# Customizing Text-to-Image Models with a Single Image Pair

Maxwell Jones[1]    Sheng-Yu Wang[1]    Nupur Kumari[1]
David Bau[2]    Jun-Yan Zhu[1]

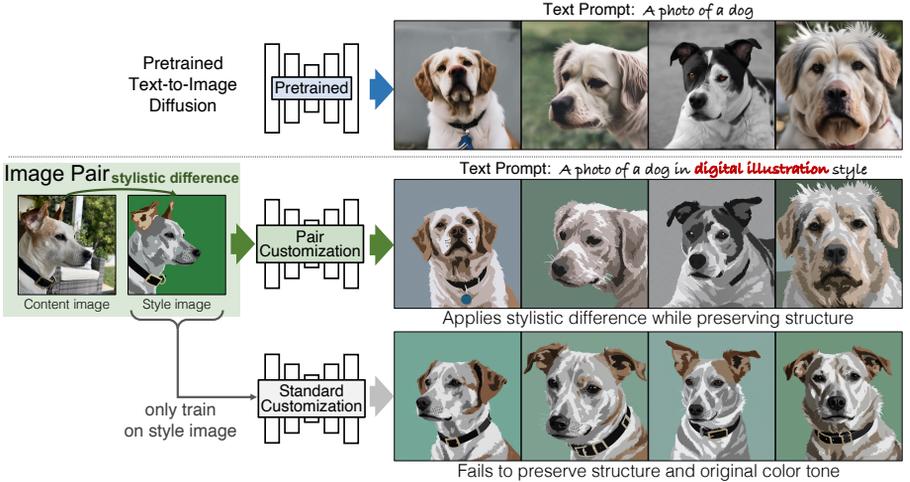Carnegie Mellon University[1]    Northeastern University[2]

**Abstract.** Art reinterpretation is the practice of creating a variation of a reference work, making a paired artwork that exhibits a distinct artistic style. We ask if such an image pair can be used to customize a generative model to capture the demonstrated stylistic difference. We propose `Pair Customization`, a new customization method that learns stylistic difference from a *single* image pair and then applies the acquired style to the generation process. Unlike existing methods that learn to mimic a single concept from a collection of images, our method captures the stylistic difference between paired images. This allows us to apply a stylistic change without overfitting to the specific image content in the examples. To address this new task, we employ a joint optimization method that explicitly separates the style and content into distinct LoRA weight spaces. We optimize these style and content weights to reproduce the style and content images while encouraging their orthogonality. During inference, we modify the diffusion process via a new style guidance based on our learned weights. Both qualitative and quantitative experiments show that our method can effectively learn style while avoiding overfitting to image content, highlighting the potential of modeling such stylistic differences from a single image pair.

## 1  Introduction

Artistic works are often inspired by a reference image, a recurring scene, or even a previous piece of art [54]. Such creations involve re-interpreting an original composition in the artist's unique style. A notable example is Van Gogh's *Repetitions* [63], in which the artist created multiple versions of the same scenes with his distinctive expressiveness, including adaptations of other artists' work. Such sets of variations allow close comparison of stylized art to a reference image, providing unique insights into an artist's detailed techniques and choices.

In our work, we explore how such *content-style image pairs* can be used to customize a generative model to capture the demonstrated stylistic difference. Our goal is to customize a pre-trained generative model to synthesize stylized images, distilling the essence of the style from as few as a single pair without fixating on specific content. We wish to create a model capable of re-interpreting a variety of different content in the style demonstrated by the paired variation.

Prior works on model customization/personalization [18,43,71] take one or a few images of a single concept to customize large-scale text-to-image models [67,

**Fig. 1:** Given a *single* image pair, we present `Pair Customization`, a method for customizing a pre-trained text-to-image model and learning a new style from the image pair's stylistic difference. Our method can apply the learned stylistic difference to new input images while preserving the input structure. Compared to Dreambooth LoRA [33, 73], a standard customization method that solely use style images, our method effectively disentangles style and content, resulting in better structure, color preservation, and style application. Style image credit: Jack Parkhouse.

69]. While they aim to learn styles without using pairs, the generated samples from these customized models often resemble the training images' content, such as specific objects, persons, and scene layouts. In Figure 1, we observe that standard single-image customization (3rd row) alters the subject, color tone, and pose of the original image (1st row). These issues arise because the artistic intent is difficult to discern from a single image: unlike image pairs that can demonstrate a style through contrasts, a singleton example will always intertwine choices of both style and content. Due to this ambiguity, the model fails to capture the artistic style accurately and, in some cases, overfits and generates the subject-specific details rather than the style, as shown in Figure 6.

On the other hand, our `Pair Customization` method exploits the contrast between image pairs to generate pairwise consistent images while better disentangling style and content. In Figure 1 (2nd row), our method accurately follows the given style, turning the background into a single color matching the original background and preserving the identity and pose for each dog. Our method achieves this by disentangling the intended style from the image pair.

Our new customization task is challenging since text-to-image models were not initially designed to generate *pairwise* content. Even when given specific text prompts like "`a portrait`" and "`a portrait with Picasso style`", a text-to-image diffusion model often struggles to generate images with consistent structure from the same noise seed. Therefore, it remains unclear how a customized model can generate stylized images while maintaining the original structure.

To address the challenges, we first propose a joint optimization method with separate sets of low-rank adaptation [33] (LoRA) weights for style and content. The optimization encourages the content LoRA to reconstruct the content image and the style LoRA to apply the style to the content. We find that the resulting style LoRA can apply the same style to other unseen content. Furthermore, we enforce row-space orthogonality [64] between style and content LoRA parameters to improve style and content disentanglement. Next, we extend the standard classifier-free guidance method [31] and propose style guidance. Style guidance integrates style LoRA predictions into the original denoising path, which aids in better content preservation and facilitates smoother control over the stylization strength. This method is more effective than the previous technique, where a customized model's strength is controlled by the magnitude of LoRA weights [74].

Our method is built upon Stable Diffusion XL [65]. We experiment with various image pairs, including different categories of content (e.g., portraits, animals, landscapes) and style (e.g., paintings, digital illustrations, filters). We evaluate our method on the above single image pairs and demonstrate the advantage of our method in preserving diverse structures while applying the stylization faithfully, compared to existing customization methods. Our code, models, and data are available on our webpage.

## 2   Related Works

**Text-to-image generative models.**    Deep generative models aim to model the data distribution of a given training set [16, 25, 30, 42, 59, 84]. Recently, large-scale text-to-image models [4, 11, 24, 38, 52, 62, 65, 67, 69, 76, 77, 96] trained on internet-scale training data [9, 78] have shown exceptional generalization. Notably, diffusion models [30, 83] stand out as the most widely adopted model class. While existing models can generate a broad spectrum of objects and concepts, they often struggle with rare or unseen concepts. Our work focuses on teaching these models to understand and depict a new style concept. Conditional generative models [8, 37, 46, 56, 61, 75, 97] learn to transform images across different domains, but the training often requires thousands to millions of image pairs. We focus on a more challenging case, where only a single image pair is available.

**Customizing generative models.**    Model customization, or personalization, aims to adapt an existing generative model with additional data, with the goal of generating outputs tailored to specific user preferences. Earlier efforts mainly focus on customizing pre-trained GANs [25, 40, 41] for smaller datasets [39, 57, 98], incorporating user edits [6, 90, 91], or aligning with text prompts [21, 58]. Recently, the focus has pivoted towards adapting large-scale text-to-image models to generate user-provided concepts, typically presented as one or a few images. Simply fine-tuning on the concept leads to overfitting. To mitigate this and enable variations via free text, several works explored different regularizations, including prior preservation [43, 71], human alignment [82], as well as parameter update restriction, where we only update text tokens [1, 15, 18, 89], attention layers [19, 27, 43], low-rank weights [33, 73, 86], or clusters of neurons [50].

More recent methods focus on encoder-based approaches for faster personalization [2, 13, 14, 20, 44, 53, 72, 80, 88, 94, 95]. Several works further focus on multiple concepts [3, 26, 43, 64, 79] instead of only a single concept. Our method takes inspiration from these techniques; however, we aim to address an inherently different task. Instead of learning concepts from an image collection, we customize the model to *learn stylistic differences* from an image pair.

**Style and content separation.**    Various past works have explored learning a style while separating it from content [12, 23, 34, 45, 85]. Our work is inspired by the seminal work Image Analogy [29], a computational paradigm that takes an image pair and applies the same translation to unseen images. Common image analogy methods include patch-wise similarity matching [29, 36, 47] and data-driven approaches [5, 60, 68, 87, 92, 99]. Different from these, we aim to exploit the text-guided generation capabilities of large-scale models so that we can *directly* use the style concept with unseen context. Recently, StyleDrop [82] has been proposed to learn a custom style for masked generative transformer models. Concurrent with our work, Hertz et al. [28] introduced a method for generating images with style consistency, offering the option of using a style reference image. In contrast, we exploit an image pair to better discern the stylistic difference. Unlike StyleDrop, we do not rely on human feedback in the process.
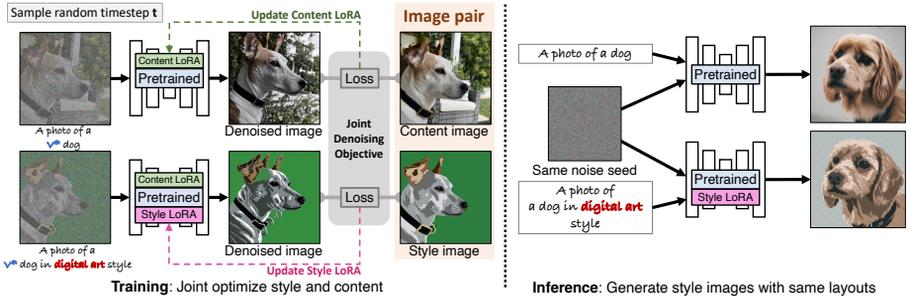
## 3    Method

Our method seeks to learn a new style from a single image pair. This task is challenging, as models tend to overfit when trained on a single image, especially when generating images in the same category as the training image (e.g., a model trained and tested on dog photos). To reduce this overfitting, we introduce a new algorithm aimed at disentangling the structure of the subject from the style of the artwork. Specifically, we leverage the image pair to learn separate model weights for style and content. At inference time, we modify the standard classifier-free guidance formulation to help preserve the original image structure when applying the learned style. In this section, we give a brief overview of diffusion models, outline our design choices, and explain the final method in detail.

