

Style Aligned Image Generation via Shared Attention

Amir Hertz^{*1}, Andrey Voynov^{*1}, Shlomi Fruchter^{†1}, and Daniel Cohen-Or^{†1,2}

¹ Google Research

² Tel Aviv University

Abstract

Large-scale Text-to-Image (T2I) models have rapidly gained prominence across creative fields, generating visually compelling outputs from textual prompts. However, controlling these models to ensure consistent style remains challenging, with existing methods necessitating fine-tuning and manual intervention to disentangle content and style. In this paper, we introduce *StyelAligned*, a novel technique designed to establish style alignment among a series of generated images. By employing minimal ‘attention sharing’ during the diffusion process, our method maintains style consistency across images within T2I models. This approach allows for the creation of style-consistent images using a reference style through a straightforward inversion operation. Our method’s evaluation across diverse styles and text prompts demonstrates high-quality synthesis and fidelity, underscoring its efficacy in achieving consistent style across various inputs.

1. Introduction

Large-scale Text-to-Image (T2I) generative models [43, 45, 51] have emerged as an essential tool across creative disciplines, such as art, graphic design, animation, architecture, gaming and more. These models show tremendous capabilities of translating an input text into an appealing visual result that is aligned with the input description.

An envisioned application of T2I models revolves around the rendition of various concepts in a way that shares a consistent style and character, as though all were created by the same artist and method (see Fig. 1). While proficient in aligning with the textual description of the style, state-of-the-art T2I models often create images that diverge significantly in their interpretations of the same stylistic descriptor, as depicted in Fig. 2.

Recent methods mitigate this by fine-tuning the T2I model over a set of images that share the same style [16, 55]. This optimization is computationally expensive and usually

^{*}Equal contribution.

[†]Equal Advising.

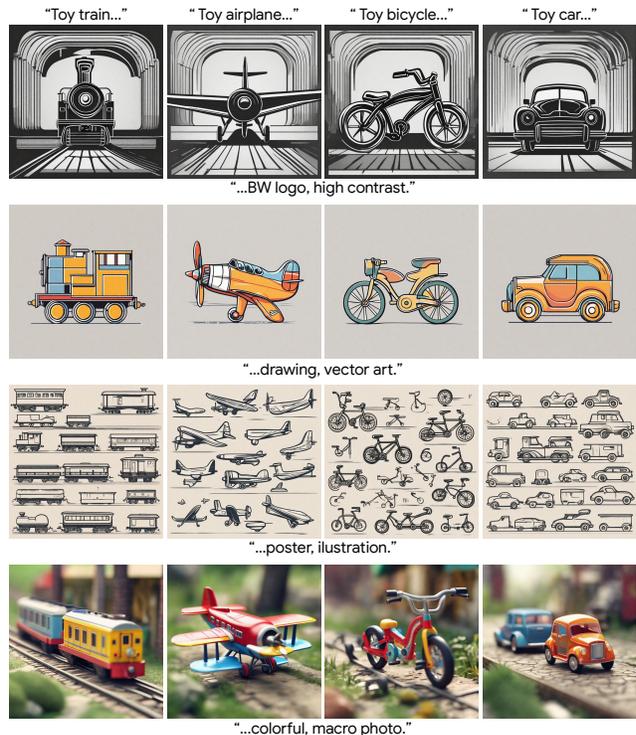


Figure 1. **Style aligned image set generation.** By fusing the features of the toy train image (left) during the diffusion process, we can generate an image set of different content that shares the style.

requires human input in order to find a plausible subset of images and texts that enables the disentanglement of content and style.

We introduce *StyelAligned*, a method that enables *consistent style interpretation* across a set of generated images (Fig. 1). Our method requires no optimization and can be applied to any attention-based text-to-image diffusion model. We show that adding minimal *attention sharing* operations along the diffusion process, from each generated image to the first one in a batch, leads to a style-consistent set. Moreover, using diffusion inversion, our method can be applied to generate style-consistent images given a reference style image, with no optimization or fine-tuning.



Figure 2. **Standard text-to-image vs. StyleAligned set generation.** Given style description of “minimal origami”, standard text-to-image generation (top) results with an unaligned image set while our method (bottom) can generate variety of style aligned content.

We present our results over diverse styles and text prompts, demonstrating high-quality synthesis and fidelity to the prompts and reference style. We show diverse examples of generated images that share their style with a reference image that can possibly be a given input image. Importantly, our technique stands as a zero-shot solution, distinct from other personalization techniques, as it operates without any form of optimization or fine-tuning. For our code and more examples, please visit the project page style-aligned-gen.github.io

2. Related Work

Text-to-image generation. Text conditioned image generative models [10, 37, 44] show unprecedented capabilities of generating high quality images from text descriptions. In particular, T2I diffusion models [41, 44, 52] are pushing the state of the art and they are quickly adopted for different generative visual tasks like inpainting [5, 50], image-to-image translation [61, 66], local image editing [12, 28], subject-driven image generation [48, 57] and more.

Attention Control in diffusion models. Hertz et al. [20] have shown how cross and self-attention maps within the diffusion process determine the layout and content of the generated images. Moreover, they showed how the attention maps can be used for controlled image generation. Other studies have leveraged modifications in attention lay-

ers to enhance the fidelity or diversity of generated images [11, 40], or apply attention control for image editing [8, 15, 36, 38, 39, 59]. However, in contrast to prior approaches that primarily enable structure-preserving image editing, our method excels at generating images with diverse structures and content while maintaining a consistent style interpretation.

Style Transfer. Transferring a style from a reference image to a target content image is well studied subject in computer graphics. Classic works [13, 14, 22, 31] rely on optimization of handcrafted features and texture resampling algorithms from an input texture image, combined with structure constrains of a content image. With the progress of deep learning research, another line of works utilizes deep neural priors for style transfer optimization using deep features of pre-trained networks [18, 58], or injecting attention features from a style image to a target one [4]. More related to our approach, Huang et al. [26] introduced a real time style transfer network based on Adaptive Instance Normalization layers (AdaIN) that are used to normalize deep features of a target image using deep features statistics of a reference style image. Follow-up works, employ the AdaIN layer for additional unsupervised learning tasks, like style-based image generation [29] and Image2Image translation [27, 34].

T2I Personalization To generalize T2I over new visual concepts, several works developed different optimization techniques over a small collection of input images that share the same concept [16, 19, 48, 62]. In instances where the collection shares a consistent style, the acquired concept becomes the style itself, affecting subsequent generations. Most close to our work is StyleDrop [55], a style personalization method that relies on fine-tuning of light weight adapter layers [24] at the end of each attention block in a non-autoregressive generative text-to-image transformer [10]. StyleDrop can generate a set of images in the same style of by training the adapter layers over a collection of images that share the same style. However, it struggles to generate a consistent image set of different content when training on a single image.

Our method can generate a consistent image set without optimization phase and without relying on several images for training. To skip the training phase, recent works developed dedicated personalization encoders [17, 32, 53, 65, 66] that can directly inject new priors from a single input image to the T2I model. However, these methods encounter challenges to disentangle style from content as they focus on generating the same subject as in the input image.

3. Method overview

In the following section we start with an overview of the T2I diffusion process, and in particular the self-attention mechanism Sec. 3.1. We continue by present-

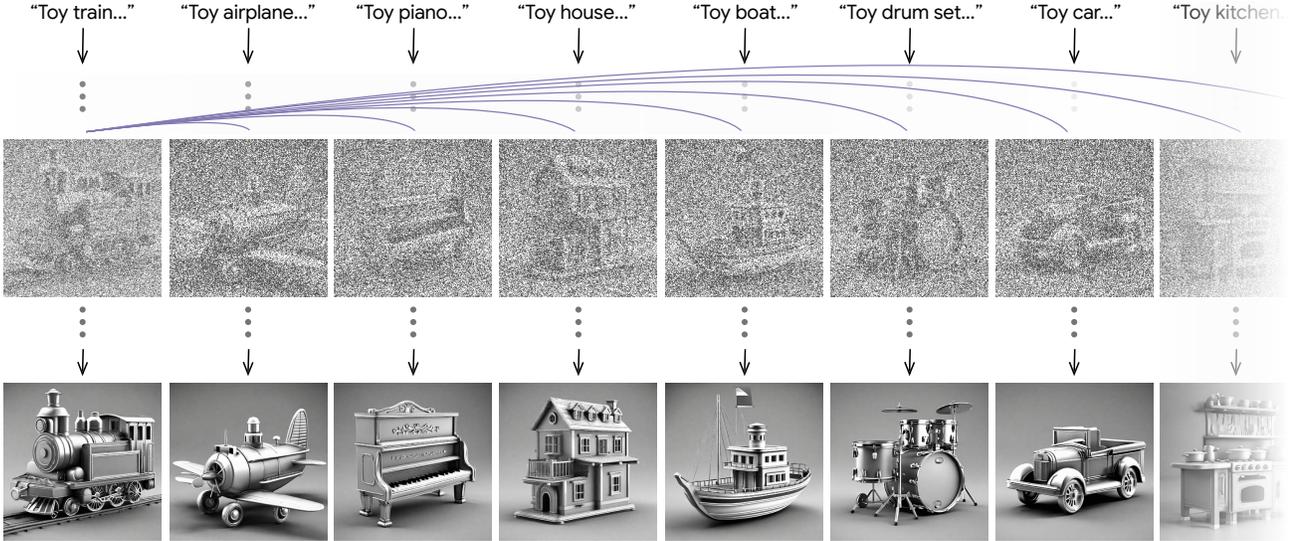


Figure 3. **Style Aligned Diffusion.** Generation of images with a style aligned to the reference image on the left. In each diffusion denoising step all the images, except the reference, perform a shared self-attention with the reference image.

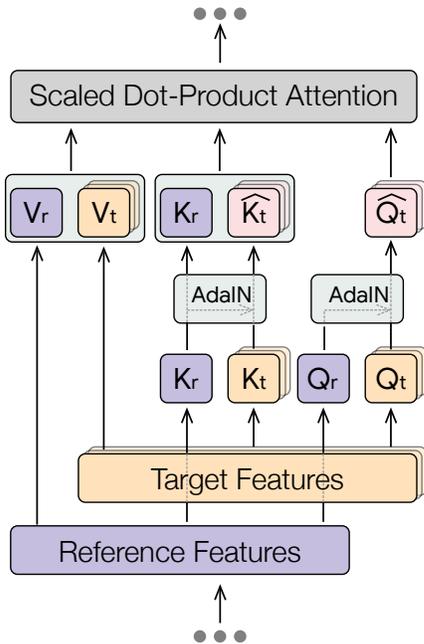


Figure 4. **Shared attention layer.** The target images attend to the reference image by applying AdaIN over their queries and keys using the reference queries and keys respectively. Then, we apply shared attention where the target features are updated by both the target values V_t and the reference values V_r .

ing our attention-sharing operation within the self-attention layers that enable style aligned image set generation.

