# ZipLoRA: Any Subject in Any Style by Effectively Merging LoRAs

Viraj Shah[*,1,2]    Nataniel Ruiz[1]    Forrester Cole[1]    Erika Lu[1]
Svetlana Lazebnik[2]    Yuanzhen Li[1]    Varun Jampani[†,1]
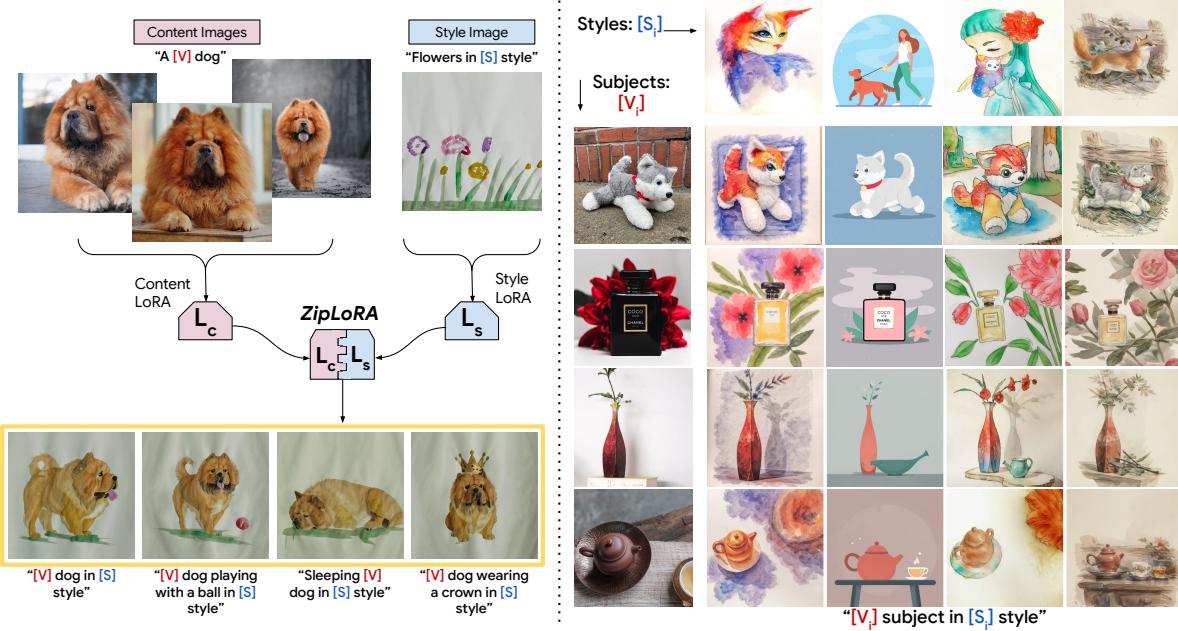[1] Google Research    [2] UIUC

Figure 1. By effectively merging independently trained style and content LoRAs, our proposed method **ZipLoRA** is able to generate *any user-provided subject in any user-provided style*, providing unprecedented control over personalized creations using diffusion models.

## Abstract

*Methods for finetuning generative models for concept-driven personalization generally achieve strong results for subject-driven or style-driven generation. Recently, low-rank adaptations (LoRA) have been proposed as a parameter-efficient way of achieving concept-driven personalization. While recent work explores the combination of separate LoRAs to achieve joint generation of learned styles and subjects, existing techniques do not reliably address the problem; they often compromise either subject fidelity or style fidelity. We propose **ZipLoRA**, a method to cheaply and effectively merge independently trained style and subject LoRAs in order to achieve generation of **any user-provided subject in any user-provided style**. Experiments on a wide range of subject and style combinations show that ZipLoRA can generate compelling results with meaningful improvements over baselines in subject and style fidelity while preserving the ability to recontextualize.*

## 1. Introduction

Recently, diffusion models [13, 28, 34] have allowed for impressive image generation quality with their excellent understanding of diverse artistic concepts and enhanced controllability due to multi-modal conditioning support (with text being the most popular mode). The usability and flexibility of generative models has further progressed with a wide variety of personalization approaches, such as DreamBooth [29] and StyleDrop [33]. These approaches fine-tune a base diffusion model on the images of a specific concept to produce novel renditions in various contexts. Such concepts can be a specific object or person, or an artistic style.

While personalization methods have been used for subjects and styles independently, a key unsolved problem is to generate a specific user-provided subject in a specific user-

Project page: https://ziplora.github.io

1

provided style. For example, an artist may wish to render a specific person in their personal style, learned through examples of their own work. A user may wish to generate images of their child's favorite plush toy, in the style of the child's watercolor paintings. Moreover, if this is achieved two problems are simultaneously solved: (1) the task of representing any given subject in any style, and (2) the problem of controlling diffusion models through images rather than text, which can be imprecise and unsuitable for certain generation tasks. Finally, we can imagine a large-scale application of such a tool, where a bank of independently learned styles and subjects are shared and stored online. The task of arbitrarily rendering *any subject in any style* is an open research problem that we seek to address.

A pitfall of recent personalization methods is that many finetune all of the parameters of a large base model, which can be costly. Parameter Efficient Fine-Tuning (PEFT) approaches allow for fine-tuning models for concept-driven personalization with much lower memory and storage budgets. Among the various PEFT approaches, Low Rank Adaptation (LoRA) [14] has emerged as a favored method for researchers and practitioners alike due to its versatility. LoRA learns low-rank factorized weight matrices for the attention layers (these learned weights are themselves commonly referred to as "LoRAs"). By combining LoRA and algorithms such as DreamBooth [29], the learned subject-specific LoRA weights enable the model to generate the subject with semantic variations.