### 3.1    Preliminary: Model Customization

**Diffusion models.**    Diffusion models [30, 81, 84], map Gaussian noise to the image distribution through iterative denoising. Denoising is learned by reversing the forward diffusion process $\mathbf{x}_0, \ldots, \mathbf{x}_T$, where image $\mathbf{x}_0$ is slowly *diffused* to random noise $\mathbf{x}_T$ over $T$ timesteps, defined by $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$ for timestep $t \in [0, T]$. Noise $\epsilon \sim \mathcal{N}(0, I)$ is randomly sampled, and $\bar{\alpha}_t$ controls the noise strength. The training objective of diffusion models is to denoise any intermediate noisy image $\mathbf{x}_t$ via noise prediction:

$$\mathbb{E}_{\epsilon, \mathbf{x}, \mathbf{c}, t}\left[w_t\|\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}, t)\|^2\right], \tag{1}$$

Fig. 2: **Method overview**. (Left) We disentangle style and content from an image pair by jointly training two low-rank adapters, StyleLoRA and ContentLoRA, representing style and content, respectively. Our training objective consists of two losses: The first loss fine-tunes ContentLoRA to reconstruct content image conditioned on a content prompt. The second loss encourages reconstructing the style image using *both* StyleLoRA and ContentLoRA conditioned on a style prompt, but we only optimize Style LoRA for this loss. (Right) At inference time, we only apply StyleLoRA to customize the model. Given the same noise seed, the customized model generates a stylized counterpart of the original pre-trained model output. V∗ is a fixed random rare token that is a prompt modifier for the content image. Style image credits: Jack Parkhouse

where $w_t$ is a time-dependent weight, $\epsilon_\theta(\cdot)$ is the denoiser that learns to predict noise, and **c** denotes extra conditioning input, such as text. At inference, the denoiser $\epsilon_\theta$ will gradually denoise random Gaussian noise into images. The resulting distribution of generated images approximates the training data distribution [30].

In our work, we use Stable Diffusion XL [65], a large-scale text-to-image diffusion model built on Latent Diffusion Models [69]. The model consists of a U-Net [70] trained on the latent space of an auto-encoder, with text conditioning from two text encoders, CLIP [66] and OpenCLIP [35].

**Model customization with low-rank adapters.** Low-Rank Adapters (LoRA) [33] is a parameter-efficient fine-tuning method [32] that applies low-rank weight changes $\Delta\theta_{\text{LoRA}}$ to pre-trained model weights $\theta_0$. For each layer with an initial weight $W_0 \in \mathbb{R}^{m \times n}$, the weight update is defined by $\Delta W_{\text{LoRA}} = BA$, a product of learnable matrices $B \in \mathbb{R}^{m \times r}$ and $A \in \mathbb{R}^{r \times n}$, where $r \ll \min(m, n)$ to enforce the low-rank constraint. The weight matrix of a particular layer with LoRA is:

$$W_{\text{LoRA}} = W_0 + \Delta W_{\text{LoRA}} = W_0 + BA. \tag{2}$$

At inference time, the LoRA strength is usually controlled by a scaling factor $\alpha \in [0, 1]$ applied to the weight update $\Delta W_{\text{LoRA}}$ [74]:

$$W_{\text{LoRA}} = W_0 + \alpha \Delta W_{\text{LoRA}}. \tag{3}$$

LoRA has been applied for customizing text-to-image diffusion models to learn new concepts with as few as three to five images [74].

## 3.2   Style Extraction from an image pair

We aim to customize a pre-trained model with an artistic style in order to stylize the original model outputs while preserving their content, as shown in Figure 2 (right). To achieve this, we introduce style LoRA weight $\theta_{\text{style}} = \theta_0 + \Delta\theta_{\text{style}}$. While a pre-trained model generates content from a noise seed and text $c$, style LoRA's goal is to generate a stylized counterpart of original content from the same noise seed and a style-specific text prompt $\mathbf{c}_{\text{style}}$, where $\mathbf{c}_{\text{style}}$ is original text $c$ appended by suffix `"in <desc> style"`. Here, `<desc>` is a placeholder for some worded description of the style (e.g., "digital art"), and style LoRA $\theta_{\text{style}}$ associates `<desc>` to the desired style.

Unfortunately, learning style LoRA $\theta_{\text{style}}$ from a single style image often leads to copying content (Figure 6). Hence, we explicitly learn disentanglement from a style and content image, denoted by $\mathbf{x}_{\text{style}}$ and $\mathbf{x}_{\text{content}}$, respectively.

**Disentangling style and content.**   We leverage the fact that the style image share the same layout and structure as the content image. Our key idea is to learn a separate content LoRA $\theta_{\text{content}} = \theta_0 + \Delta\theta_{\text{content}}$ to reconstruct the content image. By explicitly modeling the content, we can train the style LoRA to "extract" the stylistic differences between the style and content image. We apply both style and content LoRA to reconstruct the style image, i.e., $\theta_{\text{combined}} = \theta_0 + \Delta\theta_{\text{content}} + \Delta\theta_{\text{style}}$. This approach prevents leaking the content image to style LoRA, resulting in a better stylization model.

During training, we feed the content LoRA $\theta_{\text{content}}$ with a content-specific text $\mathbf{c}_{\text{content}}$, which contains a random rare token `V*`, and feed the combined model $\theta_{\text{combined}}$ with $\mathbf{c}_{\text{style}}$, where $\mathbf{c}_{\text{style}}$ is `"{`$\mathbf{c}_{\text{content}}$`} in <desc> style"`. Figure 2 (Left) summarizes our training process.

**Jointly learning style and content.**   We employ two different objectives during every training step. To learn the content of the image, we first employ the standard training objective for diffusion models as described in Section 3.1 with the content image:

$$\mathcal{L}_{\text{content}} = \mathbb{E}_{\epsilon, \mathbf{x}_{\text{content}}, t} \left[ w_t \| \epsilon - \epsilon_{\theta_{\text{content}}} (\mathbf{x}_{t,\text{content}}, \mathbf{c}_{\text{content}}, t) \|^2 \right], \qquad (4)$$

where $\epsilon_{\theta_{\text{content}}}$ is the denoiser with content LoRA applied, $\mathbf{x}_{t,\text{content}}$ is a noisy content image at timestep $t$, and $\mathbf{c}_{\text{content}}$ is text representing the content image, including some rare token `V*`. Next, we optimize the combined style and content weights to reconstruct the style image. In particular, we only train the style LoRA weights during this step, while stopping the gradient flow to the content LoRA weights via stopgrad sg[·]:

$$\theta_{\text{combined}} = \theta_0 + \text{sg}[\Delta\theta_{\text{content}}] + \Delta\theta_{\text{style}}. \qquad (5)$$

We then apply diffusion objective to train $\theta_{\text{combined}}$ to denoise $\mathbf{x}_{t,\text{style}}$, a noisy style image at timestep $t$:

$$\mathcal{L}_{\text{combined}} = \mathbb{E}_{\epsilon, \mathbf{x}_{\text{style}}, t} \left[ w_t \| \epsilon - \epsilon_{\theta_{\text{combined}}} (x_{t,\text{style}}, \mathbf{c}_{\text{style}}, t) \|^2 \right], \qquad (6)$$

**Fig. 3: (Left) Orthogonal adaptation.** Enforcing row-space orthogonality between style and content LoRA improves image quality, where the images capture the style better and have fewer visual artifacts. **(Right) Style guidance.** We compare style control capabilities between our style guidance and standard LoRA weight scaling [74]. **Blue** and **green** stand for the LoRA weight scale and style guidance scale, respectively. Style guidance better preserves content when the style is applied. More details of style guidance formulation are in Section 3.3.

where $\epsilon_{\theta_{\text{combined}}}$ is the denoiser with both LoRAs applied as in Equation 5, $\mathbf{c}_{\text{style}}$ is "{$\mathbf{c}_{\text{content}}$} in <desc> style", and <desc> is a worded description of the style (e.g., "digital art"). Finally, we jointly optimize the LoRAs with the two losses:

$$\min_{\Delta\theta_{\text{content}}, \Delta\theta_{\text{style}}} \mathcal{L}_{\text{content}} + \mathcal{L}_{\text{combined}} \qquad (7)$$

Figure 2 provides an overview of our method. Next, we discuss the regularization that promotes the disentanglement of style from content.

**Orthogonality between style and content LoRA.** To further encourage style and content LoRAs to represent separate concepts, we enforce orthogonality upon the LoRA weights. We denote by $W_0$ the original weight matrix and $W_{\text{content}}$, $W_{\text{style}}$ the LoRA modifications (layer index omitted for simplicity). Reiterating Equation 2, we decompose $W_{\text{content}}$, $W_{\text{style}}$ into low-rank matrices:

$$W_{\text{content}} = W_0 + B_{\text{content}} A_{\text{content}}; \; W_{\text{style}} = W_0 + B_{\text{style}} A_{\text{style}}. \qquad (8)$$

We initialize $B_{\text{content}}, B_{\text{style}}$ with the zero matrix and choose the rows of $A_{\text{content}}, A_{\text{style}}$ from an orthonormal basis. We then fix $A_{\text{content}}, A_{\text{style}}$ and only update $B_{\text{content}}, B_{\text{style}}$ in training. This forces the style and content LoRA updates to respond to orthogonal inputs, and empirically reduces visual artifacts, as shown in Figure 3. This technique is inspired by Po et al. [64]. While their work focuses on merging multiple customized objects after each is trained separately, we apply the method for style-content separation during joint training.