3.1. Preliminaries

Diffusion models [23, 54] are generative latent variable models that aim to model a distribution $p_\theta(x_0)$ that approximates the data distribution $q(x_0)$ and are easy to sample from. Diffusion models are trained to reverse the diffusion “forward process”:

$$x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon, \quad \epsilon \sim N(0, I),$$

where $t \in [0, \infty)$ and the values of α_t are determined by a scheduler such that $\alpha_0 = 1$ and $\lim_{t \rightarrow \infty} \alpha_t = 0$. During inference, we sample an image by gradually denoising an input noise image $x_T \sim N(0, I)$ via the reverse process:

$$x_{t-1} = \mu_{t-1} + \sigma_t z, \quad z \sim N(0, I),$$

where the value of σ_t is determined by the sampler and μ_{t-1} is given by

$$\mu_{t-1} = \frac{\sqrt{\alpha_{t-1}}x_t}{\sqrt{\alpha_t}} + \left(\sqrt{1 - \alpha_{t-1}} - \frac{\sqrt{1 - \alpha_t}}{\sqrt{\alpha_t}} \right) \epsilon_\theta(x_t, t),$$

where $\epsilon_\theta(x_t, t)$ is the output of a diffusion model parameterized by θ .

Moreover, this process can be generalized for learning a marginal distribution using an additional input condition. That leads text-to-image diffusion models (T2I), where the output of the model $\epsilon_\theta(x_t, t, y)$ is conditioned on a text prompt y .

Self-Attention in T2I Diffusion Models. State-of-the-art T2I diffusion models [7, 41, 52] employ a U-Net architecture [46] that consists of convolution layers and transformer attention blocks [60]. In these attention mechanisms, deep



Figure 5. **Ablation study – qualitative comparison.** Each pair of rows shows two sets of images generated by the same set of prompts “...in minimal flat design illustration” using different configurations of our method, and each row in a pair uses a different seed. Sharing the self-attention between all images in the set (bottom) results with some diversity loss (style collapse across many seeds) and content leakage within each set (colors from one image leak to another). Disabling the queries-keys AdaIN operation results with less consistent image sets compared to our full method (top) which keeps on both diversity between different sets and consistency within each set.

image features $\phi \in \mathbb{R}^{m \times d_h}$ attend to each other via self-attention layers and to contextual text embedding via cross-attention layers.

Our work operates at the self-attention layers where deep features are being updated by attending to each other. First, the features are projected into queries $Q \in m \times d_k$, keys $K \in m \times d_k$ and values $V \in m \times d_h$ through learned linear layers. Then, the attention is computed by the scaled

dot-product attention:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} V \right),$$

where d_k is the dimension of Q and K . Intuitively, each image feature is updated by a weighted sum of V , where the weight depends on the correlation between the projected query q and the keys K . In practice, each self-attention layer consists of several attention heads, and then the residual is computed by concatenating and projecting the attention heads output back to the image feature space d_h :

$$\hat{\phi} = \phi + \text{Multi-Head-Attention}(\phi).$$

3.2. Style Aligned Image Set Generation

The goal of our method is to generate a set of images $\mathcal{I}_1 \dots \mathcal{I}_n$ that are aligned with an input set of text prompts $y_1 \dots y_n$ and share a consistent style interpretation with each other. For example, see the garnered image set of toy objects in Fig. 3 that are style-aligned with each other and to the input text on top. A naïve way to generate a style aligned image set of different content is to use the same style description in the text prompts. As can be seen at the bottom of Fig. 2, generating different images using a shared style description of “in minimal origami style” results in an unaligned set, since each image is unaware of the exact appearance of other images in the set during the generation process.

The key insight underlying our approach is the utilization of the self-attention mechanism to allow communication among various generated images. This is achieved by sharing attention layers across the generated images.

Formally, let Q_i , K_i , and V_i be the queries, keys, and values, projected from deep features ϕ_i of \mathcal{I}_i in the set, then, the attention update for ϕ_i is given by:

$$\text{Attention}(Q_i, K_{1\dots n}, V_{1\dots n}), \quad (1)$$

$$\text{where } K_{1\dots n} = \begin{bmatrix} K_1 \\ K_2 \\ \vdots \\ K_n \end{bmatrix} \text{ and } V_{1\dots n} = \begin{bmatrix} V_1 \\ V_2 \\ \vdots \\ V_n \end{bmatrix}. \text{ However, we}$$

have noticed that by enabling full attention sharing, we may harm the quality of the generated set. As can be seen in Fig. 5 (bottom rows), full attention sharing results in content leakage among the images. For example, the unicorns got green paint from the garnered dino in the set. Moreover, full attention sharing results with less diverse sets of the same set of prompts, see the two sets in Fig. 5 in bottom rows compared to the sets above.

To restrict the content leakage and allow diverse sets, we share the attention to only one image in the generated set

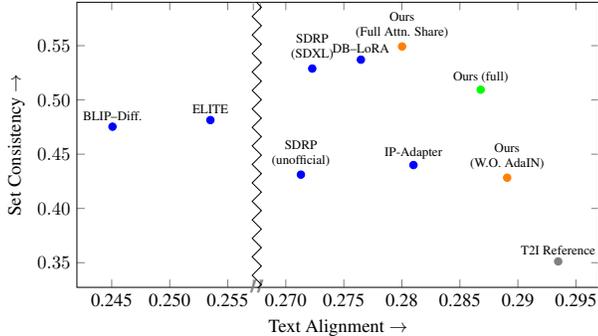


Figure 6. **Quantitative Comparison.** We compare the results of the different methods (blue marks) and our ablation experiments (orange marks) in terms of text alignment (CLIP score) and set consistency (DINO embedding similarity).

(typically the first in the batch). That is, *target* image features ϕ_t are attending to themselves and to the features of only one *reference* image in the set using Eq. 1. As can be seen in Fig. 5 (middle), sharing the attention to only one image in the set results in diverse sets that share a similar style. However, in that case, we have noticed that the style of different images is not well aligned. We suspect that this is due to low attention flow from the reference to the target image.

As illustrated in Fig. 4, to enable balanced attention reference, we normalize the queries Q_t and keys K_t of the target image using the queries Q_r and keys K_r of the reference image using the adaptive normalization operation (AdaIN) [26]:

$$\hat{Q}_t = \text{AdaIN}(Q_t, Q_r) \quad \hat{K}_t = \text{AdaIN}(K_t, K_r),$$

where the AdaIn operation is given by:

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

and $\mu(x), \sigma(x) \in \mathbb{R}^{d_k}$ are the mean and the standard deviation of queries and keys across different pixels. Finally, our shared attention is given by

$$\text{Attention}(\hat{Q}_t, Q_t K_{rt}^T, V_{rt}),$$

where $K_{rt} = \begin{bmatrix} K_r \\ \hat{K}_t \end{bmatrix}$ and $V_{rt} = \begin{bmatrix} V_r \\ V_t \end{bmatrix}$.

4. Evaluations and Experiments

We have implemented our method over Stable Diffusion XL (SDXL) [41] by applying our attention sharing overall 70 self-attention layers of the model. The generation of a four images set takes 29 seconds on a single A100 GPU. Notice that since the generation of the reference image is

Table 1. **User evaluation for style aligned image set generation.** In each question, the user was asked to select between two image sets, Which is better in terms of style consistency and match to the text descriptions (see Sec. 4). We report the percentage of judgments in favor of *StyleAligned* over 800 answers (2400 in total).

StyleDrop (unofficial MUSE)	StyleDrop (SDXL)	DreamBooth-LoRA (SDXL)
85.2 %	67.1 %	61.3%

not influenced by other images in the batch, we can generate larger sets by fixing the prompt and seed of the reference image across the set generation.

For example, see the sets in Fig. 2 and 3.

Evaluation set. With the support of ChatGPT, we have generated 100 text prompts describing different image styles over four random objects. For example, “{A guitar, A hot air balloon, A sailboat, A mountain} in papercut art style.” For each style and set of objects, we use our method to generate a set of images. The full list of prompts is provided in the appendix.

Metrics. To verify that each generated image contains its specified object, we measure the CLIP cosine similarity [42] between the image and the text description of the object. In addition, we evaluate the style consistency of each generated set, by measuring the pairwise average cosine similarity between DINO VIT-B/8 [9] embeddings of the generated images in each set. Following [47, 62], we used DINO embeddings instead of CLIP image embeddings for measuring image similarity, since CLIP was trained with class labels and therefore it might give a high score for different images in the set that have similar content but with a different style. On the other hand, DINO better distinguishes between different styles due to its self-supervised training.

4.1. Ablation Study

The quantitative results are summarized in Fig. 6, where the right-top place on the chart means better text similarity and style consistency, respectively. As a reference, we report the score obtained by generating the set of images using SDXL (T2I Reference) using the same seeds without any intervention. As can be seen, our method achieves a much higher style consistency score at the expense of text similarity. See qualitative comparison in Fig. 2.