With the growing popularity of LoRA personalization, there have been attempts to merge LoRA weights, specifically by performing a linear combination of subject and style LoRAs, with variable coefficients [30]. This allows for a control over the "strength" of each LoRA, and users sometimes are able, through careful grid search and subjective human evaluation, to find a combination that allows for accurate portrayal of the subject under the specific style. This method lacks robustness across style and subject combinations, and is also incredibly time consuming.

In our work, we propose *ZipLoRA*, a simple yet effective method to generate any subject in any style by cheaply merging independently trained LoRAs for subject and style. Our approach works consistently on a wide variety of subject and style LoRAs without enforcing any restriction on the way these are trained. This allows users and artists to easily combine publicly available subject and style LoRAs of their choice. ZipLoRA is hyperparameter-free, i.e. it does not require manual tuning of any hyperparameters or merger weights.

Our approach takes advantage of the recently released Stable Diffusion XL (SDXL) model [27] and is based on three important observations. **(1)** SDXL exhibits strong style learning properties, comparable to results shown by StyleDrop [33] on Muse [3]. Specifically, unlike previ-

ous versions of Stable Diffusion, SDXL is able to learn styles using just a single exemplar image by following a DreamBooth protocol [29] without any human feedback. **(2)** LoRA weights for different layers $\Delta W_i$ (where $i$ denotes the layer) are sparse. *i.e.*, most of the elements in $\Delta W_i$ have very small magnitude, and have little effect on generation quality and fidelity. **(3)** Columns of the weight matrices of two independently trained LoRAs may have varying levels of "alignment" between each other, as measured by cosine similarity, for example. We find that directly summing columns that are highly aligned degrades performance of the merged model.

Based on these observations, we hypothesize that a method that operates akin to a zipper, aiming to reduce the quantity of similar-direction sums while preserving the content and style generation properties of the original LoRAs will yield more robust, higher-quality merges. Much like a zipper seamlessly joins two sides of a fabric, our proposed optimization-based approach finds a disjoint set of merger coefficients for blending the two LoRAs. This ensures that the merged LoRA adeptly captures both subject and style. Our optimization process is lightweight and significantly improves the merging performance on challenging content-style combinations, where the two LoRAs are highly aligned.

We summarize our contributions as follows:

- We demonstrate some key observations about current text-to-image diffusion models and personalization methods, particularly in relation to style personalization. We further examine the sparsity of concept-personalized LoRA weight matrix coefficients and the prevalence and deleterious effect of highly aligned columns for LoRA matrices.

- Using these insights we propose **ZipLoRA**, a simple optimization method that allows for effective merging of independently trained style and subject LoRAs to allow for the generation of *any subject in any style.* ZipLoRA is a first exploration into the world of techniques that merge LoRAs to achieve new generation capabilities.

- We demonstrate the effectiveness of our approach on a variety of image stylization tasks, including content-style transfer and recontextualization. We also demonstrate that ZipLoRA outperforms existing methods of merging LoRAs as well as other baseline approaches.

## 2. Related Work

**Fine-tuning of Diffusion Models for Custom Generation.** In the evolving field of text-to-image (T2I) model personalization, recent studies have introduced various methods to fine-tune large-scale T2I diffusion models for depicting specific subjects based on textual descriptions. Techniques like Textual Inversion [8] focus on learning text embeddings, while DreamBooth [29] fine-tunes the entire T2I

model for better subject representation. Later methods aim to optimize specific parts of the networks [11, 20]. Additionally, techniques like LoRA [14] and StyleDrop [33] concentrate on optimizing low-rank approximations and a small subset of weights, respectively, for style personalization. DreamArtist [5] introduces a novel one-shot personalization method using a positive-negative prompt tuning strategy. While these fine-tuning approaches yield high-quality results, they typically are limited to learning only one concept (either subject or style). One exception is Custom Diffusion [20], which attempts to learn multiple concepts simultaneously. However, Custom Diffusion requires expensive joint training from scratch and still yields inferior results when used for stylization as it fails to disentangle the style from the subject.

**Combining LoRAs.** Combining different LoRAs remain under-explored in the literature particularly from the point of view of fusing style and the subject concepts. Ryu [30] shows a method to combine independently trained LoRAs by weighed arithmetic summation. In [10], authors discuss fusing multiple concept LoRAs, however, it is an expensive method that requires retraining as it does not merge LoRAs but rather re-trains the entire model. A concurrent work discusses a strategy to obtain Mixture of Experts by combining multiple LoRAs using a gating function [1].

**Image Stylization.** Image-based style transfer is an area of research dating back at least 20 years [6, 12]. Great advances in arbitrary style transfer was achieved by the convolutional neural network-based approaches [9, 15, 16, 22, 26]. Generative models such as GANs [17–19] can also be used as a prior for image stylization tasks [2, 24, 35]. Many recent GAN-based approaches achieve successful one-shot stylizations [4, 7, 21, 23, 25, 32, 36–39] by fine-tuning a pre-trained GAN for a given reference style. However, these methods are limited to images from only a single domain (such as faces). Further, most existing GANs do not provide any direct, text-based control over the semantics of the output, thus they cannot produce the reference subject in novel contexts. Compared to older generative models, diffusion models [13, 28, 34] offer superior generation quality and text-based control; however, to date, it has been difficult to use them for one-shot stylization driven by image examples. Ours is one of the first works demonstrating the use of diffusion models for high-quality example-based stylization combined with an ability to re-contextualize to diverse scenarios.