### 3.3 Style Guidance

A common technique to improve text-to-image model's sample quality is via classifier-free guidance [31]:

$$\hat{\epsilon}_\theta(\mathbf{x}_t, \mathbf{c}) = \epsilon_\theta(\mathbf{x}_t, \varnothing) + \lambda_{\text{cfg}}(\epsilon_\theta(\mathbf{x}_t, \mathbf{c}) - \epsilon_\theta(\mathbf{x}_t, \varnothing)), \qquad (9)$$

where $\hat{\epsilon}_\theta(\mathbf{x}_t, \mathbf{c}, t)$ is the new noise prediction, $\varnothing$ denotes no conditioning, and $\lambda_{\text{cfg}}$ controls the amplification of text guidance. For notation simplicity, we omit the timestep $t$ in this equation and subsequent ones.

To improve pairwise consistency between original and stylized content, we propose an inference algorithm that preserves the original denoising path while adding controllable style guidance:

$$
\begin{aligned}
\hat{\epsilon}_{\theta_0, \theta_{\text{style}}}(\mathbf{x}_t, \mathbf{c}, \mathbf{c}_{\text{style}}) = \; & \epsilon_{\theta_0}(\mathbf{x}_t, \varnothing) \\
& + \lambda_{\text{cfg}}(\epsilon_{\theta_0}(\mathbf{x}_t, \mathbf{c}) - \epsilon_{\theta_0}(\mathbf{x}_t, \varnothing)) \\
& + \lambda_{\text{style}}(\epsilon_{\theta_{\text{style}}}(\mathbf{x}_t, \mathbf{c}_{\text{style}}) - \epsilon_{\theta_0}(\mathbf{x}_t, \mathbf{c})),
\end{aligned}
\tag{10}
$$

where style guidance is the difference in noise prediction between style LoRA and the pre-trained model. Style guidance strength is controlled by $\lambda_{\text{style}}$, and setting $\lambda_{\text{style}} = 0$ is equivalent to generating original content. In Figure 3, we compare our style guidance against scaling LoRA weights (Equation 3), and we find that our method better preserves the layout. More details and a derivation of our style guidance are in Appendix B.
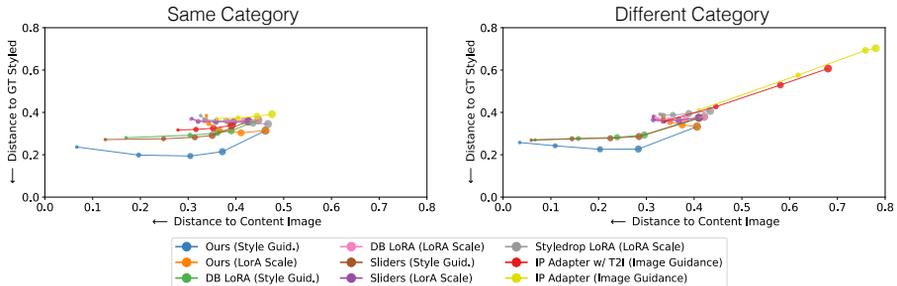
Previous works [8,49] have also considered applying multiple guidance terms with diffusion models. A major difference in our approach is that we obtain additional guidance from a customized model and apply it to the original model. Styledrop [82] considers a similar formulation with two guidance terms but for transformer-based generative models.

**Blending multiple learned styles.**    With a collection of models customized by our method, we can blend the learned styles as follows. Specifically, given some set of styles $\mathcal{S}$ and strengths $\lambda_{\text{style}_0}, \ldots, \lambda_{\text{style}_n}$, we can blend the style guidance from each model, and our new inference path is represented by

$$
\begin{aligned}
\hat{\epsilon}_{\theta_0, \theta_{\text{style}}}(\mathbf{x}_t, \mathbf{c}, \mathbf{c}_{\text{style}}) = \; & \epsilon_{\theta_0}(\mathbf{x}_t, \varnothing) \\
& + \lambda_{\text{cfg}}(\epsilon_{\theta_0}(\mathbf{x}_t, \mathbf{c}) - \epsilon_{\theta_0}(\mathbf{x}_t, \varnothing)) \\
& + \sum_{\text{style}_i \in \mathcal{S}} \lambda_{\text{style}_i}(\epsilon_{\theta_{\text{style}_i}}(\mathbf{x}_t, \mathbf{c}_{\text{style}_i}) - \epsilon_{\theta_0}(\mathbf{x}_t, \mathbf{c})),
\end{aligned}
\tag{11}
$$

We can vary the strengths of any parameter $\lambda_{\text{style}_i}$ to seamlessly increase or decrease style application while preserving content. Figure 7 gives a qualitative example of blending two different styles while preserving image content.

**Implementation details.**    We train all models using an AdamW optimizer [51] and learning rate $1 \times 10^{-5}$. For baselines, we train for 500 steps. For our method, we first train our content weights on the content image for 250 steps, and then train jointly for 500 additional steps. All image generation is performed using 50 steps of a PNDMScheduler [48]. For all methods using inference with LoRA adapters, we use SDEdit [55] to further preserve structure. Specifically, normal classifier-free guidance on the original prompt without style is used for the first 10 steps. We then apply style guidance/LoRA scale for the rest of the timesteps.

**Fig. 4: Quantitative comparison with baselines on learned style.** Given a fixed inference path, our method's pareto dominates baselines for image generation both on the same category as training (left) and when evaluated on categories different from training, e.g., trained on human portraits but tested on dog images (right). Secondly, our proposed style guidance outperforms standard LoRA weight scale guidance for our training method (blue vs. orange), DB LoRA (green vs. pink), and Sliders (brown vs. purple). In the Appendix's Figure 9, we further evaluate the diversity of generated images. We show that baselines often lose diversity, while our method leads to diverse generations while still achieving lower perceptual distance to the ground truth style. Increased marker size corresponds to an increase in style guidance scale.

## 4   Experiments

### 4.1   Dataset

In this section, we show our method's results on various image pairs and compare them with several baselines. We explain our dataset, baselines, and metrics in detail, then we present quantitative and qualitative results.

**Datasets.**   To enable large-scale quantitative evaluation, we construct a diverse set of paired style and content images as follows. First, we generate 40 content images for each class: headshots, animals, and landscapes. When generating images in the headshot class, we generate 20 images with the prompt "`A professional headshot of a man`" and 20 images with the prompt "`A professional headshot of a woman`". Similarly, we split the animal class into photos of dogs and cats. To curate synthetic pairs, we then apply image editing or image-to-image translation methods to all the content images to obtain the stylized version. For each unique prompt, we choose a *single paired instance* as training data and hold out the other pairs with the same prompt as a test set (Same Category). For each prompt, we also choose 5 pairs from each of the other prompts as a secondary test set (Different Category). We show all our synthetic training image pairs in Appendix C. By leveraging synthetic pairs for evaluation, we can train on a single synthetic pair and test our results against held out synthetic style images. Secondly, we qualitatively compare against single artist pairs in Figure 6. Next, we describe the specific methods to create the paired dataset.

**LEDITS++**   [7] is a diffusion-based image editing technique that transforms an image by updating the inference path of a diffusion model. After fine-grained

inversion, a global prompt and a set of translation prompts representing a new style or object are used to perform the image translation. We leverage LED-ITS++ on all images with the translation prompt "Impressionist style". Further, we change the word "photo" to "painting" in the original prompt when generating the style image.

**White-box cartoonization.**    Cartoonization [93] is a GAN-based image-to-image translation technique that applies a cartoon-like effect to real images. We apply the cartoonized model to our set of generated images to create image pairs.

**Stylized neural painting.**   Stylized Neural Painting [100] is a rendering based image to image translation technique where an image is reconstructed via $N$ painting strokes, where the strokes are guided by a loss function that encourages the final translated image to resemble the original. We use the Neural Painting model with $N = 1000$ to create image pairs.

**Posterization.**    Posterization is an image filtering technique that reduces the number of distinct colors in a given image to some fixed number $N$, reducing color variation and creating fixed color areas. We apply posterization to images in our training set with $N = 8$.

### 4.2   Baselines and Evaluation Metrics

**Baselines.**    We compare our method against – (1) DreamBooth LoRA [33, 74] (DB LoRA), (2) Concept Sliders [22] (3) IP-adapters [95], (4) IP-adapters w/ T2I, and (5) StyleDrop [82]. DB LoRA uses only the style image and fine-tunes low-rank adapters in all the linear layers in the attention blocks of the diffusion model. We evaluate different amounts of style applications for DB LoRA using the standard LoRA scale [74] and our style guidance. Concept sliders presents a paired image model customization method that trains a single low-rank adapter jointly on both images, with different reconstruction losses for the style and content images. We also evaluate using both the standard LoRA scale and our style guidance. IP-adapters is an encoder-based method that does not require training for every style and takes a style image as an extra condition separate from the text prompt. Increasing or decreasing the guidance from the input style image is possible by scaling the weight of the image conditioning. We consider the SDXL [65] implementation of this method. For the IP-Adapter, we also compare against the stronger baseline of providing extra conditioning of an edge map of the content image through T2I Adapters [56] to preserve the content image structure. The recently proposed Styledrop [82] technique for learning new styles is based on MUSE [11], and uses human feedback in its method. Since MUSE is not publicly available, we follow Style-Aligned Image Generation's [28] setup, and implement a version of StyleDrop on SDXL. Specifically, we train low-rank linear layers following each Feed-Forward layer in the attention blocks of SDXL. For a fair comparison, we train Styledrop without human feedback.