In addition, we compared our method to additional two variants of the shared attention as described in Sec. 3.2. The first variant uses full attention sharing (*Full Attn. Share*) where the keys and values are shared between each pair of images in the set. In the second variant (*W.A. AdaIN*) we omit the AdaIN operation over queries and keys. As expected, this *Full Attn. Share* variant, results with higher style consistency and lower text alignment. As can be seen

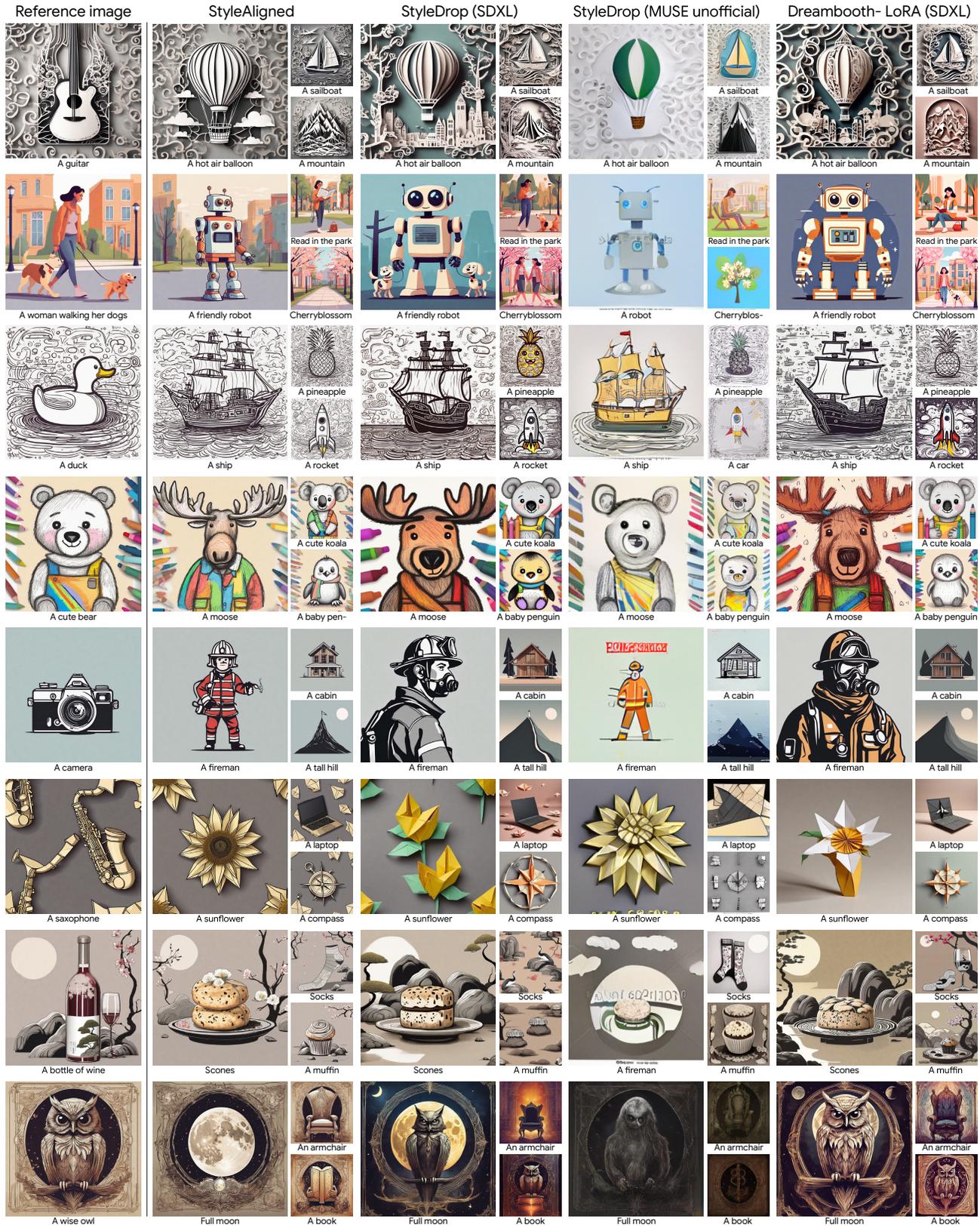


Figure 7. Qualitative comparison to personalization based methods.

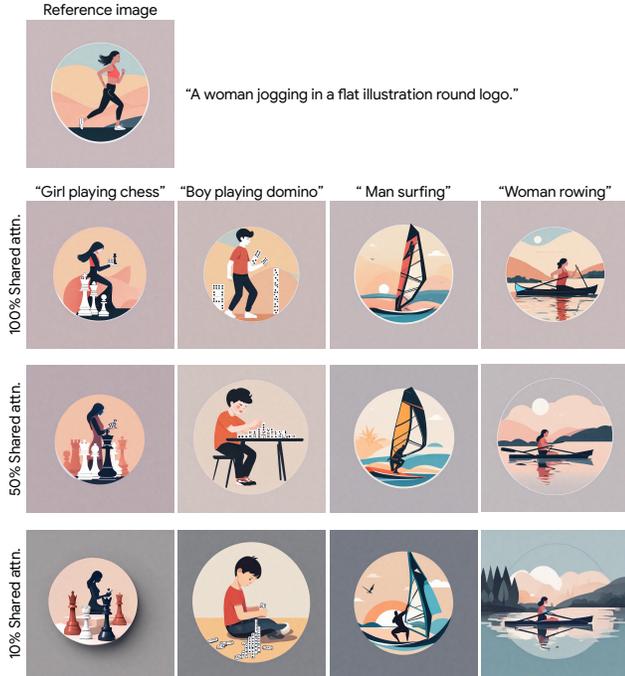


Figure 8. **Varying level of attention sharing.** By reducing the number of shared attention layers, i.e., allowing only self-attention in part of the layers, we can get more varied results (bottom rows) at the expense of style alignment (top row).

in Fig. 5, *Full Attn. Share* harms the overall quality of the image sets and diversity across sets. Moreover, our method without the use of AdaIN results in much lower style consistency. Qualitative results can be seen in Fig. 5.

4.2. Comparisons

For baselines, we compare our method to T2I personalization methods. We trained StyleDrop [55] and DreamBooth [47] over the first image in each set of our evaluation data, and use the trained personalized weights to generate the additional three images in each set. We use a public unofficial implementation of StyleDrop ¹ (SDRP-unofficial) over non-regressive T2I model. Due to the large quality gap between the unofficial MUSE model ² to the official one [10], we follow StyleDrop and implement an adapter model over SDXL (SDRP-SDXL), where we train a low rank linear layer after each Feed-Forward layer at the model’s attention blocks. For training DreamBooth, we adapt the LoRA [25, 49] variant (DB-LoRA) over SDXL using the public huggingface-diffusers implementation ³. We follow the hyperparameters tuning reported in [55] and train both SDRP-SDXL and DB-LoRA for 400 steps to prevent overfitting to the style training image.

¹github.com/aim-uofa/StyleDrop-PyTorch

²github.com/baaivision/MUSE-Pytorch

³github.com/huggingface/diffusers

As can be seen in the qualitative comparison, Fig. 7, the image sets generated by our method, are more consistent across style attributes like color palette, drawing style, composition, and pose. Moreover, the personalization-based methods may leak the content of the training reference image (on the left) when generating the new images. For example, see the repeated woman and dogs in the results of DB-LoRA and SDRP-SDXL at the second row or the repeated owl at the bottom row. Similarly, because of the content leakage, these methods obtained lower text similarity scores and higher set consistency scores compared to our method.

We also apply two encoder-based personalization methods ELITE [64], IP-Adapter [66], and BLIP-Diffusion [32] over our evaluation set. These methods receive as input the first image in each set and use its embeddings to generate images with other content. Unlike the optimization-based techniques, these methods operate in a much faster feed-forward diffusion loop, like our method. However, as can be seen in Fig. 6, their performance for style aligned image generation is poor compared to the other baselines. We argue that current encoder-based personalization techniques struggle to disentangle the content and the style of the input image. We supply qualitative results in appendix C.

User Study. In addition to the automatic evalua-

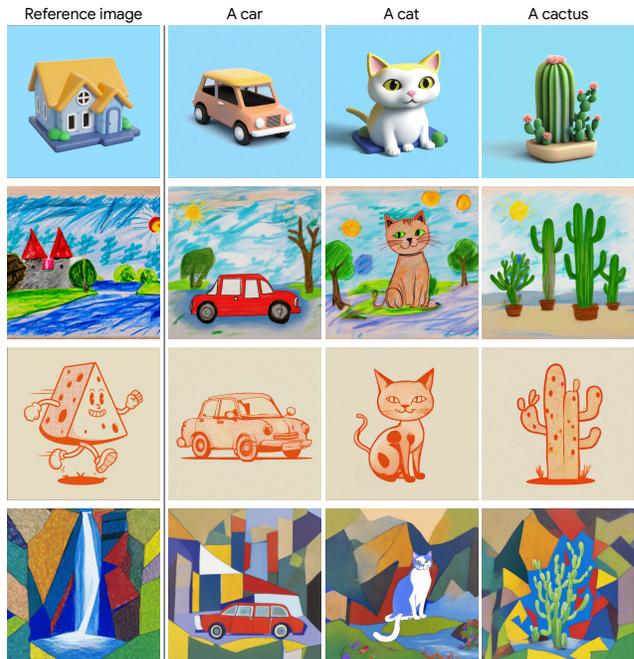


Figure 9. **Style aligned image generation to an input image.** Given an input reference image (left column) and text description, we first apply DDIM inversion over the image to get the inverted diffusion trajectory $x_T, x_{T-1} \dots x_0$. Then, starting from x_T and a new set of prompts, we apply our method to generate new content (right columns) with an aligned style to the input.

tion, we conducted a user study over the results of our method, StyleDrop (unofficial MUSE), StyleDrop (SDXL), and DreamBooth–LoRA (SDXL). In each question, we randomly sample one of the evaluation examples and show the user the 4 image set that resulted from our and another method (in a random order). The user had to choose which set is better in terms of style consistency, and text alignment. A print screen of the user study format is provided in the appendix. Overall, we collected 2400 answers from 100 users using the Amazon Mechanical Turk service. The results are summarized in Tab. 1 where for each method, we report the percentage of judgments in our favor. As can be seen, most participants favored our method by a large margin. More information about our user study can be found in appendix D.

4.3. Additional Results

Style Alignment Control. We provide means of control over the style alignment to the reference image by applying the shared attention over only part of the self-attention layers. As can be seen in Fig. 8, reducing the number of shared attention layers results with a more diverse image set, which still shares common attributes with the reference image.