# 3. Methods

## 3.1. Background

**Diffusion Models** [13, 28, 34] are state-of-the-art generative models known for their high-quality, photorealistic image synthesis. Their training comprises two phases: a for-

ward process, where an image transitions into a Gaussian noise through incremental Gaussian noise addition, and a reverse process, reconstructing the original data from the noise. The reverse process is typically learnt using an U-net with text conditioning support enabling text-to-image generation at the time of inference. In our work, we focus on widely used latent diffusion model [28] which learns the diffusion process in the latent space instead of image space. In particular, we use Stable Diffusion XL v1 [27] for all our experiments.

**LoRA Fine-tuning.** LoRA (Low-Rank Adaptation) is a method for efficient adaptation of Large Language and Vision Models to a new downstream task [14, 30]. The key concept of LoRA is that the weight updates $\Delta W$ to the base model weights $W_0 \in \mathbb{R}^{m \times n}$ during fine-tuning have a "low intrinsic rank," thus the update $\Delta W$ can be decomposed into two low-rank matrices $B \in \mathbb{R}^{m \times r}$ and $A \in \mathbb{R}^{r \times n}$ for efficient parameterization with $\Delta W = BA$. Here, $r$ represents the intrinsic rank of $\Delta W$ with $r << min(m, n)$. During training, only $A$ and $B$ are updated to find suitable $\Delta W = BA$, while keeping $W_0$ constant. For inference, the updated weight matrix $W$ can be obtained as $W = W_0 + BA$. Due to its efficiency, LoRA is widely used for fine-tuning open-sourced diffusion models.

**Problem Setup.** In this work, we aim to produce accurate renditions of a custom object in a given reference style by merging LoRA weights obtained by separately fine-tuning a given text-to-image diffusion model on a few reference images of the object/style.

We start with a base diffusion model represented as $D$ with pre-trained weights $W_0^{(i)}$ with $i$ as layer index. One can adapt the base model $D$ to any given concept by simply adding the corresponding set of LoRA weights $L_x\{\Delta W_x^{(i)}\}$ to the model weights. We represent it as: $D_{L_x} = D \oplus L_x = W_0 + \Delta W_x$. We drop the superscript $(i)$ for simplicity since our operations are applied over all the LoRA-enabled weight matrices of our base model $D$.

We are given two independently trained set of LoRAs $L_c = \{\Delta W_c^{(i)}\}$ and $L_s = \{\Delta W_s^{(i)}\}$ for our base model $D$, and we aim to find a merged LoRA $L_m = \{\Delta W_m^{(i)}\} = \mathrm{Merge}(L_c, L_s)$ that can combine the effects of both the individual LoRAs in order to stylize the given object in a desired reference style.

**Direct Merge.** LoRA is popularly used as a plug-and-play module on top of the base model, thus a most common way to combine multiple LoRAs is a simple linear combination:

$$L_m = L_c + L_s \implies \Delta W_m = w_c \cdot \Delta W_c + w_s \cdot \Delta W_s,$$
(1)

where $w_c$ and $w_s$ are coefficients of content and style LoRAs, respectively, which allow for a control over the "strength" of each LoRA. For a given subject and style
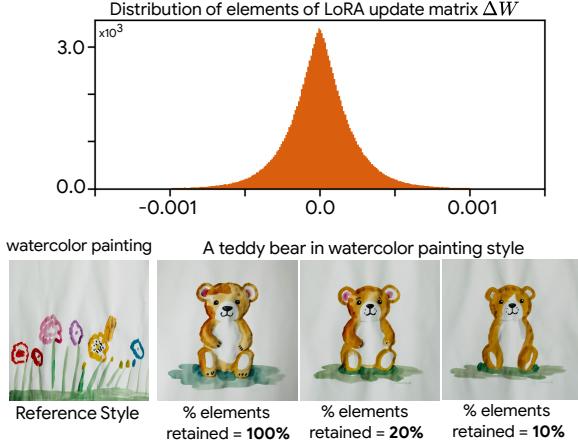
Figure 2. **LoRA weight matrices are sparse.** Most of the elements in $\Delta W$ have a magnitude very close to zero, and can be conveniently thrown away without affecting the generation quality of the fine-tuned model.

LoRA, one may be able to find a particular combination of $w_c$ and $w_s$ that allows for accurate stylization through careful grid search and subjective human evaluation, but this method is not robust and very time consuming. To this end, we propose a hyperparameter-free approach that does not require this onerous process.

### 3.2. ZipLoRA

Our approach builds on two interesting insights:

**(1) LoRA update matrices are sparse.** We observe that the update matrices $\Delta W$ for different LoRA layers are sparse, *i.e.*, most of the elements in $\Delta W$ have a magnitude very close to zero, and thus have little impact on the output of the fine-tuned model. For each layer, we can sort all the elements by their magnitude and zero out the lowest up to a certain percentile. We depict the distribution of elements of $\Delta W_i^{m \times n}$ in Fig. 2, along with samples generated after zeroing out 80% and 90% of the lowest-magnitude elements of weight update matrix $\Delta W$ for all the layers. As can be seen, the model performance is unaffected even when 90% of the elements are thrown away. This observation follows from the fact that the rank of $\Delta W$ is very small by design, thus the information contained in most columns of $\Delta W$ is redundant.

