

# CLASSDIFFUSION: MORE ALIGNED PERSONALIZATION TUNING WITH EXPLICIT CLASS GUIDANCE

Anonymous authors

Paper under double-blind review

## ABSTRACT

Recent text-to-image customization works have proven successful in generating images of given concepts by fine-tuning the diffusion models on a few examples. However, tuning-based methods inherently tend to overfit the concepts, resulting in failure to create the concept under multiple conditions (*e.g.*, headphone is missing when generating “a `<sk>` dog wearing a headphone”). Interestingly, we notice that the base model before fine-tuning exhibits the capability to compose the base concept with other elements (*e.g.*, “a dog wearing a headphone”), implying that the compositional ability only disappears after personalization tuning. We observe a semantic shift in the customized concept after fine-tuning, indicating that the personalized concept is not aligned with the original concept, and further show through theoretical analyses that this semantic shift leads to increased difficulty in sampling the joint conditional probability distribution, resulting in the loss of the compositional ability. Inspired by this finding, we present **ClassDiffusion**, a technique that leverages a **semantic preservation loss** to explicitly regulate the concept space when learning the new concept. Although simple, this approach effectively prevents semantic drift during the fine-tuning process on the target concepts. Extensive qualitative and quantitative experiments demonstrate that the use of semantic preservation loss effectively improves the compositional abilities of fine-tuning models. Lastly, we also extend our ClassDiffusion to personalized video generation, demonstrating its flexibility.

## 1 INTRODUCTION

Thanks to the rapid progress in the diffusion model [31, 46, 53, 57, 61, 63, 65, 66, 70, 91], the field of text-to-image generation has achieved significant progress in recent years. The leading text-to-image models [1, 20, 33, 37, 39, 64, 78] have been successful in generating high-fidelity images that align well with textual inputs. Recently, a significant part of the research [2, 4, 5, 8, 11, 28, 35, 58, 77, 86, 89, 90] has changed their focus from creating high-quality images to improving control over the generated images. Among these works, an important and widely explored research domain is subject-driven personalized generation, which aims to generate new images for a specific concept given some reference images of that concept.

Existing personalization methods [1, 9, 20, 21, 33, 37, 39, 49, 64, 76, 78, 79, 83] can generate images that closely resemble the concept by fine-tuning the base text-to-image model in a specific image set. However, all tuning-based models will inherently suffer from the over-fitting introduced by this process, which leads to weakening in the compositional ability of the model. For example, when generating “a `<sk>` dog wearing a headphone”, though the given dog is well reconstructed, the headphone is always missing (Fig. 1). This feature affects the diversity of the generated output in practical use. A commonly accepted explanation within the community [20, 29, 64, 73] attributes this phenomenon to overfitting given a limited number of images. However, the fundamental cause of this overfit remains unexplored. In this work, our aim is to investigate the underlying causes behind the overfitting.

Upon initial examination, it appears that the model diminishes some of its original capabilities after personalization tuning. Taking Stable Diffusion (SD) [72] as an example, from Fig. 1, we observe that the base SD model indeed has the ability to combine the concepts of a dog and a headphone. However, after fine-tuning, the model struggles to achieve compositional generation; for instance,

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

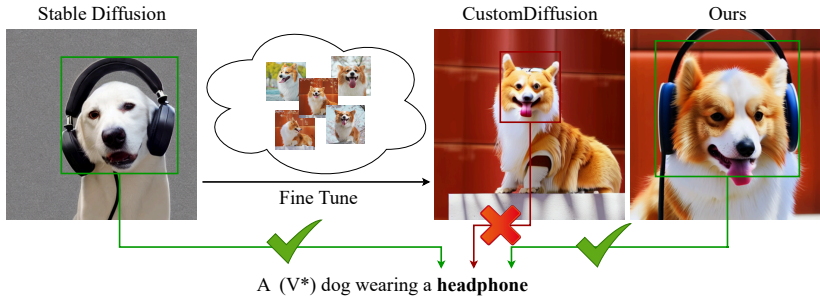


Figure 1: The base Stable Diffusion (SD) possesses the capability to compose the concept of a dog and headphone, generating a dog wearing a headphone. However, we notice that this compositional generation capability is lost during personalization tuning. For example, when using Custom Diffusion (CD) [37], the headphone is missing despite the target corgi is generated successfully. On the other hand, our method can successfully compose the target corgi with the headphone.

while the target concept  $\langle s_k s \rangle$  (dog) can be generated successfully, the headphone is missing. **We hypothesized that the decline in this compositional ability stems from the semantic drift of the target concept away from its superclass target during fine-tuning.** To better understand this, we conduct some empirical analysis by visualizing the CLIP text-space and cross-attention map activation area in Fig. 3a, 3b. In addition, we also perform theoretical analysis and find that the root cause lies in the semantic bias that reduces the entropy of the probability of the composed conditions, which significantly increases the difficulty to simultaneously sample the target concept combined with other elements.