StyelAligned from an Input Image. To generate style-aligned images to an input image, we apply DDIM inversion [56] using a provided text caption. Then, we apply our method to generate new images in the style of the input using the inverted diffusion trajectory x_T, x_{T-1}, \dots, x_0 for the reference image. Examples are shown in Fig. 9, 13, where we use BLIP captioning [33] to get a caption for each input image. For example, we used the prompt “A render of a house with a yellow roof” for the DDIM inversion of the top example and replaced the word house with other objects to generate the style-aligned images of a car, a cat, and a cactus. Notice that this method does not require any optimization. However, DDIM inversion may fail [36] or results with an erroneous trajectory [28]. More results and analysis, are provided in appendix A

Shared Self-Attention Visualization. Figure 10 depicts the self-attention probabilities from a generated target image to the reference style image. In each of the rows, we pick a point on the image and depict the associated probabilities map for the token at this particular point. Notably probabilities mapped on the reference image are semantically close to the query point location. This suggests that the self-attention tokens sharing do not perform a global style transfer, but rather match the styles in a semantically meaningful way [4]. In addition, Figure 11 visualizes the three largest components of the average shared attention maps of the rhino image, encoded in RGB channels. Note that the shared attention map is composed of both self-attention and cross-image attention to the giraffe. As can be seen, the components highlight semantically related regions like the

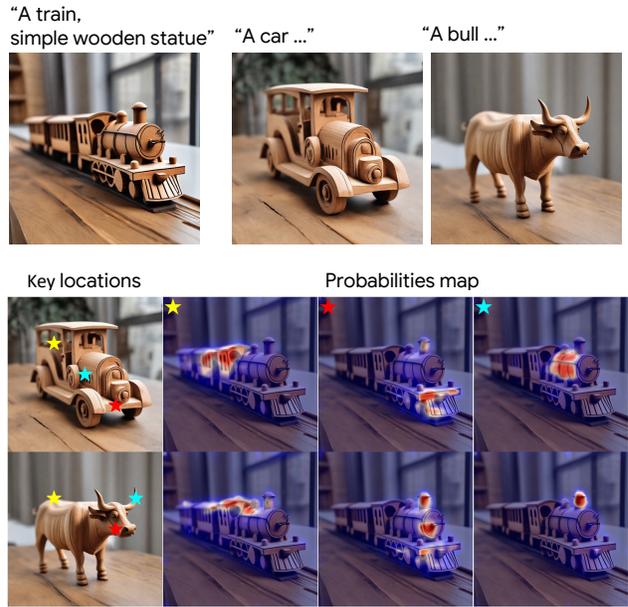


Figure 10. *Self-Attention probabilities maps from different generated image locations (**Key locations** column) to the reference train image with the target style (top-left).*

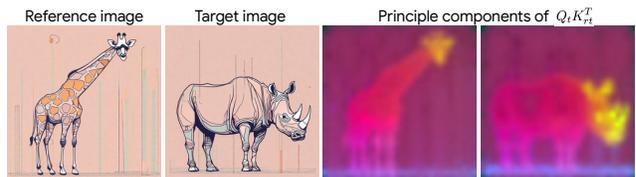


Figure 11. **Principle components of the shared attention map.** *On right, we visualize the principle components of the shared attention map between the reference giraffe and the target rhino generated images. The three largest components of the shared maps are encoded in RGB channels.*

bodies, heads, and the background in the images.

StyelAligned with Other Methods. Since our method doesn’t require training or optimization, it can be easily combined on top of other diffusion based methods to generate style-consistent image sets. Fig. 12 shows several such examples where we combine our method with ControlNet [67], DreamBooth [48] and MultiDiffusion [6]. More examples and details about the integration of StyleAligned with other methods can be found in appendix B.

5. Conclusions

We have presented StyelAligned, which addresses the challenge of achieving style-aligned image generation within the realm of large-scale Text-to-Image models. By introducing minimal attention sharing operations with AdaIN modulation during the diffusion process, our method

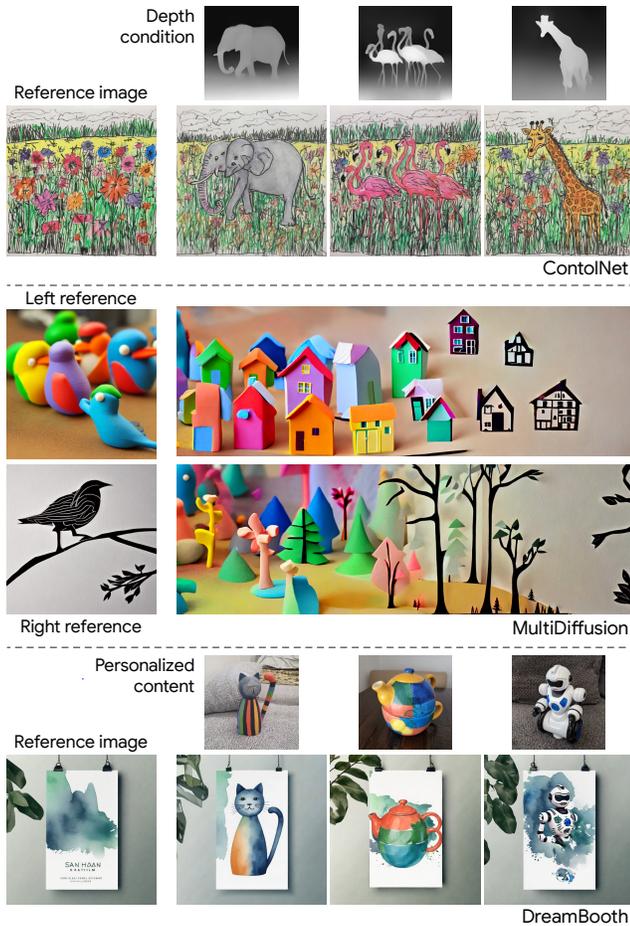


Figure 12. **StyleAligned with other methods.** On top, *StyleAligned* is combined with *ControlNet* to generate style-aligned images conditioned on depth maps. In the middle, our method combined with *MultiDiffusion* to generate panorama images that share multiple styles. On the bottom, style consistent and personalized content created by combining our method with pre-trained personalized *DreamBooth-LoRA* models.

successfully establishes style-consistency and visual coherence across generated images. The demonstrated efficacy of *StyleAligned* in producing high-quality, style-consistent images across diverse styles and textual prompts underscores its potential in creative domains and practical applications. Our results affirm *StyleAligned* capability to faithfully adhere to provided descriptions and reference styles while maintaining impressive synthesis quality.

In the future we would like to explore the scalability and adaptability of *StyleAligned* to have more control over the shape and appearance similarity among the generated images. Additionally, due to the limitation of current diffusion inversion methods, a promising direction is to leverage *StyleAligned* to assemble a style-aligned dataset which then can be used to train style condition text-to-image models.

6. Acknowledgement

We thank Or Patashnik, Matan Cohen, Yael Pritch, and Yael Vinker for their valuable inputs that helped improve this work.

References

- [1] Diffusers: controlnet-depth-sdxl-1.0. <https://huggingface.co/diffusers/controlnet-depth-sdxl-1.0>, 2023. 13
- [2] Diffusers: controlnet-openpose-sdxl-1.0. <https://huggingface.co/thibaud/controlnet-openpose-sdxl-1.0>, 2023. 13
- [3] Diffusers: Multidiffusion pipeline. <https://huggingface.co/docs/diffusers/api/pipelines/panorama>, 2023. 13
- [4] Yuval Alaluf, Daniel Garibi, Or Patashnik, Hadar Averbuch-Elor, and Daniel Cohen-Or. Cross-image attention for zero-shot appearance transfer, 2023. 2, 8
- [5] Omri Avrahami, Ohad Fried, and Dani Lischinski. Blended latent diffusion. *ACM Trans. Graph.*, 42(4), jul 2023. 2
- [6] Omer Bar-Tal, Lior Yariv, Yaron Lipman, and Tali Dekel. Multidiffusion: Fusing diffusion paths for controlled image generation. 2023. 8, 13, 18
- [7] James Betker, Gabriel Goh, Li Jing, TimBrooks, Jianfeng Wang, Linjie Li, LongOuyang, JuntangZhuang, JoyceLee, YufeiGuo, WesamManassra, PrafullaDhariwal, CaseyChu, YunxinJiao, and Aditya Ramesh. Improving image generation with better captions. 2023. 3
- [8] Mingdeng Cao, Xintao Wang, Zhongang Qi, Ying Shan, Xiaohu Qie, and Yinqiang Zheng. MasaCtrl: tuning-free mutual self-attention control for consistent image synthesis and editing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 22560–22570, October 2023. 2
- [9] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021. 5
- [10] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, José Lezama, Lu Jiang, Ming Yang, Kevin P. Murphy, William T. Freeman, Michael Rubinstein, Yanzhen Li, and Dilip Krishnan. Muse: Text-to-image generation via masked generative transformers. In *International Conference on Machine Learning*, 2023. 2, 7
- [11] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *ACM Transactions on Graphics (TOG)*, 42:1 – 10, 2023. 2
- [12] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance. In *The Eleventh International Conference on Learning Representations*, 2022. 2
- [13] Alexei A Efros and William T Freeman. Image quilting for texture synthesis and transfer. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pages 571–576. 2023. 2

Reference Image



Figure 13. Various remarkable places depicted with the style taken from Bruegel’s “*The Tower of Babel*”.
 Top row: Rome Colosseum, Rio de Janeiro, Seattle Space Needle.

- [14] Alexei A Efros and Thomas K Leung. Texture synthesis by non-parametric sampling. In *Proceedings of the seventh IEEE international conference on computer vision*, volume 2, pages 1033–1038. IEEE, 1999. 2
- [15] Dave Epstein, Allan Jabri, Ben Poole, Alexei A Efros, and Aleksander Holynski. Diffusion self-guidance for controllable image generation. *arXiv preprint arXiv:2306.00986*, 2023. 2
- [16] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *The Eleventh International Conference on Learning Representations*, 2022. 1, 2
- [17] Rinon Gal, Moab Arar, Yuval Atzmon, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. Encoder-based domain tuning for fast personalization of text-to-image models. *ACM Transactions on Graphics (TOG)*, 42(4):1–13, 2023. 2
- [18] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2414–2423, 2016. 2
- [19] Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris N. Metaxas, and Feng Yang. Svdiff: Compact parameter space for diffusion fine-tuning. *ArXiv*, abs/2303.11305, 2023. 2, 13
- [20] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. 2
- [21] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or. Prompt-to-prompt image editing with cross-attention control. In *The Eleventh International Conference on Learning Representations*, 2022. 14
- [22] Aaron Hertzmann, Charles E Jacobs, Nuria Oliver, Brian Curless, and David H Salesin. Image analogies. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pages 557–570. 2023. 2
- [23] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Proc. NeurIPS*, 2020. 3
- [24] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer

- learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR, 2019. 2
- [25] Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2021. 7, 13
- [26] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision*, pages 1501–1510, 2017. 2, 5
- [27] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 172–189, 2018. 2
- [28] Inbar Huberman-Spiegelglas, Vladimir Kulikov, and Tomer Michaeli. An edit friendly ddpn noise space: Inversion and manipulations. *arXiv preprint arXiv:2304.06140*, 2023. 2, 8
- [29] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 2
- [30] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1931–1941, 2023. 13
- [31] Hochang Lee, Sanghyun Seo, Seungtaek Ryoo, and Kyunghyun Yoon. Directional texture transfer. In *Proceedings of the 8th International Symposium on Non-Photorealistic Animation and Rendering*, pages 43–48, 2010. 2
- [32] Dongxu Li, Junnan Li, and Steven C. H. Hoi. BLIP-Diffusion: Pre-trained subject representation for controllable text-to-image generation and editing. *ArXiv*, abs/2305.14720, 2023. 2, 7, 13
- [33] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation, 2022. 8
- [34] Ming-Yu Liu, Xun Huang, Arun Mallya, Tero Karras, Timo Aila, Jaakko Lehtinen, and Jan Kautz. Few-shot unsupervised image-to-image translation. In *IEEE International Conference on Computer Vision (ICCV)*, 2019. 2
- [35] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *ICLR*, 2021. 14
- [36] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6038–6047, 2023. 2, 8
- [37] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning*, 2021. 2
- [38] Dong Huk Park*, Grace Luo*, Clayton Toste, Samaneh Azadi, Xihui Liu, Maka Karalashvili, Anna Rohrbach, and Trevor Darrell. Shape-guided diffusion with inside-outside attention. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024. 2
- [39] Gaurav Parmar, Krishna Kumar Singh, Richard Zhang, Yijun Li, Jingwan Lu, and Jun-Yan Zhu. Zero-shot image-to-image translation. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–11, 2023. 2
- [40] Or Patashnik, Daniel Garibi, Idan Azuri, Hadar Averbuch-Elor, and Daniel Cohen-Or. Localizing object-level shape variations with text-to-image diffusion models. *ArXiv*, abs/2303.11306, 2023. 2
- [41] Dustin Podell, Zion English, Kyle Lacey, A. Blattmann, Tim Dockhorn, Jonas Muller, Joe Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image synthesis. *ArXiv*, abs/2307.01952, 2023. 2, 3, 5
- [42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, 2021. 5
- [43] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 1
- [44] Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10674–10685, 2021. 2
- [45] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1
- [46] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015. 3
- [47] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. DreamBooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023. 5, 7
- [48] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Wei Wei, Tingbo Hou, Yael Pritch, Neal Wadhwa, Michael Rubinstein, and Kfir Aberman. Hyperdreambooth: Hypernetworks for fast personalization of text-to-image models. *arXiv preprint arXiv:2307.06949*, 2023. 2, 8, 13
- [49] Simo Ryu. Low-rank adaptation for fast text-to-image diffusion fine-tuning. <https://github.com/cloneofsimon/lora>, 2022. 7

- [50] Chitwan Saharia, William Chan, Huiwen Chang, Chris A. Lee, Jonathan Ho, Tim Salimans, David J. Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. *ACM SIGGRAPH 2022 Conference Proceedings*, 2021. 2
- [51] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 1
- [52] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 2, 3
- [53] Jing Shi, Wei Xiong, Zhe L. Lin, and Hyun Joon Jung. Instantbooth: Personalized text-to-image generation without test-time finetuning. *ArXiv*, abs/2304.03411, 2023. 2
- [54] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015. 3
- [55] Kihyuk Sohn, Nataniel Ruiz, Kimin Lee, Daniel Castro Chin, Irina Blok, Huiwen Chang, Jarred Barber, Lu Jiang, Glenn Entis, Yuanzhen Li, et al. Styledrop: Text-to-image generation in any style. *arXiv preprint arXiv:2306.00983*, 2023. 1, 2, 7
- [56] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2020. 8
- [57] Yoad Towel, Rinon Gal, Gal Chechik, and Yuval Atzmon. Key-locked rank one editing for text-to-image personalization. *SIGGRAPH 2023 Conference Proceedings*, 2023. 2, 13
- [58] Narek Tumanyan, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Splicing vit features for semantic appearance transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10748–10757, 2022. 2
- [59] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1921–1930, 2023. 2
- [60] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. 3
- [61] Andrey Voynov, Kfir Aberman, and Daniel Cohen-Or. Sketch-guided text-to-image diffusion models. *arXiv preprint arXiv:2211.13752*, 2022. 2



Figure 14. **Text-to-image generation with explicit style description.** Unlike our approach, this fails to produce fine and style-aligned results. See Fig. 13 to inspect our method results.

- [62] Andrey Voynov, Q. Chu, Daniel Cohen-Or, and Kfir Aberman. P+: Extended textual conditioning in text-to-image generation. *ArXiv*, abs/2303.09522, 2023. 2, 5
- [63] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. In *CVPR*, 2016. 13
- [64] Yuxiang Wei. Official implementation of ELITE. <https://github.com/csyxwei/ELITE>, 2023. 7, 13
- [65] Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. ELITE: Encoding visual concepts into textual embeddings for customized text-to-image generation. *ArXiv*, abs/2302.13848, 2023. 2
- [66] Hu Ye, Jun Zhang, Siyi Liu, Xiao Han, and Wei Yang. IP-Adapter: Text compatible image prompt adapter for text-to-image diffusion models. *ArXiv*, abs/2308.06721, 2023. 2, 7, 13
- [67] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, 2023. 8, 13

Appendix

A. StyleAligned from an Input Image

Figure 13 shows our techniques being applied for style transfer for the Peter Bruegels' "The Tower of Babel" to multiple places around the world. As for the prompt we always use the places' followed by "Pieter Bruegel Painting", e.g. "Rome Coliseum, Pieter Bruegel Painting". Even though the original masterpiece is known to model, it fails to reproduce its style with only text guidance. Fig. 14 shows some of the places generated with the direct instruction to resemble the original painting, without self-attention sharing. Notably, the model fails to produce an accurate style alignment with the original picture.

Further examples of style transferring from real examples are presented in Figures 17 and 18.

We also noticed that once the style transfer is performed from an extremely famous image, the default approach may sometimes completely ignore the target prompt, generating an image almost identical to the reference. We suppose that this happens because the outputs of the denoising model for the famous reference image have very high confidence and activations magnitudes. Thus in the shared self-attention, most of the attention is taken by the reference keys. To compensate for it, we propose the simple trick of the attention scores rescaling. In the self-attention sharing mechanism, for some fixed scale $\lambda < 1$, we rescale the queries and keys products conducting the new scores $\lambda \cdot \langle Q, K_{\text{target}} \rangle$. We apply this only to the reference image keys. First, this suppresses extra-high keys. Also, this makes the attention scores more uniformly distributed, encouraging the generated image to capture style aggregated from the whole reference image. Fig. 15 demonstrates the rescaling factor variation effect for the particularly popular reference "Starr Night" by Van Gogh. Notably, without rescaling, the model generates an image almost identical to the reference, while the scale relaxation produces a plausible transfer.

B. Integration with Other Methods

Below, we show different examples where our method can provide style aligned image generation capability on top of different diffusion-based image generation methods.

Style Aligned Subject Driven Generation. To use our method on top of a personalized diffusion model, first, given a collection of images (3-6) of the personalized content, we follow DreamBooth–LoRA training [25,48] where the layers of the attention layers are fine-tuned via low-rank adaptation weights (LoRA). Then, during inference, we apply our method by sharing the attention of personalized generated images with a generated reference style image. During this process, the LoRA weights are used only for the generation of personalized content. Examples of style aligned personalized images are shown in Fig. The results of our method on top of different personalized models are shown in Fig. 19 where in each column we fine-tuned the SDXL model over the image collection on top and generated the personalized content with the reference images on the left. It can be seen that in some cases, like in the backpack photos on the right, the subject in the image remains in the same style as in the original photos. This is a known limitation of training-based personalization methods [57] which we believe can be improved by applying our method over other T2I personalization techniques [19,30] or more careful search for training hyperparameters that allow better generalization of the personalized model to different styles.

Style Aligned MultiDiffusion Image Generation. Bar et al. [6] presented MultiDiffusion, a method for generating images in any resolution by aggregating diffusion predictions of overlapped squared crops. Our method can be

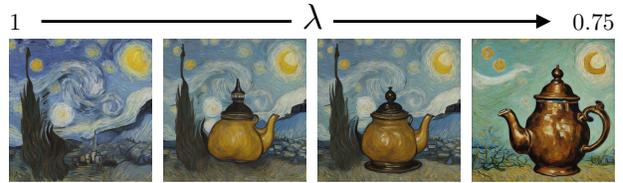


Figure 15. Reference attention rescaling factor variation used for extremely popular reference image assets.

used on top of MultiDiffusion by enabling our shared attention between the crops to a reference image that is generated in parallel. Fig. 20 shows style aligned panorama images generated with MultiDiffusion in conjunction with our method using the public implementation of MultiDiffusion over Stable Diffusion V2 [3]. Notice that compared to a *vanilla* MultiDiffusion image generation (small images in 20), our method not only enables the generation of style aligned panoramas but also helps to preserve the style within each image.

StyleAligned with Additional Conditions. Lastly, we show how our method can be combined with ControlNet [67] which enriches the conditioning signals of diffusion text-to-image generation to include additional inputs, like depth map and pose. ControlNet injects the additional information by predicting residual features that are added to the diffusion image features outputs of the down and middle U-Net blocks. Similar to previous modifications, we apply StyleAligned image generation by sharing the attention of ControlNet conditioned images to a reference image that isn't conditioned on additional input. Fig. 21 shows style aligned image set (different rows) that are conditioned on depth maps (different columns) using ControlNet depth encoder over SDXL [1]. Fig. 22 shows style aligned image set (different rows) that are conditioned on pose estimation obtained by OpenPose [63] (different columns) using ControlNet pose encoder over SDXL [2].

C. Additional Comparisons

We provide additional comparisons of our method to encoder-based text-to-image personalization methods and editing approaches over the evaluation set presented in Section 4 in the main paper. Table 2 summarized the full quantitative results presented in the paper and here.

Encoder Based Approaches As reported in the paper, we compare our method to encoder-based text-to-image personalization methods: BLIP-Diffusion [32], ELITE [64], and IP-Adapter [66]. These methods train an image encoder and fine-tune the T2I diffusion model to be conditioned on visual input. Fig. 23 shows a qualitative comparison on the same set shown in the paper (Fig. 7). As can be seen, our image sets are more consistent and aligned to the reference. Notice that, currently, only IP-Adapter provides an encoder model for Stable Diffusion XL (SDXL). Nev-

Table 2. **Full quantitative comparison for style aligned image generation.** We evaluate the generated image sets in terms of text alignment (CLIP score) and set consistency (DINO embedding similarity). $\pm X$ denotes the standard deviation of the score across 100 image set results.

Method	Text Alignment (CLIP \uparrow)	Set Consistency (DINO \uparrow)
StyleDrop (SDXL)	0.272 \pm 0.04	0.529 \pm 0.15
StyleDrop (unofficial MUSE)	0.271 \pm 0.04	0.301 \pm 0.14
DreamBooth-LoRA (SDXL)	0.276 \pm 0.03	0.537 \pm 0.17
IP-Adapter (SDXL)	0.281 \pm 0.03	0.44 \pm 0.13
ELITE (SD 1.4)	0.253 \pm 0.03	0.481 \pm 0.13
BLIP-Diffusion (SD 1.4)	0.245 \pm 0.04	0.475 \pm 0.12
Prompt-to-Prompt (SDXL)	0.283 \pm 0.03	0.454 \pm 0.18
SDEdit (SDXL)	0.274 \pm 0.03	0.453 \pm 0.16
StyleAligned (SDXL)	0.287 \pm 0.03	0.51 \pm 0.14
StyleAligned (W.O. AdaIN)	0.289 \pm 0.03	0.428 \pm 0.14
StyleAligned (Full Attn.)	0.28 \pm 0.03	0.55 \pm 0.15

ertheless, BLIP-Diffusion and ELITE struggle to produce consistent image sets that match the text descriptions.