**(2) Highly aligned LoRA weights merge poorly.** Columns of the weight matrices of two independently trained LoRAs may contain information that is not disentangled, *i.e.*, the cosine similarity between them can be non-zero. We observe that the extent of alignment between the columns of LoRA weights plays a significant role in determining the quality of resulting merge: if we directly add the columns with non-zero cosine similarity to each other, it leads to superimposition of their information about the individual con-
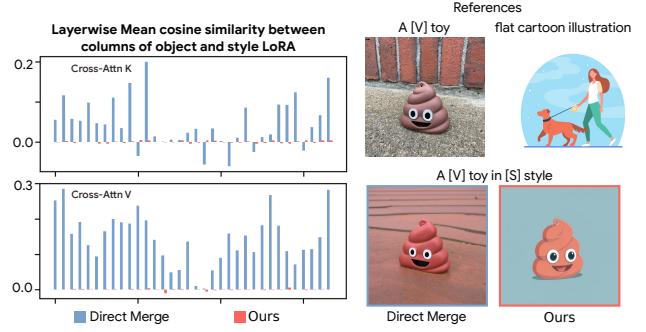


Figure 3. **Highly aligned LoRA weights merge poorly.** When LoRA weight columns are highly aligned, a direct merge obtains subpar results. Instead, our approach minimizes the mean cosine similarity between the columns of the LoRA updates across the layers.

cepts, resulting in the loss of the ability of the merged model to synthesize input concepts accurately. We further observe that such loss of information is avoided when the columns are orthogonal to each other with cosine similarity equal to zero.

Note that each weight matrix represents a linear transformation defined by its columns, so it is intuitive that the merger would retain the information available in these columns only when the columns that are being added are orthogonal to each other. For most content-style LoRA pairs the cosine similarities are non-zero, resulting in signal interference when they are added directly. In Fig. 3 we show the mean cosine similarity values for each layer of the last U-net block for a particular content-style pair before and after applying ZipLoRA. One can see high non-zero cosine similarity values for the direct merge which results in poor stylization quality. On the other hand, ZipLoRA reduces the similarity values significantly to achieve a superior result.

To prevent signal interference during the merger, we multiply each column with a learnable coefficient such that the orthogonality between the columns can be achieved. The fact that LoRA updates are sparse allows us to neglect certain columns from each LoRA, thus facilitating the task of minimizing interference. As shown in Fig. 4, we introduce a set of merger coefficient vectors $m_c$ and $m_s$ for each LoRA layer of the content and style LoRAs, respectively:

$$L_m = \text{Merge}(L_c, L_s, m_c, m_s)$$
$$\implies \Delta W_m = m_c \otimes \Delta W_c + m_s \otimes W_s, \quad (2)$$

where $\otimes$ represents element-wise multiplication between $\Delta W$ and broadcasted merger coefficient vector $m$ such that $j^{th}$ column of $\Delta W$ gets multiplied with $j^{th}$ element of $m$. The dimensionalities of $m_c$ and $m_s$ are equal to the number of columns in corresponding $\Delta W$, thus each element of the merger coefficient vector represents the contribution of the
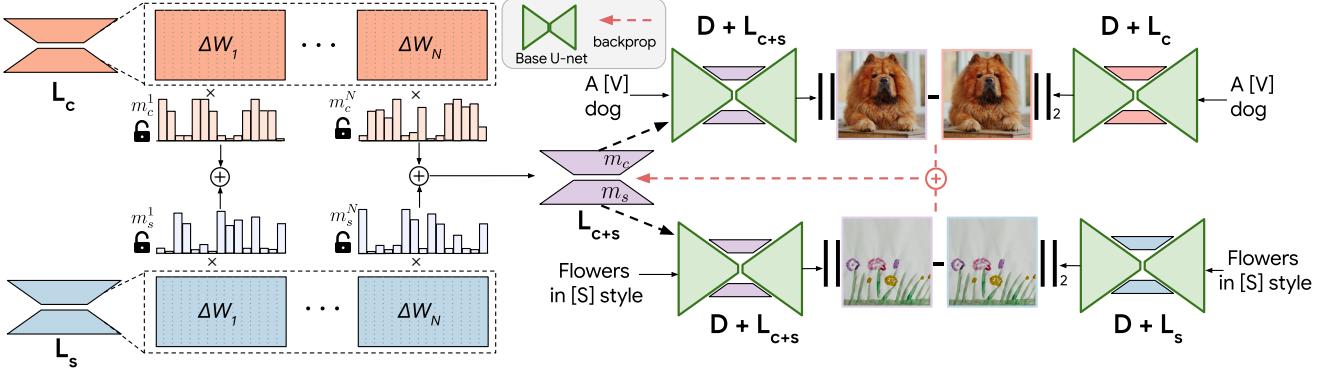
Figure 4. **Overview of ZipLoRA**. Our method learns mixing coefficients for each column of $\Delta W_i$ for both style and subject LoRAs. It does so by (1) minimizing the difference between subject/style images generated by the mixed LoRA and original subject/style LoRA models, while (2) minimizing the cosine similarity between the columns of content and style LoRAs. In essence, the zipped LoRA tries to conserve the subject and style properties of each individual LoRA, while minimizing signal interference of both LoRAs.

corresponding column of the LoRA matrix $\Delta W$ to the final merge.