**Evaluation metrics.**    When evaluating the performance of each method, we consider two quantitative metrics: perceptual distance to ground truth style images and structure preservation from the original image. A better customization
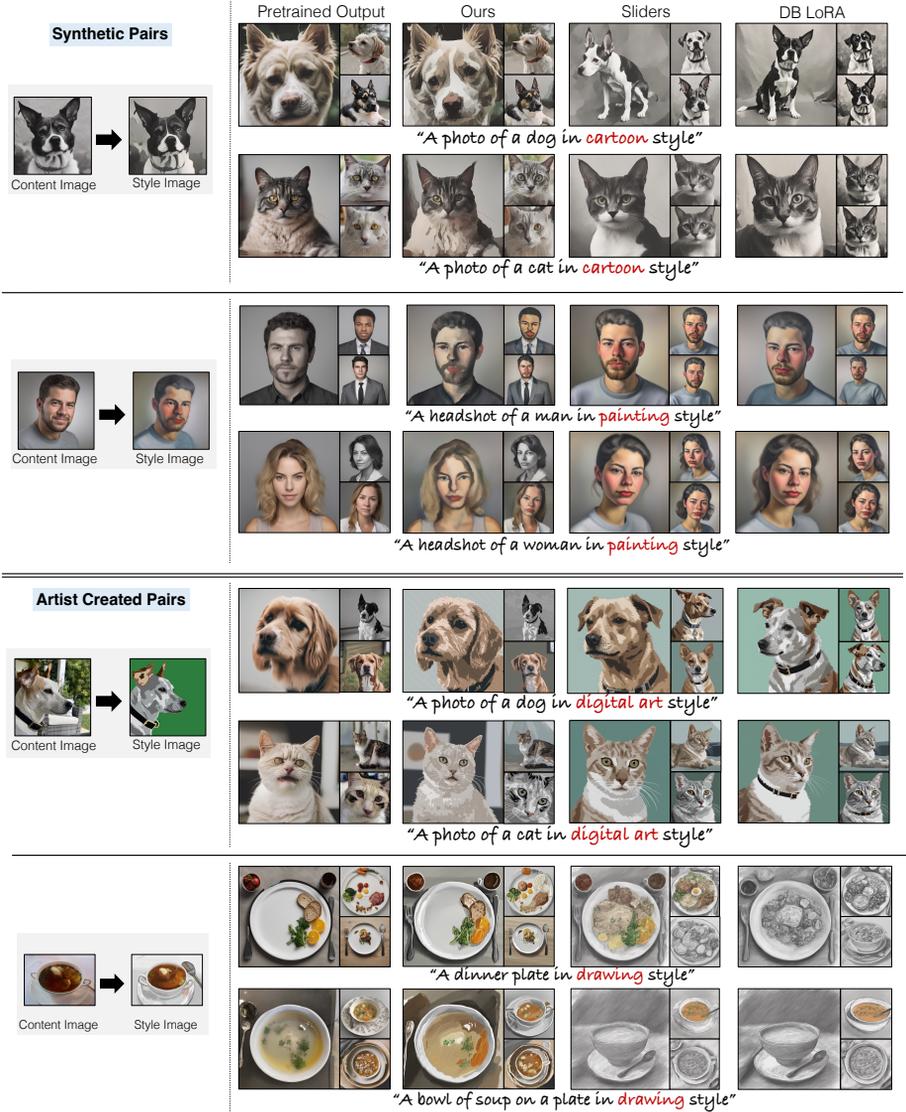
**Fig. 5: Human preference study.** Our method is preferred over the baselines ($\geqslant$ 60%). Further, our full method, including orthogonal weight matrices (Section 3.2), is preferred over the one w/o orthogonal weight matrices, specifically for the same category as training pair, e.g., trained on a headshot of a man and tested on other headshots of man. The Gray dashed line denotes 50% chance performance.

method will have a low perceptual distance to the ground truth style images while still preserving content of the original image before adding style. We measure these using – (1) *Distance to GT Styled*: given holdout ground truth style images, we measure the perceptual distance between our styled outputs and the ground truth style images using DreamSim [17], a recent method for measuring the perceptual distance between images. DreamSim image embeddings are comprised of an ensemble of image embedding models, including CLIP [66] and DINO [10], which are then fine-tuned so the final embeddings respect human perception. We measure DreamSim distance as (1 - cosine similarity) between DreamSim embeddings, where a lower value implies that the images are perceptually more similar. (2) *Distance to Content Image*: to measure content preservation after style application, we measure the perceptual distance of our generated style image to the original content image with no style guidance. We again use DreamSim, this time comparing styled and content images. Note here that a perceptual distance of zero to the content image is undesirable, as this would require no style to be applied. However, a better-performing method should obtain a better tradeoff between the two distances. (3) We also perform a *human preference study* of our method against baselines.

## 4.3 Results

**Quantitative evaluation.** We show quantitative results in Figure 4. Increased marker size (circles) indicates the higher application of style, and line color determines the method. When evaluating style similarity vs. structure preservation in Figure 4, we see that our training method's Pareto dominates all baselines, yielding lower perceptual distance to style images while still being perceptually similar to the original content image. Secondly, style guidance outperforms the LoRA scale for Ours (blue vs orange), DB LoRA (green vs. pink), and Concept Sliders (brown vs. purple), highlighting the effectiveness of both parts of our method.

**Fig. 6:** Result of our method compared to the strongest baselines. When only training with the style image as in DB LoRA, the image structure is not preserved and overfitting occurs. While Concept Slider's training scheme [22] uses both style and content images, it still exhibits overfitting and loss of structure in many cases. Our method preserves the structure of the input mage while faithfully applying the desired style. We use style guidance strength 3 and classifier guidance strength 5. Style image credits: Jack Parkhouse (Third row) and Aaron Hertzmann (Fourth row)
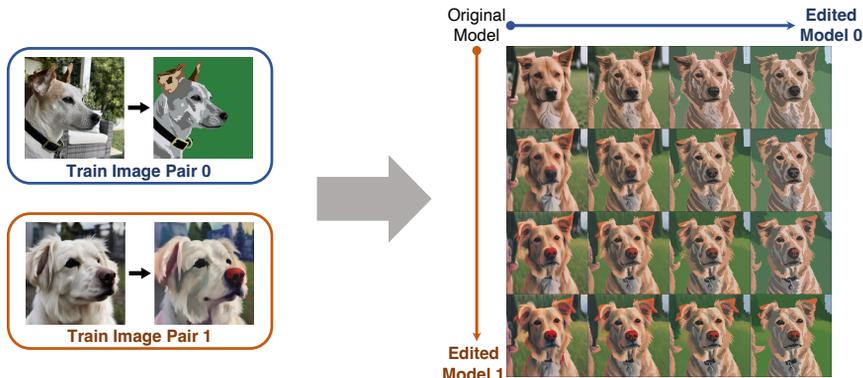
**Qualitative evaluation.**    We compare our method with the highest perform-
ing baselines in Figure 6. The finetuning-based methods DB LoRA [33, 74] and
Concept Sliders  [22] outperform the encoder-based method [95] for our task.
Hence, we compare against that in Figure 6. For both baselines, we modulate
style application with LoRA scale (Equation 3). We observe that DB LoRA often
fails to generate the style-transformed version of the original image and overfits
to the training pair image when generating similar concepts. There are two main
reasons why this may occur. First, we are in a challenging case where there is
only 1 training image instead of the usual $3-5$ images that customization meth-
ods use. Second, we are prompting the model on the same or very similar text
prompts to the training prompt, and the baseline method overfits to the train-
ing image for these prompts. Our method preserves the structure of the original
image while applying the learned style. Moreover, applying our style guidance
instead of the LoRA scale benefits the baseline method as well (Figure 6, last 2
columns), as it can better preserve the structure of the original image, though
it still tends to overfit to the content of the training image. We observe a simi-
lar issue for other baselines as well. We show a qualitative comparison with the
other baselines in Appendix A. We also compare with baselines using our style
guidance for style application at inference time in Appendix A.

**User preference study.**    We perform a user preference study using Amazon
Mechanical Turk. We test our method against all baselines, as well as a version of
our method trained without orthogonality constraint. Specifically, we test on all
datasets in Section 4.1. When evaluating against DB LoRA and Concept Sliders,
we consider inference with both LoRA scale as in Equation 3 and style guidance
as in Equation 10. For each method, we pick a single style strength that performs
most optimally according to quantitative metrics as in Figure 4. Full details are
available in Appendix C. We collect 400 responses per paired test of ours vs the
other method. The user is shown an image generated via our method and an
image generated via the other method and asked to select the image that best
applies the given style to the new content image. We provide a detailed setup of
the user study in Appendix C. As shown in Figure 5, our method is favored by
users in comparison to baselines, whether evaluating images generated within the
same category as the training image pair or across different categories. Secondly,
users prefer our full method to ours without the orthogonality constraint.

**Blending learned styles.**    We show that we can blend the learned styles by
applying a new inference path, defined in Equation 11. In Figure 7, we show
the results of blending two models. We can seamlessly blend the two styles at
varying strengths while still preserving the content.

## 5    Discussion and Limitations

In this work, we have introduced a new task: customizing a text-to-image model
with a single image pair. To address this task, we have developed a customiza-
tion method that explicitly disentangles style and content through both training
objectives and a separated parameter space. Our method enables us to grasp

**Fig. 7: Blending multiple style guidances.** We can compose multiple customized models by directly blending each style guidance together. Adjusting blending strength of each model allows us to acquire a smooth style transition. Train Image Pair 0 Style image credits: Jack Parkhouse



**Fig. 8: Limitations.** *Left*: our method struggles with categories when they significantly differ from the training categories. Here, our method fails to transfer the artistic style of landscape image pairs to human portraits. *Right*: our method can cause structure changes in some instances, like change of body position or background changes.

the style concept without memorizing the content of input examples. While our approach outperforms existing customization methods, it still exhibits several limitations, as discussed below.

**Limitations.**   First, while our method is able to transfer the style from a pair of dog images to cat photos, it struggles to handle completely different categories from the training image pair, particularly when the test category significantly differs from the training. As shown in Figure 8 (left), our method falls short of faithfully replicating the style of landscape paintings in some human images.

Second, our current method relies on test-time optimization, which takes around 15 minutes on a single A5000 GPU. This can be computationally demanding if we need to process many image styles. Leveraging encoder-based approaches [2, 72] for predicting style and content weights in a feed-forward manner could potentially speed up the customization process.

Finally, our method may occasionally fail to completely maintain input structure, as demonstrated in Figure 8 (right).