Based on our experimental findings and theoretical analysis, we introduce ClassDiffusion to address the issue of weakening compositional capacity after fine-tuning. Fig. 2 shows the performance of our method. Our method uses semantic preservation loss to explicitly guide the model to restore the semantic imbalance that arises during the fine-tuning stage. In particular, it narrows the gap between the text embeddings of the target concept and its respective superclass in the textual space. Despite its simplicity, the proposed loss can successfully recover the compositional ability as shown in Fig. 1. Therefore, distinct from prevalent loss design in the community which seek to migrate the overfitting in tuning-based models, our method enhances the model’s capacity of following the text prompt while maintaining the concept of a customized subject. Extensive experiments have demonstrated the effectiveness of our method in restoring the compositional generation capability of the base model. Furthermore, we explore the potential of our approach in personalized video synthesis, showcasing its ability in recovering the semantical space of the generative model. In addition, we found that the CLIP-T metric can hardly reflect the actual performance of personalized generation. Therefore, we introduce the BLIP2-T metric, a more equitable and effective evaluation metric for this particular domain. To summarize, the contributions of our work are:

- We offer a thorough examination to understand why existing tuning-based subject-driven personalized methods inherently suffer from the loss of compositional ability. This is elucidated through both experimental observations and theoretical analysis.
- We propose ClassDiffusion, a simple technique to recover the compositional capabilities lost during personalized tuning.
- Extensive experiments demonstrate that the proposed technique achieves improved personalization ability in image and video generation tasks.

## 2 RELATED WORK

**Text-to-Image Generation and Its Control** Text-to-image generation is designed to generate high-quality, high-fidelity images that are aligned with textual prompts. This field has been under research for an extended period. Recently, the field of Text-to-image generation has made significant progress with extensive research in Generative Adversarial Network (GAN) [23, 24, 51, 93], Variational Autoencoder (VAE) [6, 13, 36, 74], and Diffusion models [31, 46, 53, 57, 61, 63, 65, 66, 70, 91].





Figure 2: A qualitative result of two small stories produced by our model. The above showcases a bear’s literary journey: from reading a book to ultimately earning a Nobel Literature Prize. The below shows the fate of a sunglasses. Finally, the bear gets the sunglasses. It shows a potential real-world application due to our model’s high performance.

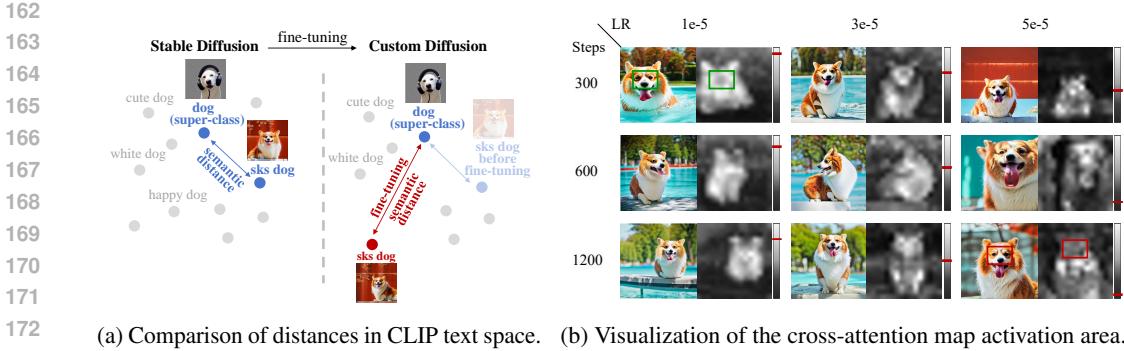
The diffusion models achieve a new state-of-the-art (SOTA) in unconditional image generation. Numerous works [12, 44, 53, 59–61, 65] have been done to make the image generated by diffusion models more aligned with the textual prompts. Among them, Stable Diffusion [63] is a widely recognized model in the field, utilizes a cross-attention mechanism to integrate textual conditions into the image generation process and employs the Latent Diffusion Model, which maps the image to latent space [63]. Our research is based on the Stable Diffusion framework due to its adaptability and wide use in the community. Furthermore, different methods exist for controlling generative models. The primary categories for controlling generative models typically include Text-guided [7, 18, 22, 58, 62], Image-guided [1, 20, 33, 37, 39, 64, 78], Additional Sparse conditions [2, 4, 5, 8, 11, 28, 35, 47, 58, 77, 86, 89, 90], Brain-guided [3, 10, 19, 45, 52, 55, 72], Sound-Guided [58, 87], and some universe control [41, 47, 88]. Text-guided control utilizes textual descriptions to directly influence the outcome, guiding the model based on specific verbal instructions. Our method focused on the text-guided controllable generative model.