Our ZipLoRA approach has two goals: (1) to minimize the interference between content and style LoRAs, defined by the cosine similarity between the columns of content and style LoRAs while (2) conserving the capability of the merged LoRA to generate the reference subject and style independently by minimizing the difference between subject/style images generated by the mixed LoRA and original subject/style LoRAs. To ensure that the columns that are merged with each other minimize signal interference, our proposed loss seeks to minimize the cosine similarity between the merge vectors $m_c$ and $m_s$ of each layer. Meanwhile, we wish to ensure that the original behavior of both the style and the content LoRAs is preserved in the merged model. Therefore, as depicted in Fig. 4, we formulate an optimization problem with following loss function:

$$
\begin{aligned}
\mathcal{L}_{merge} =& \|(D \oplus L_m)(x_c, p_c) - (D \oplus L_c)(x_c, p_c)\|_2 \\
& + \|(D \oplus L_m)(x_s, p_s) - (D \oplus L_s)(x_s, p_s)\|_2 \\
& + \lambda \sum_i |m_c^{(i)} \cdot m_s^{(i)}|,
\end{aligned} \tag{3}
$$

where the merged model $L_m$ is calculated using $m_c$ and $m_s$ as per Eq. 2; $p_c, p_s$ are text conditioning prompts for content and style references respectively, and $\lambda$ is an appropriate multiplier for the cosine-similarity loss term. Note that the first two terms ensure that the merged model retains the ability to generate individual style and content, while the third term enforces an orthogonality constraint between the columns of the individual LoRA weights. Importantly, we keep the weights of the base model and the individual LoRAs frozen, and update only the merger coefficient vectors. As seen in the next section, such a simple optimization method is effective in producing strong stylization of custom subjects. Further, ZipLoRA requires only 100 gradient

updates which is $10\times$ faster compared to joint training approaches.

## 4. Experiments

**Datasets.** We choose a diverse set of content images from the DreamBooth dataset [29], which provides 30 image sets each containing 4-5 images of a given subject. Similarly, a diverse set of style reference images is selected from the data provided by authors of StyleDrop [33]. We use only a single image for each style. The attribution and licence information for all the content and style images used are available in the DreamBooth and StyleDrop manuscripts/websites.

**Experimental Setup.** We perform all our experiments using the SDXL v1.0 [27] base model. We use DreamBooth fine-tuning with LoRA of rank 64 for obtaining all the style and content LoRAs. We update the LoRA weights using Adam optimizer for 1000 steps with batch size of 1 and learning rate of 0.00005. We keep the text encoders of SDXL frozen during the LoRA fine-tuning. For ZipLoRA, we use $\lambda = 0.01$ in Eq. 3 for all our experiments, and run the optimization until cosine similarity drops to zero with a maximum number of gradient updates set to 100.

### 4.1. Style-tuning behavior of SDXL model

As discussed in Sec. 3, we observe, surprisingly, that a pretrained SDXL model exhibits strong style learning when fine-tuned on only one reference style image. We show style-tuning results on SDXL model in Fig. 5. For each reference image, we apply LoRA fine-tuning of SDXL model using DreamBooth objective with LoRA rank= 64. For fine-tuning, we follow a similar prompt formation as provided in StyleDrop: "an <object> in the <style description> style". Once fine-tuned, SDXL is able to represent diverse set of concepts in the reference style by cap-
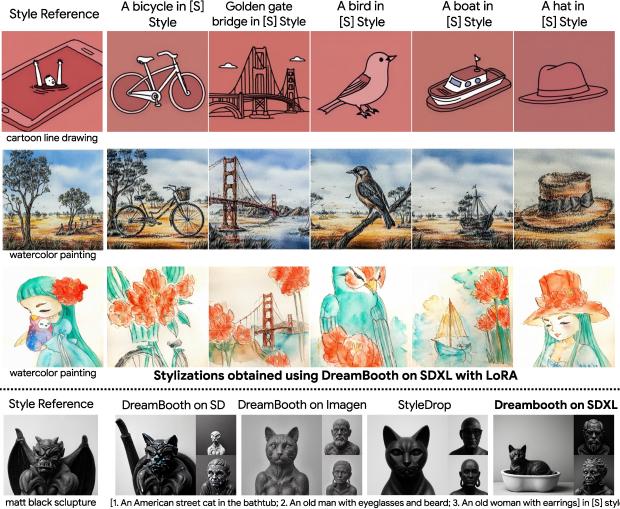
Figure 5. **Style Learning using DreamBooth on SDXL.** Top: SDXL model learns to produce stylized outputs when fine-tuned on a single example of a reference style using LoRA with a Dream-Booth objective. Bottom: The stylizations produced by fine-tuned SDXL model are highly competent, compared to those of other models. Note that unlike StyleDrop, SDXL DreamBooth fine-tuning does not require human feedback.

turing the nuances of painting style, lighting, colors, and geometry accurately. The question of why this model exhibits this strong style learning performance, as opposed to the lesser performance of previous SD versions [28] (or Imagen [31]) is left open and can have many answers including training data, model architecture and training schemes.

We also provide comparisons with StyleDrop on Muse [3], DreamBooth on Imagen, and DreamBooth on Stable Diffusion in Fig. 5. We observe that SDXL style-tuning performs significantly better than the competing methods. Note that StyleDrop requires iterative training with human feedback whereas SDXL style-tuning does not. This behavior of SDXL makes it the perfect candidate for investigating the merging of style LoRAs with subject LoRAs to achieve personalized stylizations. Thus, we choose to use it as a base model for all of our experiments.