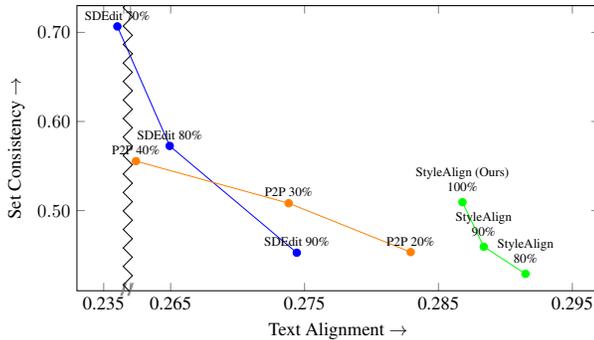


Figure 16. **Quantitative Comparison to zero shot editing approaches.** We compare the results of the different methods in terms of text alignment (CLIP score) and set consistency (DINO embedding similarity).

Zero Shot Editing Approaches Other baselines that can be used for style aligned image set generation are diffusion-based editing methods when applied over the reference images. However, unlike our method, these methods assume structure preservation of the input image. We report the results of two diffusion-based editing approaches: SDEdit [35] and Prompt-to-Prompt (P2P) [21] in Fig. 16. Notice that similar to our method, these methods provide a level of control that trade-off between alignment to text and alignment to the input image. To get higher text alignment, SDEdit can be applied over an increased percentage of diffusion steps, and P2P can reduce the number of attention injection steps. Our method can achieve higher text alignment, as described in Section 4 in the main paper, by using our shared attention over only a subset of self-attention layers. Fig. 16 presents the trade-off of the results over the different methods. As can be seen, only our method

can achieve text alignment while preserving high set consistency.

D. User Study and Evaluation Settings

As described in the main paper, we generate the images for evaluation using a list of 100 text prompts where each prompt describes 4 objects in the same style. The full list is provided at the end of supplementary materials D. We evaluated the results of the different methods using the automatic CLIP and DINO scores and through user evaluation. The format of the user study is provided in Fig. 24 where the user has to select between the results of two methods. For each method from StyleDrop (SDXL), StyleDrop (unofficial Muse), and DreamBooth-LoRA (SDXL), we collected 800 answers compared to our results. In total, we collected 2400 answers from 100 participants.



Figure 17. Samples of the proposed style transfer techniques applied for a variety of different images and target prompts.

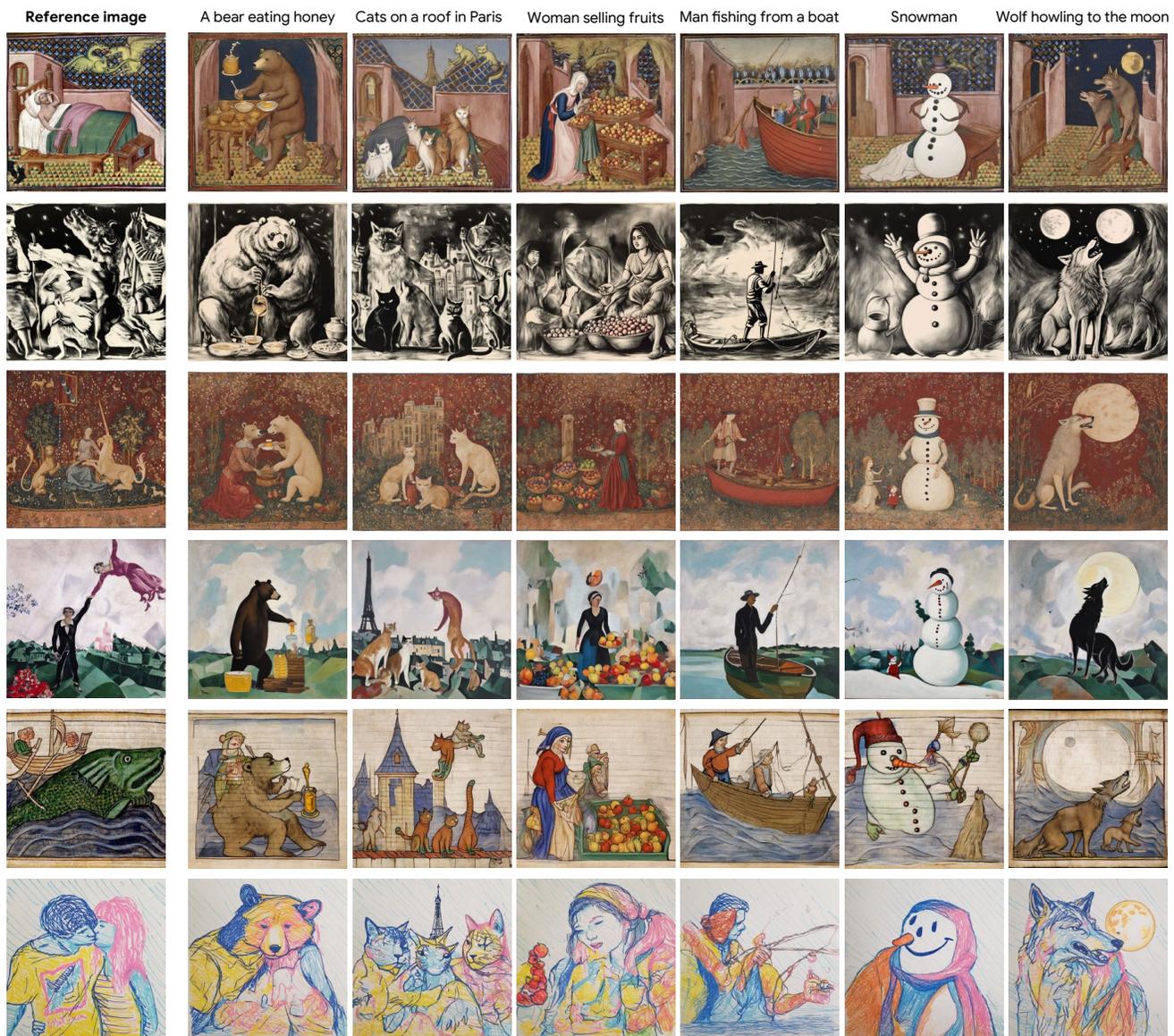


Figure 18. Samples of the proposed style transfer techniques applied for a variety of different images and target prompts.



Figure 19. **Personalized T2I diffusion with StyleAligned.** Each row shows style aligned image st using the reference image on the left, applied on different personalized diffusion models, fine-tuned over the personalized content on top. The top two rows were generated using the prompt "[my subject] in the style of a beautiful papercut art." The bottom two rows were generated using the prompt "[my subject] in beautiful flat design." where [my subject] is replaced with the subject name.

Reference image



"A poster in a flat design style."



"Houses in a flat design style."



"Mountains in a flat design style."



"Girrafes in a flat design style."

Reference image



"A poster in a papercut art style."



"A village in a papercut art style."



"Futuristic city scape in a papercut art style."



"A jungle in a papercut art style."

Figure 20. **MultiDiffusion with StyleAligned.** The panoramas were generated with MultiDiffusion [6] using the text prompt beneath and the left image as reference. The small images in the bottom right corners are the results of MultiDiffusion results without our method.

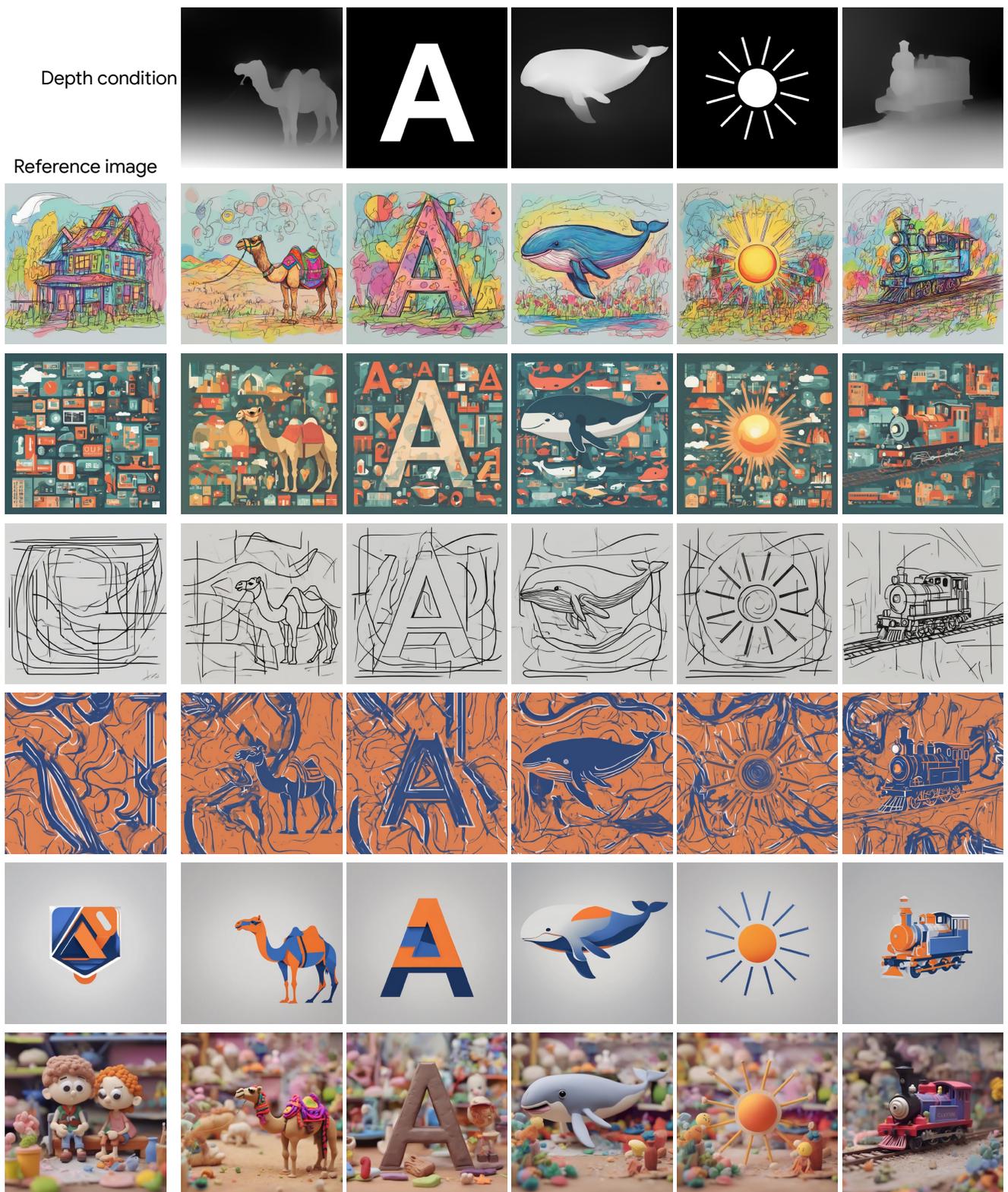


Figure 21. ControlNet Depth with StyleAligned.