**Subject-Driven Personalized Generation** Subject-driven personalized generation is focused on creating images based on reference images. Recent works [1, 9, 20, 21, 33, 37, 39, 49, 56, 64, 71, 76, 78, 79, 83, 85] have explored techniques for producing striking resemblance images in multiple ways. One of the primary ways is to fine-tune the base text-to-image models. Furthermore, there has been a significant effort in research [14, 29, 32, 34, 37, 38, 69, 78, 82] aimed at integrating various concepts in personalization. While striking resemblance images are produced, fine-tuning the base model on a small set of images leads to overfitting, resulting in unexpected issues. One prevalent issue discussed in prior research is the decrease in diversity. Recent studies have proposed various methods to address these issues. For instance, DreamBooth [64] introduced the Class-Specific Prior Loss, which effectively addresses diversity reduction by recovering the class’s prior knowledge. However, it does not effectively maintain the ability of the model to follow the text prompt. Another common issue is the inability to generate images under multiple conditions. Some Recent Research [29, 30, 39, 48, 73] proposed some methods to migrate this appearance. However, the underlying reasons for this phenomenon have not been thoroughly investigated. Our research endeavors to explore these reasons and develop solutions to overcome this challenge.

## 3 METHOD

### 3.1 PRELIMINARY

**Text-to-Image Diffusion Model** Stable Diffusion [63] is widely used in image generation task. For any input image, Stable Diffusion first transforms it into a latent representation  $x$  using the encoder  $\varepsilon$  of a variant auto-encoder [36]. For any input image, Stable Diffusion first transforms it into a latent representation  $x$  using the encoder  $\varepsilon$  of a variant auto-encoder [36]. The diffusion process then operates on  $x$  by incrementally introducing noise, resulting in a fixed-length Markov chain represented as  $x_1, x_2, \dots, x_T$ , where  $T$  is the chain’s length. Stable Diffusion uses a UNet architecture



(a) Comparison of distances in CLIP text space. (b) Visualization of the cross-attention map activation area.

Figure 3: (a) Each dot represents the position of a phrase combining an adjective and "dog" in the CLIP text-space. After fine-tuning, customized concepts move further away from the the distribution of super-class. (b) Visualization results of cross-attention map activation maps corresponding to the dog token. The bar chart on the right shows the average activation level in the dog area. Experiments show that the activation strengths of the corresponding classes decrease with the increase of the learning rate and the total number of training steps. These demonstrate that the customized concepts likely no longer belong to the super-class, resulting in a loss of super-class semantic information, such as wearing a headphone.

to learn the reverse of this diffusion process, predicting a denoised version of the latent input  $x_t$  at each timestep  $t$  from 1 to  $T$ . In the context of text-to-image generation, the text prompts' conditioning information  $y$  is encoded into an intermediate representation  $\tau_\theta(y) = c$ , where  $\tau_\theta$  is a pre-trained CLIP [59] text encoder. The primary objective in training this text-to-image diffusion model involves optimizing this transformation and prediction process, and it can be expressed as:

$$\mathcal{L}_{recon} = \mathbb{E}_{x,y,\epsilon,t} \left[ \|\epsilon - \epsilon_\theta(x_t, t, \tau_\theta(y))\|_2^2 \right] \tag{1}$$

where  $\epsilon$  and  $\epsilon_\theta$  represent the noise samples from the standard Gaussian noise  $\epsilon \sim \mathcal{N}(0, 1)$  and predicted noise residual, respectively.