# References

1. Alaluf, Y., Richardson, E., Metzer, G., Cohen-Or, D.: A neural space-time representation for text-to-image personalization. In: SIGGRAPH Asia (2023) 3
2. Arar, M., Gal, R., Atzmon, Y., Chechik, G., Cohen-Or, D., Shamir, A., H. Bermano, A.: Domain-agnostic tuning-encoder for fast personalization of text-to-image models. In: SIGGRAPH Asia 2023 Conference Papers. pp. 1–10 (2023) 4, 14
3. Avrahami, O., Aberman, K., Fried, O., Cohen-Or, D., Lischinski, D.: Break-a-scene: Extracting multiple concepts from a single image. In: SIGGRAPH Asia 2023 Conference Papers. pp. 1–12 (2023) 4
4. Balaji, Y., Nah, S., Huang, X., Vahdat, A., Song, J., Kreis, K., Aittala, M., Aila, T., Laine, S., Catanzaro, B., et al.: ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. arXiv preprint arXiv:2211.01324 (2022) 3
5. Bar, A., Gandelsman, Y., Darrell, T., Globerson, A., Efros, A.A.: Visual prompting via image inpainting. In: Adv. Neural Inform. Process. Syst. (2022) 4
6. Bau, D., Liu, S., Wang, T., Zhu, J.Y., Torralba, A.: Rewriting a deep generative model. In: Eur. Conf. Comput. Vis. (2020) 3
7. Brack, M., Friedrich, F., Kornmeier, K., Tsaban, L., Schramowski, P., Kersting, K., Passos, A.: Ledits++: Limitless image editing using text-to-image models. arXiv preprint arXiv:2311.16711 (2023) 9
8. Brooks, T., Holynski, A., Efros, A.A.: Instructpix2pix: Learning to follow image editing instructions. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18392–18402 (2023) 3, 8, 23
9. Byeon, M., Park, B., Kim, H., Lee, S., Baek, W., Kim, S.: Coyo-700m: Image-text pair dataset. https://github.com/kakaobrain/coyo-dataset (2022) 3
10. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: Int. Conf. Comput. Vis. (2021) 11
11. Chang, H., Zhang, H., Barber, J., Maschinot, A., Lezama, J., Jiang, L., Yang, M.H., Murphy, K., Freeman, W.T., Rubinstein, M., et al.: Muse: Text-to-image generation via masked generative transformers. In: Proceedings of the 40th International Conference on Machine Learning. pp. 4055–4075 (2023) 3, 10
12. Chen, T.Q., Schmidt, M.: Fast patch-based style transfer of arbitrary style. arXiv preprint arXiv:1612.04337 (2016) 4
13. Chen, W., Hu, H., Li, Y., Rui, N., Jia, X., Chang, M.W., Cohen, W.W.: Subject-driven text-to-image generation via apprenticeship learning. In: Adv. Neural Inform. Process. Syst. (2023) 4

14. Chen, X., Huang, L., Liu, Y., Shen, Y., Zhao, D., Zhao, H.: Anydoor: Zero-shot object-level image customization. arXiv preprint arXiv:2307.09481 (2023) 4

15. Daras, G., Dimakis, A.G.: Multiresolution textual inversion. arXiv preprint arXiv:2211.17115 (2022) 3

16. Dinh, L., Sohl-Dickstein, J., Bengio, S.: Density estimation using real nvp. In: Int. Conf. Learn. Represent. (2017) 3

17. Fu, S., Tamir, N., Sundaram, S., Chai, L., Zhang, R., Dekel, T., Isola, P.: Dreamsim: Learning new dimensions of human visual similarity using synthetic data. arXiv preprint arXiv:2306.09344 (2023) 11, 21

18. Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A.H., Chechik, G., Cohen-or, D.: An image is worth one word: Personalizing text-to-image generation using textual inversion. In: The Eleventh International Conference on Learning Representations (2022) 1, 3

19. Gal, R., Arar, M., Atzmon, Y., Bermano, A.H., Chechik, G., Cohen-Or, D.: Designing an encoder for fast personalization of text-to-image models. ACM Transactions on Graphics (TOG) (2023) 3

20. Gal, R., Arar, M., Atzmon, Y., Bermano, A.H., Chechik, G., Cohen-Or, D.: Encoder-based domain tuning for fast personalization of text-to-image models. ACM Transactions on Graphics (TOG) 42(4), 1–13 (2023) 4

21. Gal, R., Patashnik, O., Maron, H., Bermano, A.H., Chechik, G., Cohen-Or, D.: Stylegan-nada: Clip-guided domain adaptation of image generators. ACM Transactions on Graphics (TOG) 41(4), 1–13 (2022) 3

22. Gandikota, R., Materzynska, J., Zhou, T., Torralba, A., Bau, D.: Concept sliders: Lora adaptors for precise control in diffusion models. arXiv preprint arXiv:2311.12092 (2023) 10, 12, 13

23. Gatys, L.A., Ecker, A.S., Bethge, M.: A neural algorithm of artistic style. arXiv preprint arXiv:1508.06576 (2015) 4

24. Gokaslan, A., Cooper, A.F., Collins, J., Seguin, L., Jacobson, A., Patel, M., Frankle, J., Stephenson, C., Kuleshov, V.: Commoncanvas: An open diffusion model trained with creative-commons images. arXiv preprint arXiv:2310.16825 (2023) 3

25. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial networks. Communications of the ACM 63(11), 139–144 (2020) 3

26. Gu, Y., Wang, X., Wu, J.Z., Shi, Y., Chen, Y., Fan, Z., Xiao, W., Zhao, R., Chang, S., Wu, W., et al.: Mix-of-show: Decentralized low-rank adaptation for multi-concept customization of diffusion models. Advances in Neural Information Processing Systems 36 (2024) 4

27. Han, L., Li, Y., Zhang, H., Milanfar, P., Metaxas, D., Yang, F.: Svdiff: Compact parameter space for diffusion fine-tuning. In: Int. Conf. Comput. Vis. (2023) 3

28. Hertz, A., Voynov, A., Fruchter, S., Cohen-Or, D.: Style aligned image generation via shared attention. arXiv preprint arXiv:2312.02133 (2023) 4, 10, 21, 22, 23

29. Hertzmann, A., Jacobs, C.E., Oliver, N., Curless, B., Salesin, D.H.: Image analogies (2001) 4

30. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. In: Adv. Neural Inform. Process. Syst. (2020) 3, 4, 5

31. Ho, J., Salimans, T.: Classifier-free diffusion guidance. arXiv preprint arXiv:2207.12598 (2022) 3, 7, 26

32. Houlsby, N., Giurgiu, A., Jastrzebski, S., Morrone, B., De Laroussilhe, Q., Gesmundo, A., Attariyan, M., Gelly, S.: Parameter-efficient transfer learning for nlp. In: International Conference on Machine Learning. pp. 2790–2799. PMLR (2019) 5

33. Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: Lora: Low-rank adaptation of large language models. In: Int. Conf. Learn. Represent. (2021) 2, 3, 5, 10, 13

34. Huang, X., Belongie, S.: Arbitrary style transfer in real-time with adaptive instance normalization. In: Int. Conf. Comput. Vis. (2017) 4

35. Ilharco, G., Wortsman, M., Wightman, R., Gordon, C., Carlini, N., Taori, R., Dave, A., Shankar, V., Namkoong, H., Miller, J., Hajishirzi, H., Farhadi, A., Schmidt, L.: Openclip (Jul 2021). https://doi.org/10.5281/zenodo.5143773, https://doi.org/10.5281/zenodo.5143773, if you use this software, please cite it as below. 5

36. Irony, R., Cohen-Or, D., Lischinski, D.: Colorization by example. In: Eurographics Conference on Rendering Techniques (2005) 4

37. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: IEEE Conf. Comput. Vis. Pattern Recog. (2017) 3

38. Kang, M., Zhu, J.Y., Zhang, R., Park, J., Shechtman, E., Paris, S., Park, T.: Scaling up gans for text-to-image synthesis. In: IEEE Conf. Comput. Vis. Pattern Recog. (2023) 3

39. Karras, T., Aittala, M., Hellsten, J., Laine, S., Lehtinen, J., Aila, T.: Training generative adversarial networks with limited data. In: Adv. Neural Inform. Process. Syst. (2020) 3

40. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: IEEE Conf. Comput. Vis. Pattern Recog. (2019) 3

41. Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of stylegan. In: IEEE Conf. Comput. Vis. Pattern Recog. (2020) 3

42. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. In: Int. Conf. Learn. Represent. (2014) 3

43. Kumari, N., Zhang, B., Zhang, R., Shechtman, E., Zhu, J.Y.: Multi-concept customization of text-to-image diffusion. In: 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1931–1941. IEEE Computer Society (2023) 1, 3, 4

44. Li, D., Li, J., Hoi, S.: Blip-diffusion: Pre-trained subject representation for controllable text-to-image generation and editing. Advances in Neural Information Processing Systems **36** (2024) 4

45. Li, Y., Wang, N., Liu, J., Hou, X.: Demystifying neural style transfer. arXiv preprint arXiv:1701.01036 (2017) 4

46. Li, Y., Liu, H., Wu, Q., Mu, F., Yang, J., Gao, J., Li, C., Lee, Y.J.: Gligen: Openset grounded text-to-image generation. In: IEEE Conf. Comput. Vis. Pattern Recog. (2023) 3

47. Liao, J., Yao, Y., Yuan, L., Hua, G., Kang, S.B.: Visual attribute transfer through deep image analogy. ACM Trans. Graph. **36**(4) (jul 2017) 4

48. Liu, L., Ren, Y., Lin, Z., Zhao, Z.: Pseudo numerical methods for diffusion models on manifolds. In: Int. Conf. Learn. Represent. (2022) 8

49. Liu, N., Li, S., Du, Y., Torralba, A., Tenenbaum, J.B.: Compositional visual generation with composable diffusion models. In: European Conference on Computer Vision. pp. 423–439. Springer (2022) 8

50. Liu, Z., Feng, R., Zhu, K., Zhang, Y., Zheng, K., Liu, Y., Zhao, D., Zhou, J., Cao, Y.: Cones: Concept neurons in diffusion models for customized generation. In: Int. Conf. Mach. Learn. (2023) 3

51. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: International Conference on Learning Representations (2018) 8
52. Luo, S., Tan, Y., Huang, L., Li, J., Zhao, H.: Latent consistency models: Synthesizing high-resolution images with few-step inference. arXiv preprint arXiv:2310.04378 (2023) 3
53. Ma, J., Liang, J., Chen, C., Lu, H.: Subject-diffusion: Open domain personalized text-to-image generation without test-time fine-tuning. arXiv preprint arXiv:2307.11410 (2023) 4
54. Markus, C.: How six different artists have re-interpreted da vinci's 'mona lisa'. https://www.parkwestgallery.com/six-different-artists-da-vinci-mona-lisa/ (11 2019) 1
55. Meng, C., He, Y., Song, Y., Song, J., Wu, J., Zhu, J.Y., Ermon, S.: Sdedit: Guided image synthesis and editing with stochastic differential equations. In: Int. Conf. Learn. Represent. (2022) 8
56. Mou, C., Wang, X., Xie, L., Wu, Y., Zhang, J., Qi, Z., Shan, Y.: T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 38, pp. 4296–4304 (2024) 3, 10
57. Nitzan, Y., Aberman, K., He, Q., Liba, O., Yarom, M., Gandelsman, Y., Mosseri, I., Pritch, Y., Cohen-Or, D.: Mystyle: A personalized generative prior. In: SIGGRAPH ASIA (2022) 3
58. Nitzan, Y., Gharbi, M., Zhang, R., Park, T., Zhu, J.Y., Cohen-Or, D., Shechtman, E.: Domain expansion of image generators. In: IEEE Conf. Comput. Vis. Pattern Recog. (2023) 3
59. Van den Oord, A., Kalchbrenner, N., Espeholt, L., Vinyals, O., Graves, A., et al.: Conditional image generation with pixelcnn decoders. In: Adv. Neural Inform. Process. Syst. (2016) 3
60. Park, T., Efros, A.A., Zhang, R., Zhu, J.Y.: Contrastive learning for unpaired image-to-image translation. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IX 16. pp. 319–345. Springer (2020) 4
61. Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y.: Semantic image synthesis with spatially-adaptive normalization. In: IEEE Conf. Comput. Vis. Pattern Recog. (2019) 3
62. Peebles, W., Xie, S.: Scalable diffusion models with transformers. In: Int. Conf. Comput. Vis. (2023) 3
63. Phillips Collection, T.: Vah gogh repetitions. https://www.phillipscollection.org/event/2013-10-11-van-gogh-repetitions (10 2013) 1
64. Po, R., Yang, G., Aberman, K., Wetzstein, G.: Orthogonal adaptation for modular customization of diffusion models. arXiv preprint arXiv:2312.02432 (2023) 3, 4, 7
65. Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., Rombach, R.: Sdxl: Improving latent diffusion models for high-resolution image synthesis. arXiv preprint arXiv:2307.01952 (2023) 3, 5, 10
66. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: Int. Conf. Mach. Learn. (2021) 5, 11
67. Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., Chen, M.: Hierarchical text-conditional image generation with clip latents. arXiv preprint arXiv:2204.06125 (2022) 1, 3
68. Reed, S., Zhang, Y., Zhang, Y., Lee, H.: Deep visual analogy-making. In: Adv. Neural Inform. Process. Syst. (2015) 4

69. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: IEEE Conf. Comput. Vis. Pattern Recog. (2022) 1, 3, 5
70. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical image computing and computer-assisted intervention. pp. 234–241. Springer (2015) 5
71. Ruiz, N., Li, Y., Jampani, V., Pritch, Y., Rubinstein, M., Aberman, K.: Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 22500–22510 (2023) 1, 3
72. Ruiz, N., Li, Y., Jampani, V., Wei, W., Hou, T., Pritch, Y., Wadhwa, N., Rubinstein, M., Aberman, K.: Hyperdreambooth: Hypernetworks for fast personalization of text-to-image models (2023) 4, 14
73. Ryu, S.: Lora-stable diffusion. https://github.com/cloneofsimo/lora (2023) 2, 3
74. Ryu, S.: Low-rank adaptation for fast text-to-image diffusion fine-tuning. https://github.com/cloneofsimo/lora (2023) 3, 5, 7, 10, 13
75. Saharia, C., Chan, W., Chang, H., Lee, C., Ho, J., Salimans, T., Fleet, D., Norouzi, M.: Palette: Image-to-image diffusion models. In: ACM SIGGRAPH 2022 Conference Proceedings. pp. 1–10 (2022) 3
76. Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E., Ghasemipour, S.K.S., Ayan, B.K., Mahdavi, S.S., Lopes, R.G., et al.: Photorealistic text-to-image diffusion models with deep language understanding. In: NeurIPS (2022) 3
77. Sauer, A., Karras, T., Laine, S., Geiger, A., Aila, T.: Stylegan-t: Unlocking the power of gans for fast large-scale text-to-image synthesis. In: International conference on machine learning. pp. 30105–30118. PMLR (2023) 3
78. Schuhmann, C., Vencu, R., Beaumont, R., Kaczmarczyk, R., Mullis, C., Katta, A., Coombes, T., Jitsev, J., Komatsuzaki, A.: Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. arXiv preprint arXiv:2111.02114 (2021) 3
79. Shah, V., Ruiz, N., Cole, F., Lu, E., Lazebnik, S., Li, Y., Jampani, V.: Ziplora: Any subject in any style by effectively merging loras. arXiv preprint arXiv:2311.13600 (2023) 4
80. Shi, J., Xiong, W., Lin, Z., Jung, H.J.: Instantbooth: Personalized text-to-image generation without test-time finetuning. arXiv preprint arXiv:2304.03411 (2023) 4
81. Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., Ganguli, S.: Deep unsupervised learning using nonequilibrium thermodynamics. In: Int. Conf. Mach. Learn. (2015) 4
82. Sohn, K., Ruiz, N., Lee, K., Chin, D.C., Blok, I., Chang, H., Barber, J., Jiang, L., Entis, G., Li, Y., et al.: Styledrop: Text-to-image generation in any style. arXiv preprint arXiv:2306.00983 (2023) 3, 4, 8, 10
83. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. In: Int. Conf. Learn. Represent. (2021) 3
84. Song, Y., Sohl-Dickstein, J., Kingma, D.P., Kumar, A., Ermon, S., Poole, B.: Score-based generative modeling through stochastic differential equations. In: ICLR (2021) 3, 4
85. Tenenbaum, J., Freeman, W.: Separating style and content. Advances in neural information processing systems 9 (1996) 4
86. Tewel, Y., Gal, R., Chechik, G., Atzmon, Y.: Key-locked rank one editing for text-to-image personalization. ACM Transactions on Graphics (TOG) (2023) 3

87. Upchurch, P., Snavely, N., Bala, K.: From a to z: supervised transfer of style and content using deep neural network generators. arXiv preprint arXiv:1603.02003 (2016) 4

88. Valevski, D., Lumen, D., Matias, Y., Leviathan, Y.: Face0: Instantaneously conditioning a text-to-image model on a face. In: SIGGRAPH Asia 2023 Conference Papers. pp. 1–10 (2023) 4

89. Voynov, A., Chu, Q., Cohen-Or, D., Aberman, K.: $p+$: Extended textual conditioning in text-to-image generation. arXiv preprint arXiv:2303.09522 (2023) 3

90. Wang, S.Y., Bau, D., Zhu, J.Y.: Sketch your own gan. In: Int. Conf. Comput. Vis. (2021) 3

91. Wang, S.Y., Bau, D., Zhu, J.Y.: Rewriting geometric rules of a gan. ACM SIGGRAPH (2022) 3

92. Wang, X., Wang, W., Cao, Y., Shen, C., Huang, T.: Images speak in images: A generalist painter for in-context visual learning. In: IEEE Conf. Comput. Vis. Pattern Recog. (2023) 4

93. Wang, X., Yu, J.: Learning to cartoonize using white-box cartoon representations. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 8090–8099 (2020) 10

94. Wei, Y., Zhang, Y., Ji, Z., Bai, J., Zhang, L., Zuo, W.: Elite: Encoding visual concepts into textual embeddings for customized text-to-image generation. In: Int. Conf. Comput. Vis. (2023) 4

95. Ye, H., Zhang, J., Liu, S., Han, X., Yang, W.: Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. arXiv preprint arXiv:2308.06721 (2023) 4, 10, 13

96. Yu, J., Xu, Y., Koh, J.Y., Luong, T., Baid, G., Wang, Z., Vasudevan, V., Ku, A., Yang, Y., Ayan, B.K., et al.: Scaling autoregressive models for content-rich text-to-image generation. Transactions on Machine Learning Research (2022) 3

97. Zhang, L., Agrawala, M.: Adding conditional control to text-to-image diffusion models. In: Int. Conf. Comput. Vis. (2023) 3, 21

98. Zhao, S., Liu, Z., Lin, J., Zhu, J.Y., Han, S.: Differentiable augmentation for data-efficient gan training. In: Adv. Neural Inform. Process. Syst. (2020) 3

99. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: Int. Conf. Comput. Vis. (2017) 4

100. Zou, Z., Shi, T., Qiu, S., Yuan, Y., Shi, Z.: Stylized neural painting. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 15689–15698 (2021) 10

# Appendix

In Section A, we evaluate our method against baselines on the diversity metric, showing that our method leads to more diverse generations comparatively. We also show more qualitative results along with a comparison to the concurrent work of Style Aligned Image Generation [28]. In Section B, we then present details of our style guidance formulation. Finally, in Section C, we provide more implementation details, including the setup for our human preference study and the full synthetic training dataset used for evaluation.