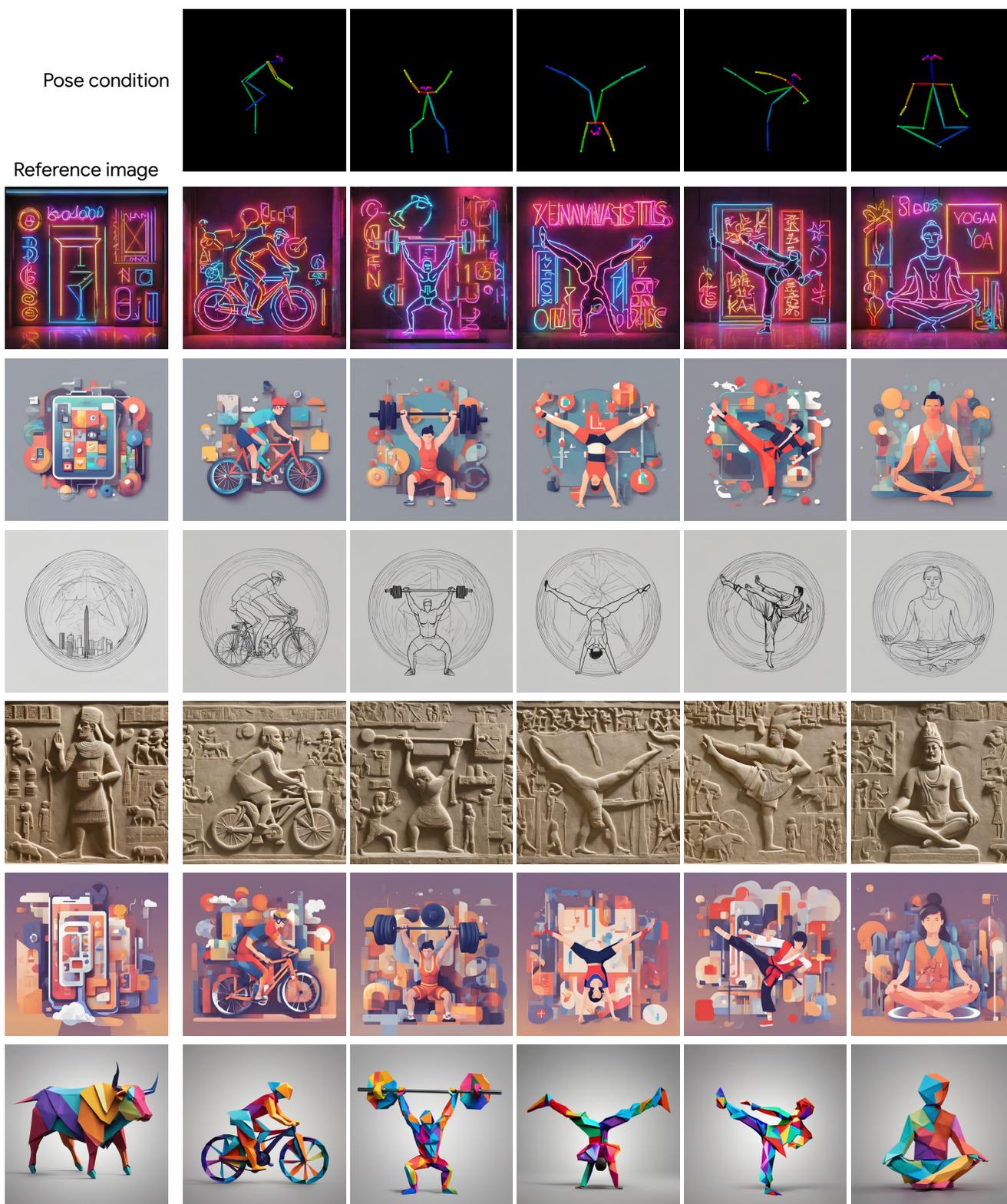


Figure 22. ControlNet pose with StyleAligned.

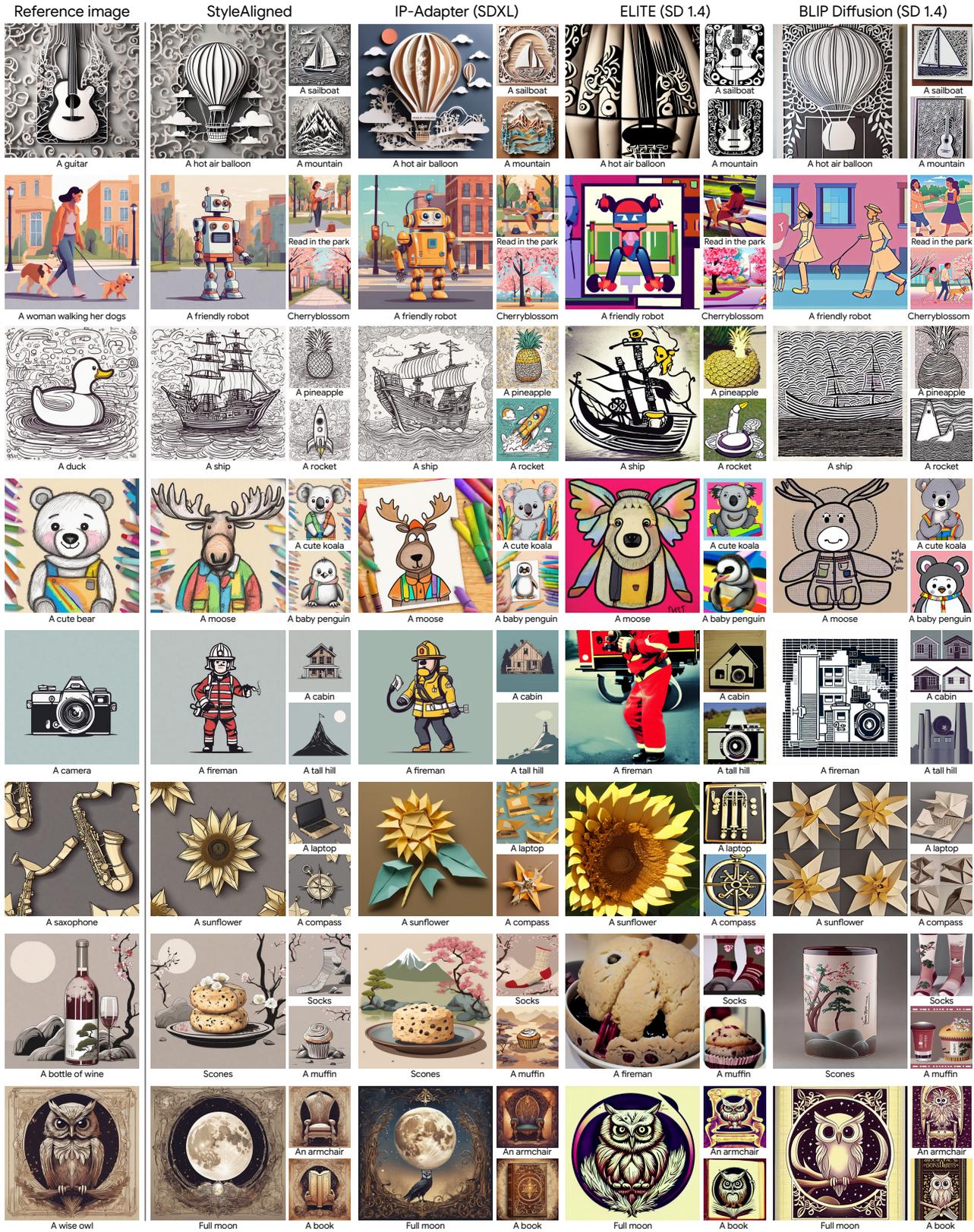


Figure 23. Qualitative comparison to encoders based personalization methods.

In which row below, the images better **share the same style** while **matching the text above**?
Consider **consistency, alignment to the texts, and overall quality** of the images in the row.

- Top row
- Bottom row

Continue

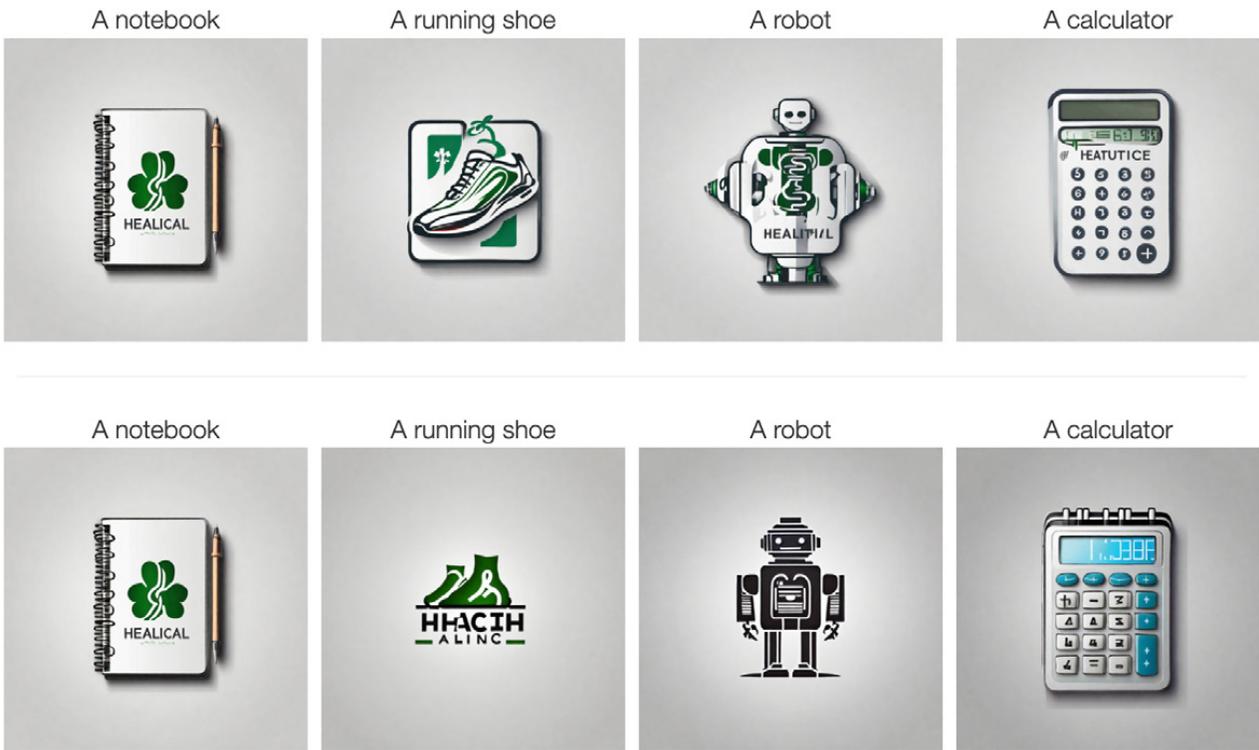


Figure 24. Screenshot from the user study. Each row of images represents the result obtained by different method. The user had to assess which row is better in terms of style alignment and text alignment.