**Subject-driven Diffusion Model** Although text-to-image models have achieved remarkable performance, their controllability is limited. To personalize the generated outputs, DreamBooth[64] fine-tunes the diffusion U-Net to fit several target concept images. Custom Diffusion[37] introduces a new modifier token  $V^*$  in front of the category name and optimizes only the key and value matrices in the cross-attention layers, thereby improving efficiency.

### 3.2 EXPERIMENTAL ANALYSIS

We begin by observing simple experimental test cases to realize that the loss of compositional ability after personalization tuning is a common phenomenon. We then analyze the underlying logic through visualizations of the CLIP text-space and cross-attention strength map. Finally, we conduct a theoretical analysis to support our hypothesis.

**Simple experimental test cases.** As shown in Fig. 1, we observe a loss of compositional ability after fine-tuning. We then conduct additional test cases and find that the headphone concept is not the only one affected; other concepts also experience a similar loss. Furthermore, this situation occurs in both the dog case and other classes.

**Semantic drift in CLIP text-space.** To elucidate the reasons, we project text into the CLIP text space and use two dimensions for simplified visualization in Fig. 3a. Each dot represents phrases composing the super-class ("dog" in this case) with different adjectives (e.g., "a cute dog", "a white dog" etc.). Before fine-tuning, the customized concepts (<sks> dog) have no special meaning and are at a similar distance from the super-class as other words. After fine-tuning, we observe a significant increase in the distance of the target concept from its corresponding super-class. This indicates that the semantics of the personalized concepts change during fine-tuning. In short, the model increasingly fails to recognize that personalized concepts belong to the dog category. This shift may lead to an

inability to access the knowledge associated with the super-class (like wearing a headphone). More details can be found in Appendix D, E.

**Reduction of cross-attention activation strength.** We further investigate the model by visualizing the cross-attention layers in Fig. 3b. The attention maps indicate the activation area of super-class words in cross-attention layers. It shows that while the "dog" token activates the relevant region in the image, its activation level is notably lower than that of the pre-trained model. Furthermore, its activation level decreases with the increase in epochs and the learning rate. These findings align with our observations in the CLIP text space and provide support for hypothesis that the customized concepts are increasingly not recognized as part of the super-class during fine-tuning.

Next, we theoretically analyze why the semantic drift of personalized phrases results in the weakening of the compositional ability from the respective reduction in the entropy of composable conditional probability.

### 3.3 THEORETICAL ANALYSIS

Drawing on the insights of [43], a trained diffusion model can be seen as implicitly defining an Energy-Based Model (EBM) [17]. This perspective allows us to build on prior research in composing EBMs and adapting them for use in diffusion models. Building on the work of [16] in the context of generating images with multiple attributes and Bayes' theorem, the conditional probability can be decomposed as:

$$p(x|c_{class}, c_1, c_2, \dots, c_i) \propto p(x, c_{class}, c_1, c_2, \dots, c_3) = p(c_{class}|x)p(x) \prod_{i \in T} p(c_i|x) \quad (2)$$

$$= p(c_{class}|x)p(x) \prod_{i \in T} \frac{p(c_i)p(x|c_i)}{p(x)} \quad (3)$$

where  $T$  is all the set of conditions in prompts except for the class,  $p(c_i)$  represents the probability of occurrence of condition  $c_i$  in the training dataset and can be regarded as a constant for large-scale pre-training models.  $p(c_i|x)$  represents an implicit classifier, denoting the probability of categorizing a concept as  $c_{class}$ . Specifically,  $p(c_{class}|x)$  represents a specific implicit classifier for the super-category. Thus, we have:

$$p(x|c_{class}, c_1, c_2, \dots, c_i) \propto p(c_{class}|x)p(x) \prod_{i \in T} \frac{p(c_i)p(x|c_i)}{p(x)} \quad (4)$$

Denoted  $p(x) \prod_{i \in T} \frac{p(c_i)p(x|c_i)}{p(x)}$  as  $d(x)$ ,  $p(c_{class}|x)$  as  $q(x)$ , and  $p(x|c_1, c_2, \dots, c_i)$  as  $a(x)$ . The entropy of  $a$  is calculated as:

$$H(a) = - \sum_x q(x)d(x) [\log(q(x)) + \log(d(x))] \quad (5)$$

After fine-tuning, the components of  $d(x)$  change only slightly and can be treated as unchanged, and the implicit classifier  $p_\theta(c_{class}|x)$  changes to  $p_{\theta'}(c_{class}|x)$ . Thus the difference in entropy before and after can be expressed as:

$$\Delta H = \sum_x q_\theta(x)d(x) [\log(q_\theta(x)) + \log(d(x))] \quad (6)$$

$$\begin{aligned} & - \sum_x q_{\theta'}(x)d(x) [\log(q_{\theta'}(x)) + \log(d(x))] \\ & = d(x) \sum_x \{ [q_\theta(x) \log q_\theta(x) - q_{\theta'}(x) \log q_{\theta'}(x)] \\ & \quad + \log d(x) [q_\theta(x) - q_{\theta'}(x)] \} \end{aligned} \quad (7)$$