### 4.2. Personalized Stylizations

To start with, we obtain the style LoRAs following the style-tuning on SDXL as described in Sec. 4.1, and obtain object LoRAs by applying DreamBooth fine-tuning on the subject references. Fig. 1 and Fig. 6 show the results of our approach for combining various style and content LoRAs. Our method succeeds at both preserving the identity of the reference subject and capturing the unique characteristics of the reference style.

We also present qualitative comparisons with other approaches in Fig. 6. As a baseline, we compare with the

direct arithmetic merge obtained through Eq. 1 with $w_c$ and $w_s$ set to 1. Such direct addition results in loss of information captured in each LoRA and produces inferior results with distorted object and/or style.

We additionally compare our method with joint training of subject and style using a multi-subject variant of Dream-Booth with multiple rare unique identifiers. As shown, joint training fails to learn the disentanglement between object and style and produces poor results. It also is the least flexible method since it does not allow the use of pre-trained LoRAs, neither can it be used as a style-only or content-only LoRA. Further, it requires $10\times$ as many training steps as ZipLoRA.

StyleDrop [33] proposes a DreamBooth+StyleDrop approach for achieving personalized stylizations, where a StyleDrop method is applied on a DreamBooth model fine-tuned on the reference object. Our comparisons show that its performance is not ideal, considering the high compute cost and human feedback requirements. It also requires adjusting the object and style model weights $w_c$ and $w_s$ similar to the direct merge in order to produce reasonable outputs, while our method is free from any such hyperparameter tuning.

**Quantitative results.** We conduct user studies for a quantitative comparison of our method with existing approaches. In our study, each participant is shown a reference subject and a reference style along with outputs of two methods being compared, in a random order, and asked which output best depicts the reference style while preserving the reference subject fidelity. We conducted separate user studies for ZipLoRA vs. each of the three competing approaches, and received 360 responses across 45 users for each case. We show the results in Tab. 1. As we can see, ZipLoRA re-

Table 1. **User Preference Study**. We compare the user preference of accurate stylization and subject fidelity between our approach and competing methods. Users generally prefer our approach.

| % **Preference for ZipLoRA over:** | | |
|---|---|---|
| Direct Merge | Joint Training | StyleDrop |
| 82.7% | 71.1% | 68.0% |

Table 2. **Image-alignment and Text-alignment Scores.** We compare cosine similarities between CLIP (for style and text) and DINO features (for subject) of the output and reference style, subject, and prompt respectively. ZipLoRA provides superior subject and text fidelity while also maintaining the style-alignment.

| | ZipLoRA | Joint Training | Direct Merge |
|---|---|---|---|
| Style-alignment ↑ | 0.699 | 0.680 | 0.702 |
| Subject-alignment ↑ | 0.420 | 0.378 | 0.357 |
| Text-alignment ↑ | 0.303 | 0.296 | 0.275 |

Figure 6. **Qualitative Comparison:** We compare samples from our method (Ours), versus direct arithmetic merge, joint training and StyleDrop [33]. We observe that our method achieves strong style and subject fidelity that surpasses competing methods.

ceives higher user preference in all three cases owing to its high-quality stylization while preserving subject integrity.

Following DreamBooth [29], we also provide comparisons using image-alignment and text-alignment scores in Tab. 2. We employ three metrics: for style-alignment, we use CLIP-I scores of image embeddings of output and the style reference; for subject-alignment, we employ DINO features for the output and the reference subject; and for text-alignment, we use CLIP-T embeddings of the output and the text prompt. In all three cases, we use cosine-similarity as the metric and calculate averages over 4 subjects in 8 styles each. ZipLoRA results in competitive style-alignment scores as compared to joint training and direct merge, while achieving significantly better scores for subject-alignment. This highlights ZipLoRA's superiority in maintaining the subject fidelity. ZipLoRA also outper-

forms the other two in text-alignment, implying that it preserves the text-to-image generation capability, and also expresses the designated style and subject better (since these are also part of the text prompt). One should note that these metrics are not perfect, particularly when it comes to measuring style alignment, since they lack the ability to capture subtle stylistic details, and are entangled with semantic properties of images, such as the overall content.