List of prompts for our evaluation set generation:

1. {A house, A temple, A dog, A lion} in sticker style.
2. {Flowers, Golden Gate bridge, A chair, Trees, An airplane} in watercolor painting style.
3. {A village, A building, A child running in the park, A racing car} in line drawing style.
4. {A phone, A knight on a horse, A train passing a village, A tomato in a bowl} in cartoon line drawing style.
5. {Slices of watermelon and clouds in the background, A fox, A bowl with cornflakes, A model of a truck} in 3d rendering style.
6. {A mushroom, An Elf, A dragon, A dwarf} in glowing style.
7. {A thumbs up, A crown, An avocado, A big smiley face} in glowing 3d rendering style.
8. {A bear, A moose, A cute koala, A baby penguin} in kid crayon drawing style.
9. {An orchid, A Viking face with beard, A bird, An elephant} in wooden sculpture.
10. {A portrait of a person wearing a hat, A portrait of a woman with a long hair, A person dancing, A person fishing} in oil painting style.
11. {A woman walking a dog, A friendly robot, A woman reading in the park, Cherryblossom} in flat cartoon illustration style.
12. {A birthday cake, The letter A, An espresso machine, A Car} in abstract rainbow colored flowing smoke wave design.
13. {A flower, A piano, A butterfly, A guitar} in melting golden 3d rendering style.
14. {A train, A car, A bicycle, An airplane} in minimalist round BW logo.
15. {A rocket, An astronaut, A man riding a snowboard, A pair of rings} in neon graffiti style.
16. {A teapot, A teacup, A stack of books, A cozy armchair} in vintage poster style.
17. {A mountain range, A bear, A campfire, A pine forest} in woodblock print style.
18. {A surfboard, A beach shack, A wave, A seagull} in retro surf art style.
19. {A paintbrush, A sunflower field, A scarecrow, A rustic barn} in a minimal origami style.
20. {A cityscape, Hovering vehicles, Dragons, Boats} in cyberpunk art style.
21. {A treasure box, A pirate ship, A parrot, A skull} in tattoo art style.
22. {Music stand, A vintage microphone, A turtle, A saxophone} in art deco style.
23. {A tropical island, A mushroom, A palm tree, A cocktail} in vintage travel poster style.
24. {A carousel, Cotton candy, A ferris wheel, Balloons} in retro amusement park style.
25. {A serene river, A rowboat, A bridge, A willow tree} in 3D render, animation studio style.
26. {A retro guitar, A jukebox, A chess piece, A milkshake} in 1950s diner art style.
27. {A snowy cabin, A sleigh, A snowman, A winter forest} in Scandinavian folk art style.
28. {A bowl with apples, A pencil, A big armor, A magical sunglasses} in fantasy poison book style.
29. {A kiwi fruit, A set of drums, A hammer, A tree} in Hawaiian sunset painting style.
30. {A guitar, A hot air balloon, A sailboat, A mountain} in papercut art style.
31. {A coffee cup, A typewriter, A pair of glasses, A vintage camera} in retro hipster style.
32. {A board of backgammon, A shirt and pants, Shoes, A cocktail} in vintage postcard style.
33. {A roaring lion, A soaring eagle, A dolphin, A galloping horse} in tribal tattoo style.
34. {A pizza, Candles and roses, A bottle, A chef} in Japanese ukiyo-e style.
35. {A wise owl, A full moon, A magical chair, A book of spells} in fantasy book cover style.
36. {A cozy cabin, Snow-covered trees, A warming fireplace, A steaming cup of cocoa} in hygge style.
37. {A bottle of wine, A scone, A muffin, Pair of socks} in Zen garden style.
38. {A diver, Bowl of fruits, An astronaut, A carousel} in celestial artwork style.
39. {A horse, A castle, A cow, An old phone} in medieval fantasy illustration style.
40. {A mysterious forest, Bioluminescent plants, A graveyard, A train station} in enchanted 3D rendering style.
41. {A globe, An airplane, A suitcase, A compass} in travel agency logo style.
42. {A persian cat playing with a ball of wool, A man skiing down the hill, A train at the station, A bear eating honey} in cafe logo style.
43. {A book, A quill pen, An inkwell, An umbrella} in educational institution logo style.
44. {A hat, A strawberry, A screw, A giraffe} in mechanical repair shop logo style.
45. {A notebook, A running shoe, A robot, A calculator} in healthcare and medical clinic logo style.
46. {A rubber duck, A pirate ship, A rocket, A pineapple} in doodle art style.
47. {A trumpet, A fishbowl, A palm tree, A bicycle} in abstract geometric style.
48. {A teapot, A kangaroo, A skyscraper, A lighthouse} in mosaic art style.
49. {A ninja, A hot air balloon, A submarine, A watermelon} in paper collage style.
50. {A saxophone, A sunflower, A compass, A laptop} in origami style.
51. {A penguin, A bicycle, A tornado, A pineapple} in abstract graffiti style.
52. {A magician's hat, A UFO, A roller coaster, A beach ball} in street art style.
53. {A cactus, A shopping cart, A child playing with cubes, A camera} in mixed media art style.
54. {A snowman, A surfboard, A helicopter, A cappuccino} in abstract expressionism style.
55. {A robot, A cupcake, A woman playing basketball, A sunflower} in digital glitch art style.
56. {A treehouse, A disco ball, A sailing boat, A cocktail} in psychedelic art style.
57. {A football helmet, A playmobil, A truck, A watch} in street art graffiti style.
58. {A cabin, A leopard, A squirrel, A rose} in pop art style.
59. {A bus, A drum, A rabbit, A shopping mall} in minimalist surrealism style.
60. {A frisbee, A monkey, A snake, skates} in abstract cubism style.
61. {A piano, A villa, A snowboard, A rubber duck} in abstract impressionism style.
62. {A laptop, A man playing soccer, A woman playing tennis, A rolling chair} in post-modern art style.
63. {A cute puppet, A glass of beer, A violin, A child playing with a kite} in neo-futurism style.
64. {A dog, A brick house, A lollipop, A woman playing on a guitar} in abstract constructivism style.
65. {A kite surfing, A pizza, A child doing homework, A person doing yoga} in fluid art style.
66. {Ice cream, A vintage typewriter, A pair of reading glasses, A handwritten letter} in macro photography style.
67. {A gourmet burger, A sushi, A milkshake, A pizza} in professional food photography style for a menu.
68. {A crystal vase, A pocket watch, A compass, A leather-bound journal} in vintage still life photography style.
69. {A sake set, A stack of books, A cozy blanket, A cup of hot cocoa} in miniature model style.
70. {A retro bicycle, A sunhat, A picnic basket, A kite} in outdoor lifestyle photography style.
71. {A group of hikers on a mountain trail, A winter evening by the fire, A hen, A person enjoying music} in realistic 3D render.
72. {A tent, A person knitting, A rural farm scene, A basket of fresh eggs} in retro music and vinyl photography style

73. {A giraffe, A blanket, A fork and knife, A pile of candies} in cozy winter lifestyle photography style.
74. {A wildflower, A ladybug, An igloo in antarctica, A person running} in bokeh photography style.
75. {A coffee machine, A laptop, A person working, A plant on the desk} in minimal flat design style.
76. {A camera, A fireman, A wooden house, A tall hill} in minimal vector art style.
77. {A person texting, A person scrawling, A cozy chair, A lamp} in minimal pastel colors style.
78. {A smartphone, A book, A dinner table, A glass of wine} in minimal digital art style.
79. {A brush, An artist painting, A girl holding umbrella, a pool table} in minimal abstract illustration style.
80. {A pair of running shoes, A motorcycle, Keys, A fitness machine} in minimal monochromatic style.
81. {A compass rose, A cactus, A zebra, A blizzard} in woodcut print style.
82. {A lantern, A tricycle, A seashell, A swan} in chalk art style.
83. {Magnifying glass, Gorilla, Airplane, Swing} in pixel art style.
84. {Hiking boots, Kangaroo, Ice cream cone, Hammock} in comic book style.
85. {Horseshoe, Vintage typewriter, Snail, Tornado} in vector illustration style.
86. {A lighthouse, A hot air balloon, A cat, A cityscape} in isometric illustration style.
87. {A compass, A violin, A palm tree, A koala} in wireframe 3D style.
88. {Beach umbrella, Rocket ship, Fox, Waterfall} in paper cutout style.
89. {Tree stump, Harp, Chameleon, Canyon} in blueprint style.
90. {Elephant, UFO toy, Flamingo, Lightning bolt} in retro comic book style.
91. {Robot, Temple, Jellyfish, Sofa} in infographic style.
92. {Microscope, Giraffe, Laptop, Rainbow} in geometric shapes style.
93. {Teapot, Dragon toy, Skateboard, Storm cloud} in cartoon line drawing style.
94. {Crystal ball, Carousel horse, Hummingbird, Glacier} in watercolor and ink wash style.
95. {Feather quill, Satellite dish, Deer, Desert scene} in dreamy surreal style.
96. {Map, Saxophone, Mushroom, Dolphin} in steampunk mechanical style.
97. {Anchor, Clock, Globe, Bicycle} in 3D realism style.
98. {Clock, Helicopter, Whale, Starfish} in retro poster style.
99. {Binoculars, Bus, Pillow, Cloud} in bohemian hand-drawn style.
100. {Rhino, Telescope, Stool, Panda} in vintage stamp style.