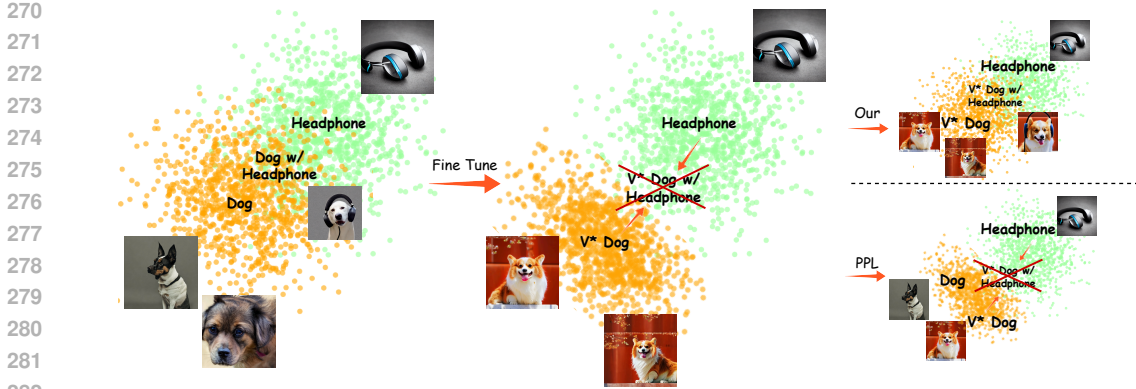


Figure 4: The orange and green point sets represent the distributions of dogs and headphones, respectively, and their overlapping regions represent their joint probability distributions. During the tuning process, the conditional distribution of dogs and headphones shrinks, which gradually increases the difficulty of sampling. Unlike the Prior Preservation Loss (PPL) in DreamBooth [64], which aims to maintain class diversity, our proposed Semantic Preservation Loss (SPL) focuses on recovering the semantic space of the customized concept. This approach enables our method to synthesize images that are more consistent with the text prompt.

Based on our observations in Fig. 3a, 3b we can show that  $q_{\theta}(x) > q_{\theta'}(x)$ , combining the properties of probability theory and the monotonically decreasing nature of  $x \log x$  at  $(0,1)$ , we have:

$$q_{\theta}(x) \log q_{\theta}(x) - q_{\theta'}(x) \log q_{\theta'}(x) < 0; \log d(x) < 0 \tag{8}$$

Thus, we have:

$$\Delta H(a) < 0 \tag{9}$$

As a result, it is more difficult to sample from our demanded conditional distributions under  $c_{\text{class}}, c_1, \dots, c_i$  conditions than before the fine-tuning, leading to the phenomenon that the combining ability is weakened after the fine-tuning. We will discuss the theoretical reasons here in more detail in Appendix B. Fig. 4 illustrates the changes in the distribution during this process that lead to a weakening of the compositional generation capability.

The diminished combinatorial capacity is due to the increased difficulty in sampling from joint conditional probabilities caused by shifts in the semantics of customisation concepts. Therefore, in order to reduce the difficulty of sampling in joint conditional probability distributions, we need to recover the original semantic space of the text, i.e. to recover the semantic distance between custom concepts and superclasses in the semantic space that has been distanced by the fine-tuning process. So, in the next section, we present our proposed semantic preservation loss to mitigate the semantic drift that occurred during fine-tuning.

### 3.4 SEMANTIC PRESERVATION LOSS

As analyzed above, the key challenge lies in preserving semantic information during fine-tuning to reduce the difficulty of sampling from the joint conditional distribution. To address this, we propose a novel loss function aimed at constraining semantic variation throughout the fine-tuning process.

