, mask, example results for InST [12], AdaIN [4], AdaAttN [6], SANet [10], StyTr2 [2], and our PHDiffusion. Best viewed in color and zoom in.