## A    More Quantitative and Qualitative Results

**Diversity metric.**    To measure the overfitting behavior of our method and baselines, we consider a diversity metric. Concretely, we measure the DreamSim [17] perceptual distance between any two images trained with the same style image pair and generated with the same prompt and average results over training pairs and prompts. More formally, we let
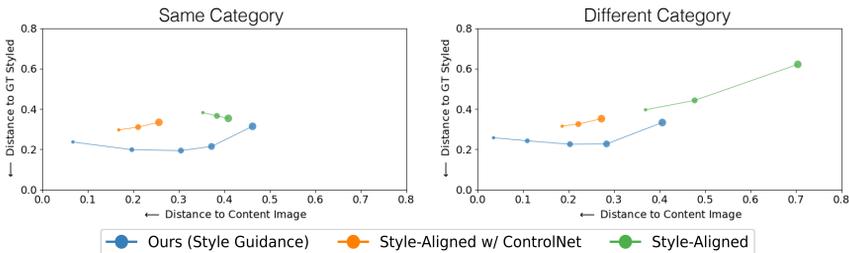
$$\text{DreamSim Diversity} = \mathbb{E}_{S \in \mathcal{S}, P \in \mathcal{P}} \left[ \mathbb{E}_{i_1, i_2 \in \text{data}_{S,P}} \text{DreamSim}(i_1, i_2) \right] \qquad (12)$$

where $\mathcal{S}$ is the set of style image pairs, $\mathcal{P}$ is the set of prompts, and $\text{data}_{S,P}$ is the set of images generated with prompt $P$ by a model customized on style $S$. $\text{DreamSim}(\cdot, \cdot)$ is DreamSim perceptual distance. A decrease in DreamSim Diversity indicates that all images in a certain domain are becoming perceptually similar, which may indicate overfitting to the style training image. Methods that do not overfit the style training image should have higher diversity scores while also having a low perceptual distance to the ground truth testing style images. We present our findings in Figure 9. Our method is able to achieve a low perceptual distance to style ground-truth images while maintaining higher diversity scores. As shown in Figure 6 in the main paper, the baseline results mode collapses to the training image, thus lowering their diversity score as they all become perceptually similar to each other.

**Style Aligned Image Generation [28] Baseline**    This is a recent work for zero-shot style-consistent image generation from an exemplar style image. Given the exemplar style image, it is first inverted to a noise map; then for a new text prompt, the image is generated by attending to both its own self-attention map and the self-attention map from the style exemplar at every denoising step. We compare against this baseline by using the style image in our training image pair as an exemplar and generating a new style image with a new text prompt using this method. Optionally, we condition this generation on the edge map of the newly generated image without attention sharing using ControlNet [97] to help with content preservation. We show the qualitative results of our method compared to all the variants of this baseline in Figure 12. Figures 10 and 11 show quantitative comparison, where our method outperforms this baseline in terms of both style similarity and diversity metric. We achieve lower perceptual
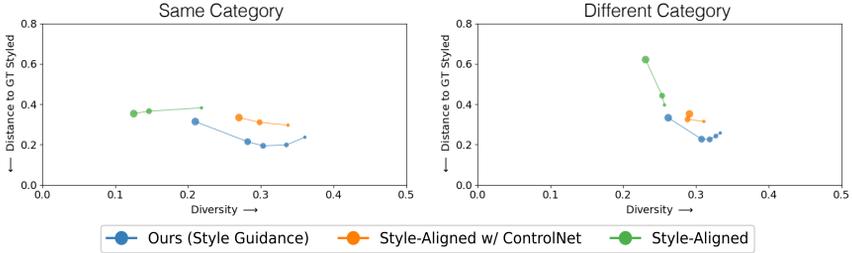
Fig. 9: **Quantitative comparison on Diversity metric.** Our method with style guidance has high diversity and low perceptual distance to ground truth style images both on the same category as training (left) and when evaluated on categories different from training, e.g., trained on human portraits but tested on dog images (right). Methods without edge control tend to lose diversity indicating overfitting, and methods with edge control have similar/higher diversity, but much worse style application. Increased marker size corresponds to an increase in style guidance scale.



Fig. 10: **Style similarity with Style Aligned [28].** Our method Pareto dominates both versions of Style Aligned Image Generation for image generation both on the same category as training (left) and when evaluated on categories different from training, e.g., trained on human portraits but tested on dog images (right).

distance to the style ground-truth images, low perceptual distance from content images, and high diversity.

**Extra Qualitative Evaluation**     We compare our method to non finetuning-based methods in Figure 12. We observe that these methods perform worse than finetuning-based methods, especially when generating images in a different category to the training style image. Secondly, we compare our method with the highest-performing baselines, but use our style guidance (Equation 10 ) to apply stylization during inference for these baselines. We present our results in Figure 13. First, we notice that using style guidance for adding style allows the baseline methods to better preserve original content over LoRA scale (Figure 13 vs Figure 6). While adding our style guidance is better able to preserve content while applying style for baseline methods, our full method is still able to outperform baselines with style guidance applied.

**Fig. 11: Image diversity with Style Aligned [28] on learned style (Diversity).** Our method has high diversity and low perceptual distance to ground truth style images both on the same category as training (left) and when evaluated on categories different from training, e.g., trained on human portraits but tested on dog images (right) as compared to both versions of Style Aligned Image Generation.

## B    Style Guidance Details

In this section, we derive our style guidance formulation. We consider the probability of latent $\mathbf{x}$ with multiple conditionings [8], i.e., the text prompt $c_t$ and a class of style images $c_{\text{style}}$. First, we apply Bayes' rule:

$$P(\mathbf{x}|c_t, c_{\text{style}}) = \frac{P(\mathbf{x}, c_t, c_{\text{style}})}{P(c_t, c_{\text{style}})} = \frac{P(c_{\text{style}}|c_t, \mathbf{x})P(c_t|\mathbf{x})P(\mathbf{x})}{P(c_t, c_{\text{style}})} \quad (13)$$

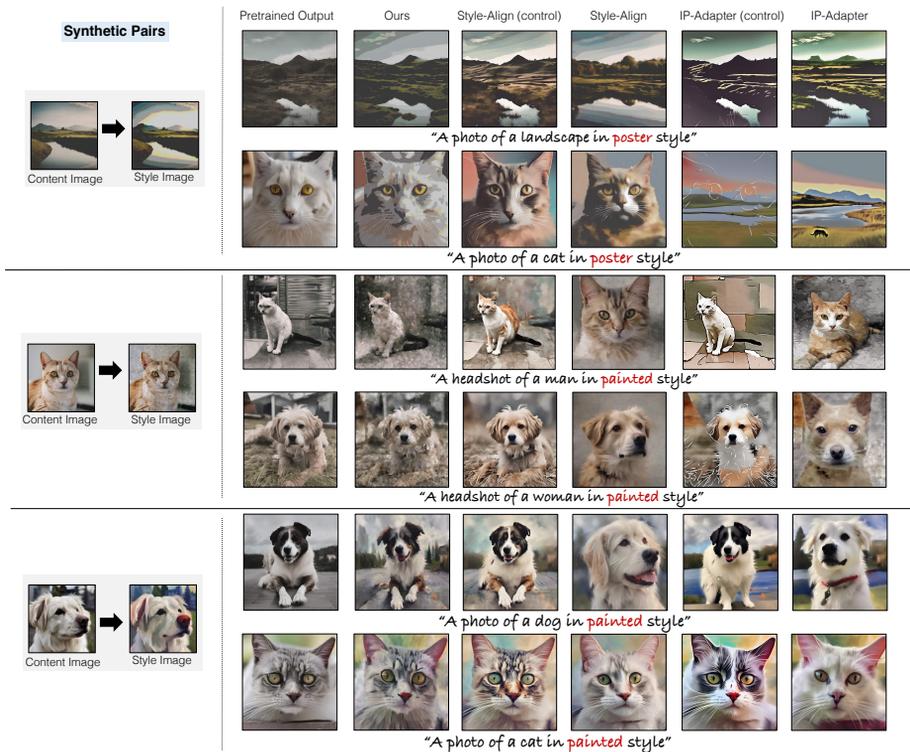Applying logarithm on both sides, we get:

$$\begin{aligned} \log(P(\mathbf{x}|c_t, c_{\text{style}})) \\ = \log(P(c_{\text{style}}|c_t, \mathbf{x})) + \log(P(c_t|\mathbf{x})) + \log(P(\mathbf{x})) \\ - \log(P(c_t, c_{\text{style}})) \end{aligned} \quad (14)$$

Next, we take the derivative with respect to $\mathbf{x}$:

$$\begin{aligned} &\nabla_{\mathbf{x}} \log(P(\mathbf{x}|c_t, c_{\text{style}})) \\ =&\nabla_{\mathbf{x}} \log(P(c_{\text{style}}|c_t, \mathbf{x})) + \nabla_{\mathbf{x}} \log(P(c_t|\mathbf{x})) + \nabla_{\mathbf{x}} \log(P(\mathbf{x})) \\ =& \nabla_{\mathbf{x}} \log\left(\frac{P(c_{\text{style}}, c_t, \mathbf{x})}{P(c_t, \mathbf{x})}\right) + \nabla_{\mathbf{x}} \log\left(\frac{P(c_t, \mathbf{x})}{P(\mathbf{x})}\right) + \nabla_{\mathbf{x}} \log(P(\mathbf{x})) \\ =& (\nabla_{\mathbf{x}} \log P(c_{\text{style}}, c_t, \mathbf{x}) - \nabla_{\mathbf{x}} \log P(c_t, \mathbf{x})) \\ &+ (\nabla_{\mathbf{x}} \log P(c_t, \mathbf{x}) - \nabla_{\mathbf{x}} \log P(\mathbf{x})) \\ &+ (\nabla_{\mathbf{x}} \log(P(\mathbf{x}))) \end{aligned} \quad (15)$$
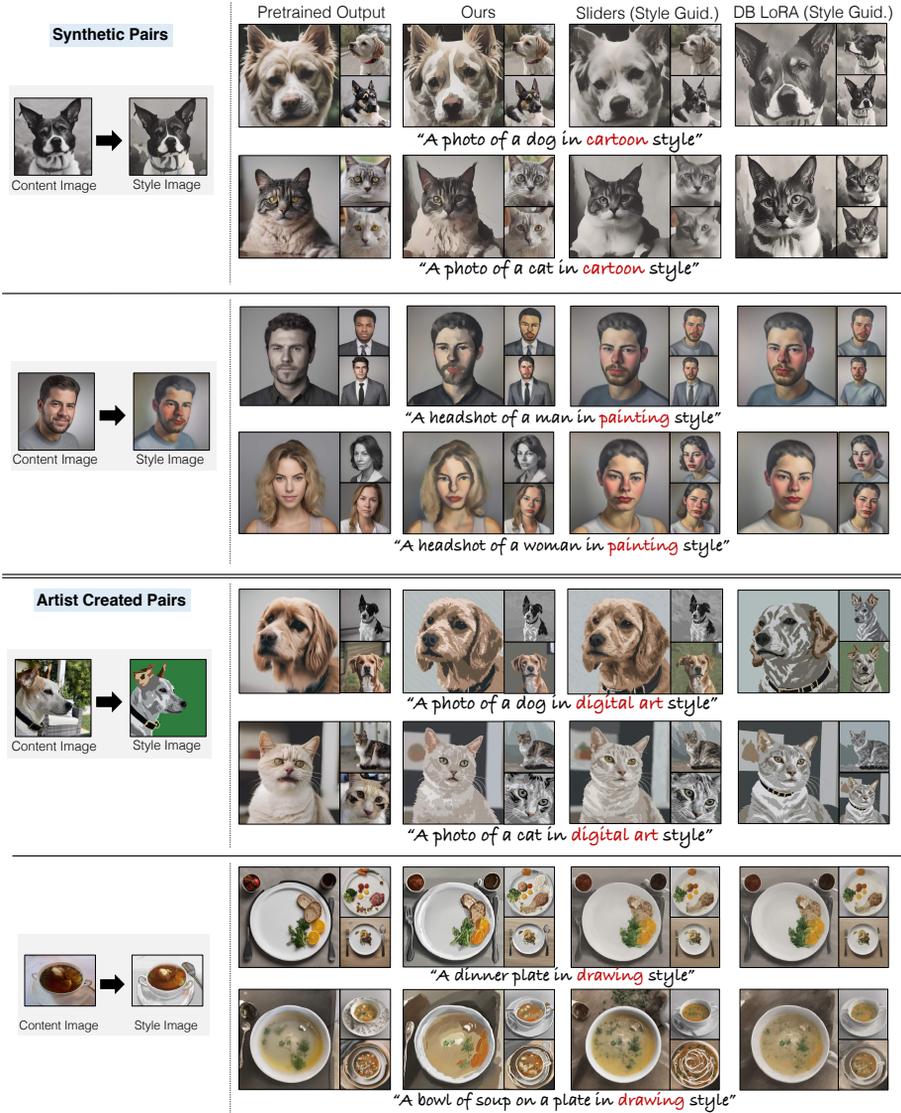
As usual, we approximate $\nabla_{\mathbf{x}} (\log P(c_t, \mathbf{x}))$ via $\epsilon_\theta(\mathbf{x}_t, c_t)$ and $\nabla_{\mathbf{x}} \log(P(\mathbf{x}))$ via $\epsilon_\theta(\mathbf{x}_t, \varnothing)$. **Importantly, we approximate**

$$\nabla_{\mathbf{x}} \log(P(\mathbf{x}_t|c_t, \mathbf{c}_{\text{style}})) \approx \epsilon_{\theta_{\text{style}}}(\mathbf{x}_t, c_{t,\text{style}}) \quad (16)$$

**Fig. 12:** Result of our method compared to the methods without finetuning (zoom in for best viewing). For all methods, we consider adding the edgemap from the pretrained output as an extra conditioning using ControlNet. Without this edgemap, other methods tend to lose the structure of the pretrained output. In some cases, however, an additional edgemap can overly constrain the output of a model, like in the second and fourth stylistic image pairs. Our method preserves the structure of the Stable Diffusion image, while faithfully applying the desired style. We use style guidance strength 3 and classifier guidance strength 5 for our method and set the IP-adapter scale and style-alignment scale to 0.5.

where $c_t$ is the original text prompt, $\mathbf{c}_{\text{style}}$ is the class of stylized images from the training style, $\theta_{\text{style}}$ is the UNet with style LoRA adapters applied, and $c_{t,\text{style}} = $ "$\{c_t\}$ in <desc> style". Here, we use $c_t$ to push the prediction in the text direction, and both text conditioning ("in <desc> style") and low-rank adapters ($\theta_{\text{style}}$) to push the prediction into the class of images in the artist's
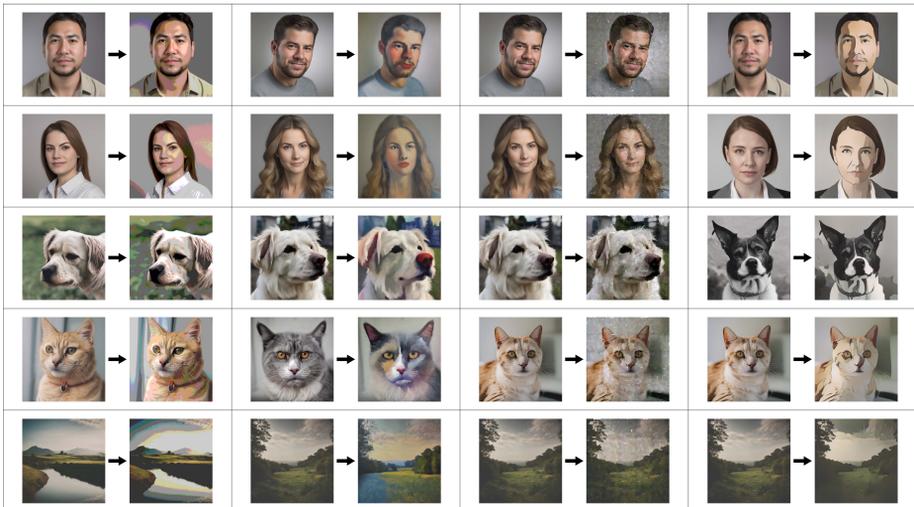
**Fig. 13:** Result of our method compared to the strongest baselines, but replacing LoRA scale (Eq. 3) with our style guidance (Eq. 10) for the baselines. While our style guidance increases baseline performance over LoRA scale images displayed in Figure 6, our method is still superior in terms of preserving content while applying style.

style denoted by $c_{\text{style}}$. Following this, our new score estimate is:

$$\hat{\epsilon}_\theta(\mathbf{x}_t, c_t, \mathbf{c}_{\text{style}}) = \epsilon_\theta(\mathbf{x}_t, \varnothing)$$
$$+ \lambda_{\text{cfg}}(\epsilon_\theta(\mathbf{x}_t, c_t) - \epsilon_\theta(\mathbf{x}_t, \varnothing)) \tag{17}$$
$$+ \lambda_{\text{style}}(\epsilon_{\theta_{\text{style}}}(\mathbf{x}_t, c_{t,\text{style}}) - \epsilon_\theta(\mathbf{x}_t, c_t)) \tag{18}$$
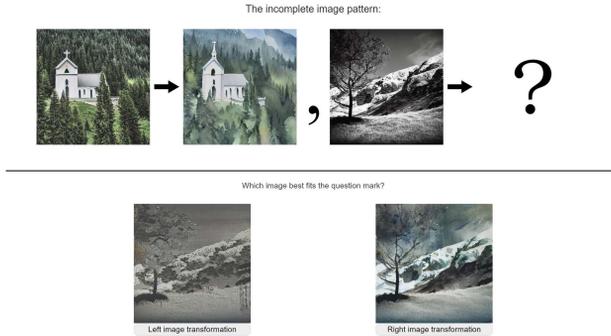
Fig. 14: **Training Data.** We present the synthetic training data set used for evaluation, where each pair is used as a single training instance. Each column corresponds to a different style, and each row corresponds to a different content category.

$\lambda_{\text{cfg}}$ and $\lambda_{\text{style}}$ are guidance scales that can be varied as in classifier free guidance [31]. Given a fixed $\lambda_{\text{cfg}}$, we can vary the $\lambda_{\text{style}}$ term as desired to generate an original guidance $\lambda_{\text{cfg}}$ image with varying amounts of style. Notice that at $\lambda_{\text{cfg}} = \lambda_{\text{style}}$, the $\epsilon_\theta(\mathbf{x}_t, \varnothing)$ terms cancel and we are left with the original classifier guidance.

## C   Implementation Details

**Training data.**   We present our full training set of 20 different style transformations in Figure 14. Each image pair is a standalone training instance used in our method. We consider four different styles (posterization, impressionist, neural painting, cartoonization), with each column corresponding to a single style. For each style, we consider five categories for training (man, woman, dog, cat, landscape).

**Mechanical Turk details.**    When running Amazon Mechanical Turk, we prompt users with an analogy-style interface. First, we provide the training pair of images, followed by the testing content image, and two options for possible styled examples. After viewing both images, users choose either the left or right image. Figure 15 shows an example. Each individual user is presented with four training examples, as in Figure 15, followed by 16 random testing examples comparing our method with one of our baselines. We survey 75 users for each of the 16 individual studies and use bootstrapping to obtain variance estimates. In total, we collect 19200 user samples. For each method, we pick a stylization hyperparameter based on Figure 4. For details, see Table 1

**Fig. 15:** Mturk User Interface

| Method | Hyperparameter value | |
|---|---|---|
| | Same Category | Different Category |
| Ours (Style Guid.) | 3 | 4 |
| Ours w/ Orthog (Style Guid.) | 3 | 4 |
| DB LoRA (Style Guid.) | 2 | 4 |
| DB LoRA (LoRA Scale) | 0.4 | 0.8 |
| Concept Sliders (Style Guid.) | 2 | 4 |
| Concept Sliders (LoRA Scale) | 0.6 | 0.8 |
| StyleDrop LoRA (LoRA Scale) | 0.6 | 1 |
| IP Adapter w/T2I (Image Guidance) | 0.5 | 0.5 |
| IP Adapter (Image Guidance) | 0.5 | 0.5 |

**Table 1: Experiment Hyperparameters.** We choose a fixed stylization hyperparameter for our own model and each baseline when generating images for Mechanical Turk. When picking a hyperparameter, we try and optimize tradeoffs between style application and content preservation, informed by Figure 4 in the main body. Our style guidance (Equation 10) generally takes values from 0 to $\lambda_{\text{cfg}} = 5$, while all other stylization hyperparameters generally take values 0 to 1.