Specifically, considering that there are  $N$  special tokens, each representing a customized concept. During the process, the embeddings associated with these tokens are fine-tuned to align with the target concepts. Our loss function is designed to minimize the semantic distance between the phrase containing the special token (e.g., a photo of a  $V^*$  dog) and the phrase containing only the class word (e.g., a photo of a dog). By labeling the training prompt as  $P_{tp}$  (e.g., a photo of a  $V^*$  dog), and the class prompt as  $P_{cp}$  (e.g., a photo of a dog), we use a text encoder to get their embeddings  $E_{tp}$  and  $E_{cp}$  in Stable Diffusion’s semantic space. The semantic preservation loss (SPL) is calculated by the sum of the cosine distance of their text embeddings. Formally speaking, our proposed SPL can be

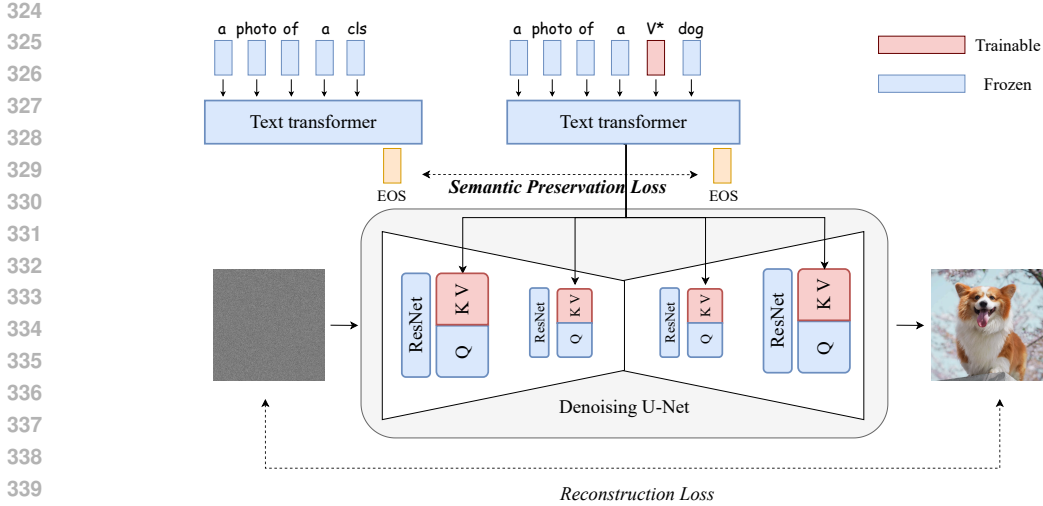


Figure 5: The framework of ClassDiffusion. The personalization fine-tuning strategy is based on Custom Diffusion [38], which primarily fine-tunes the K and V parameters in the transformer block. Our **semantic preservation loss (SPL)** is calculated by measuring the cosine distance between text features extracted from the same text transformer (using EOS tokens as text features following CLIP) for phrases with personalized tokens and phrases with only super-class.

expressed by the following equation:

$$\mathcal{L}_{sp} = \sum^N \sum^B \sum^L D_c(E_{SC}, E_C) \tag{10}$$

where  $B$  represents the batch size,  $L$  denotes the hidden dimensions of the text encoder, and  $D_c$  implies the cosine distance. We can represent the final training objective as:

$$\mathcal{L} = \mathcal{L}_{recon} + \lambda \mathcal{L}_{sp} \tag{11}$$

The overview of our proposed model is shown in Fig. 5

## 4 EXPERIMENTS

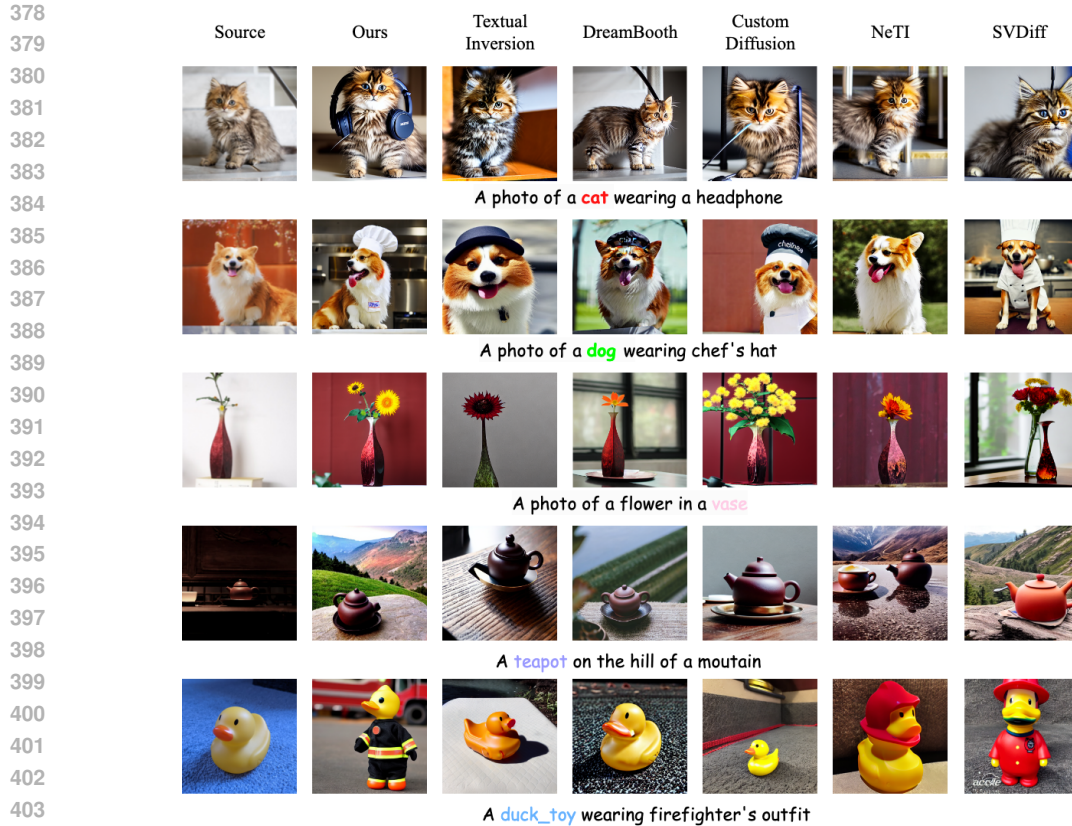
### 4.1 EXPERIMENT DETAILS

**Implementation details** Our method is built on Stable Diffusion V1.5, with a learning rate  $10^{-6}$ , and batch size 2 for fine-tuning. We used 500 optimization steps for a single concept and 800 for multiple concepts, respectively. During inference, the guidance scale is set to 6.0 and the inference steps are set to 100. The semantical preservation loss weight is set to 1.0 during all experiments. All experiments are conducted on  $2 \times \text{RTX4090}$  GPUs. Our method uses  $\sim 6$  min for single concept generation and  $\sim 11$  min for multiple concept generation.

To better preserve the semantic space, we compute SPL between text embeddings embedded in the semantic space of the Stable Diffusion model. Therefore, we utilize the CLIP [59] text encoder from Stable Diffusion v1.5 [61], specifically clip-vit-large-patch14 [45], to extract the text embeddings of phrases. Following common practice, we use the End of Sequence (EOS) token to represent the semantics of embeddings.

**Baselines** We compare our method with state-of-the-art (SOTA) competitors, including DreamBooth [64], Textual Inversion [20], Custom Diffusion [37], NeTI [1], SVDiff [29]. For DreamBooth, CustomDiffusion, and Textual Inversion, we used the diffusers [75] version of the implementation. For NeTI, we use its official implementation. Given that SVDiff does not have an official open-source repository. For SVDiff, we use the implementation of [67]. All training parameters follow the recommendations of the official paper. To ensure fairness of comparison, all these baselines are built on Stable Diffusion V1.5.





405 Figure 6: Qualitative comparison between our method and baselines with single given concept. Our  
 406 method generates images that align with the prompts, surpassing all baselines.

407 **Datasets** Following previous work [29, 64, 73], we conduct quantitative experiments on DreamBooth  
 408 Dataset [64]. It contains 30 objects including both live objects and non-live objects. In addition,  
 409 we used images from the Textual Inversion Dataset [20] and CustomConcept101 [37] in qualitative  
 410 experiments.