**Ability to re-contextualize.** The merged ZipLoRA model can recontextualize reference objects in diverse contexts and with semantic modifications while maintaining stylization quality. As shown in Fig. 7, our method preserves the base model's text-to-image generation capabilities while accurately stylizing the entire image in the reference style. Such ability is highly valuable in various artistic use cases that requires controlling contexts, subject identities, and
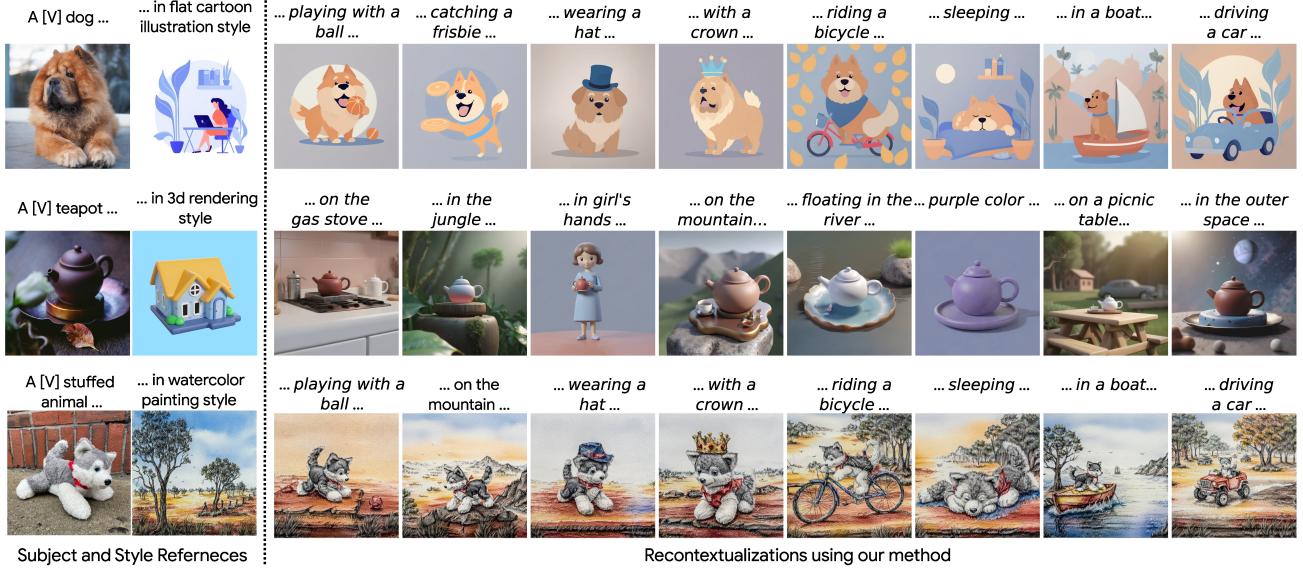
Figure 7. Our method successfully re-contextualizes the reference subject while preserving the stylization in the given style.
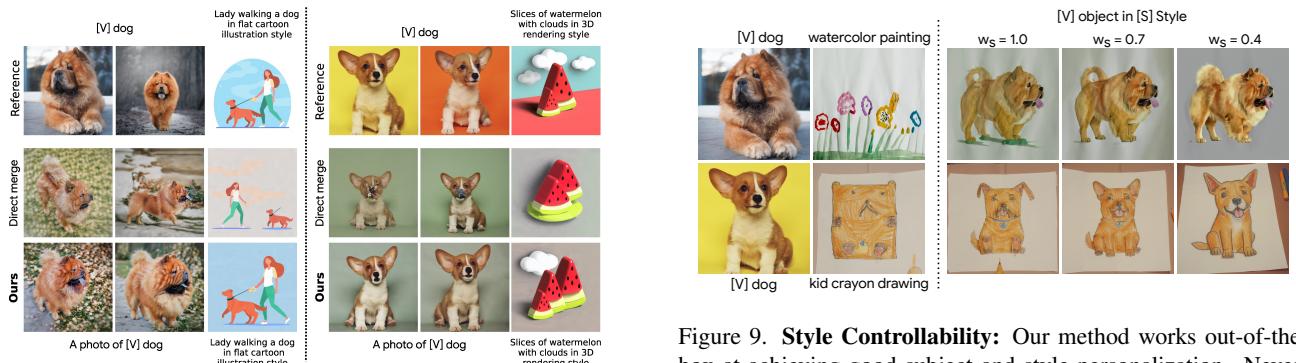


Figure 8. Our method does not lose the ability to generate individual concepts, unlike the direct merge approach.
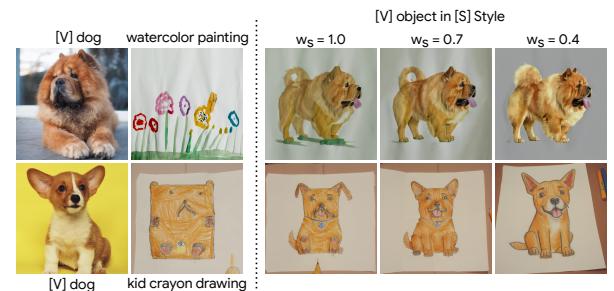


Figure 9. **Style Controllability:** Our method works out-of-the-box at achieving good subject and style personalization. Nevertheless, varying the merging weights $w_s$ allows for controlling the extent of stylization.

styles.

**Controlling the extent of stylization.** Our optimization-based method directly provides a scalar weight value for each column of the LoRA update, thus eliminating a need for tuning and adjustments for obtaining reasonable results. However, we can still allow the strength of object and style content to be varied for added controllability. One can attenuate the style layer weights by multiplying them with an additional scalar multiplier $w_s$ to limit the contribution of the style in the final output. As shown in Fig. 9, this allows for a smooth control over the extent of stylization as $w_s$ varies between 0 to 1.

**Ability to produce the reference object and the style.** Apart from producing accurate stylizations, an ideal LoRA merge should also preserve the ability to generate individual object and style correctly. This way, a merged LoRA model can also be used as a replacement of both the indi-

vidual LoRAs, or as a Mixture-of-Expert model. As shown in Fig. 8, our approach retains the original behavior of both the models and can accurately generate specific structural and stylistic elements of each constituent LoRA, while direct merge fails.

## 5. Conclusion

In this paper, we have introduced **ZipLoRA**, a novel method for seamlessly merging independently trained style and subject LoRAs. Our approach unlocks the ability to generate **any subject in any style** using sufficiently powerful diffusion models like SDXL. By leveraging key insights about pre-trained LoRA weights, we surpass existing methods for this task. ZipLoRA offers a streamlined, cheap, and hyperparameter-free solution for simultaneous subject and style personalization, unlocking a new level of creative controllability for diffusion models.

# References

[1] Anonymous. MoLE: Mixture of loRA experts. In *Submitted to The Twelfth International Conference on Learning Representations*, 2023. under review. 3

[2] Ashish Bora, Ajil Jalal, Eric Price, and Alexandros G. Dimakis. Compressed sensing using generative models. In *ICML*, 2017. 3

[3] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023. 2, 6

[4] Min Jin Chong and David A. Forsyth. Jojogan: One shot face stylization. *CoRR*, abs/2112.11641, 2021. 3

[5] Ziyi Dong, Pengxu Wei, and Liang Lin. Dreamartist: Towards controllable one-shot text-to-image generation via positive-negative prompt-tuning, 2023. 3

[6] Alexei A Efros and William T Freeman. Image quilting for texture synthesis and transfer. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 341–346, 2001. 3

[7] Rinon Gal, Or Patashnik, Haggai Maron, Gal Chechik, and Daniel Cohen-Or. Stylegan-nada: Clip-guided domain adaptation of image generators. *ArXiv*, abs/2108.00946, 2021. 3

[8] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion, 2022. 2

[9] 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. 3

[10] Yuchao Gu, Xintao Wang, Jay Zhangjie Wu, Yujun Shi, Chen Yunpeng, Zihan Fan, Wuyou Xiao, Rui Zhao, Shuning Chang, Weijia Wu, Yixiao Ge, Shan Ying, and Mike Zheng Shou. Mix-of-show: Decentralized low-rank adaptation for multi-concept customization of diffusion models. *arXiv preprint arXiv:2305.18292*, 2023. 3

[11] Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris Metaxas, and Feng Yang. Svdiff: Compact parameter space for diffusion fine-tuning. *arXiv preprint arXiv:2303.11305*, 2023. 3

[12] Aaron Hertzmann, Charles E. Jacobs, Nuria Oliver, Brian Curless, and D. Salesin. Image analogies. *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, 2001. 3

[13] Jonathan Ho, Ajay Jain, and P. Abbeel. Denoising diffusion probabilistic models. *ArXiv*, abs/2006.11239, 2020. 1, 3

[14] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. 2, 3

[15] 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. 3

[16] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, 2016. 3

[17] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4396–4405, 2019. 3

[18] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. *ArXiv*, abs/2006.06676, 2020.

[19] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8107–8116, 2020. 3

[20] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. In *CVPR*, 2023. 3

[21] Gihyun Kwon and Jong-Chul Ye. One-shot adaptation of gan in just one clip. *ArXiv*, abs/2203.09301, 2022. 3

[22] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Universal style transfer via feature transforms. In *Advances in Neural Information Processing Systems*, 2017. 3

[23] Mingcong Liu, Qiang Li, Zekui Qin, Guoxin Zhang, Pengfei Wan, and Wen Zheng. Blendgan: Implicitly gan blending for arbitrary stylized face generation. *ArXiv*, abs/2110.11728, 2021. 3

[24] Sachit Menon, Alexandru Damian, Shijia Hu, Nikhil Ravi, and Cynthia Rudin. Pulse: Self-supervised photo upsampling via latent space exploration of generative models. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2434–2442, 2020. 3

[25] Utkarsh Ojha, Yijun Li, Jingwan Lu, Alexei A. Efros, Yong Jae Lee, Eli Shechtman, and Richard Zhang. Few-shot image generation via cross-domain correspondence. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10738–10747, 2021. 3

[26] Dae Young Park and Kwang Hee Lee. Arbitrary style transfer with style-attentional networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5880–5888, 2019. 3

[27] 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

[28] 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. 1, 3, 6

[29] 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. 1, 2, 5, 7

[30] Simo Ryu. Merging loras. `https://github.com/cloneofsimo/lora`. 2, 3

[31] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, Seyedeh Sara Mahdavi, Raphael Gontijo Lopes, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. *ArXiv*, abs/2205.11487, 2022. 6

[32] Viraj Shah, Ayush Sarkar, Sudharsan Krishnakumar Anita, and Svetlana Lazebnik. Multistylegan: Multiple one-shot image stylizations using a single gan. *arXiv*, 2023. 3

[33] 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, 3, 5, 6, 7

[34] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *ArXiv*, abs/2010.02502, 2020. 1, 3

[35] Xintao Wang, Yu Li, Honglun Zhang, and Ying Shan. Towards real-world blind face restoration with generative facial prior. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9164–9174, 2021. 3

[36] Yue Wang, Ran Yi, Ying Tai, Chengjie Wang, and Lizhuang Ma. Ctlgan: Few-shot artistic portraits generation with contrastive transfer learning. *ArXiv*, abs/2203.08612, 2022. 3

[37] Ceyuan Yang, Yujun Shen, Zhiyi Zhang, Yinghao Xu, Jiapeng Zhu, Zhirong Wu, and Bolei Zhou. One-shot generative domain adaptation, 2021.

[38] Shuai Yang, Liming Jiang, Ziwei Liu, and Chen Change Loy. Pastiche master: Exemplar-based high-resolution portrait style transfer, 2022.

[39] Peihao Zhu, Rameen Abdal, John Femiani, and Peter Wonka. Mind the gap: Domain gap control for single shot domain adaptation for generative adversarial networks. In *International Conference on Learning Representations*, 2022. 3