411 **Evaluation metrics** We assess our approach using three metrics: CLIP-I, CLIP-T, and DINO-I.  
 412 CLIP-I calculates the visual similarity between the produced images and the target concept images  
 413 by utilizing CLIP [59] visual features. CLIP-T evaluates the similarity between text prompts and  
 414 images. If one baseline contains the special token  $S^*$ , it will be replaced with a prior class word. In  
 415 the case of DINO-I, we evaluate the cosine similarity between the ViT-S/16 DINO [54] embeddings  
 416 of the generated images and the concept images. Further, we note the impact of CLIP’s outdated  
 417 performance on the fairness of the evaluation. Therefore, we introduce the BLIP2-T Score, which  
 418 calculates the similarity between text features extracted from BLIP2’s Q-former and image features  
 419 extracted from Vision Encoder as a score. **This metric is designed by calculating the similarity  
 420 between image and text embeddings extracted by the BLIP2 model.** Our approach involved utilizing  
 421 the Transformer [80] implementation and the fine-tuned weights of BLIP2-IMT on CoCo [42], with  
 422 ViT/L [15]. This new metric aims to offer a more equitable and efficient evaluation measure for  
 423 future studies in this field. Empirical findings from various studies [25, 39, 40, 68, 81, 84] indicate  
 424 that BLIP2 outperforms CLIP significantly in the assessment of text-image alignment.

## 425 4.2 QUALITATIVE & QUANTITATIVE EXPERIMENTS

427 **Qualitative Experiments** We compare our method with DreamBooth [64], Textual Inversion [20],  
 428 NeTI [1], SVDiff [29], and Custom Diffusion [37] on challenging prompts. The results depicted  
 429 in Fig. 6 demonstrate the outcomes obtained from these prompts. Fig. 2 shows a story of a given  
 430 dog and sunglasses. The experimental findings indicate a substantial superiority of our approach  
 431 over other techniques regarding alignment with text prompts, without any decline in similarity  
 to the specified concept. More qualitative results are shown in Appendix. J. Also, we conduct



Figure 7: Qualitative comparison between our method and Custom Diffusion(CD) in multiple concepts. Our method has better text alignment than custom diffusion.

qualitative experiments with multiple concepts on the combinations <cat>, <sunglasses>; <bear>, <barn>; <dog>, <backpack>; Fig. 7 shows the results of the experiments. The experiments show that our method can be aligned with prompts better than custom diffusion in multi-concept generation.

**Quantitative Experiments** Following the previous work, we used 20 concepts for quantitative experiments. For Single Concept text similarity metrics (CLIP-T, BLIP2-T), we followed the 25 prompts used in [64], sampling 20 images per prompt. The results of the experiment are shown in Tab. 1. Experimental results show that our method obtains new SOTA on each text similarity metric, indicating that we have good compositional generation capability.

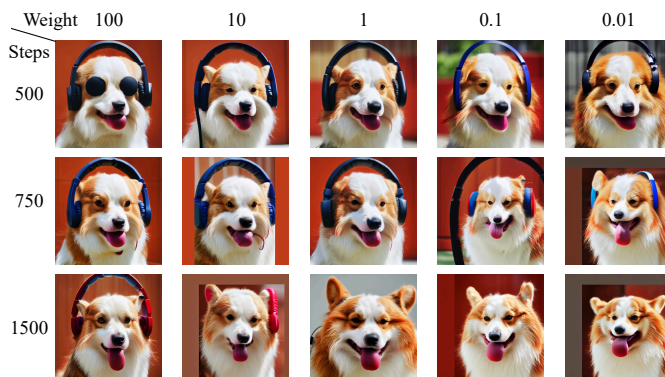
Table 1: Quantitative Results on all Metrics and Results of the User Study. The last two columns display the win rates of our method compared to other approaches in the user study, evaluated in terms of text similarity and image similarity. A win rate exceeding 50% (highlight by ✓) indicates that our method outperforms the compared methods on the corresponding metric, as judged from a human perspective.

	Method	CLIP-T↑	CLIP-I↑	DINO-I↑	BLIP2-T↑	TIFA↑	User-T Win Rate ↑	User-I Win Rate ↑
<b>Single Concept</b>	DreamBooth [64]	0.249	<b>0.855</b>	<b>0.700</b>	0.295	<b>0.559</b>	95.4% ✓	42.1%
	Textual Inversion [20]	0.242	0.825	0.631	0.308	<b>0.505</b>	95.1% ✓	75.0% ✓
	Custom Diffusion [37]	0.286	0.837	0.693	0.416	<b>0.746</b>	79.1% ✓	40.0%
	NeTI [1]	0.290	0.838	0.648	0.329	<b>0.607</b>	78.8% ✓	70.0% ✓
	SVDiff [29]	0.293	0.834	0.606	0.418	<b>0.835</b>	56.6% ✓	95.8% ✓
	Our	<b>0.300</b>	0.828	0.673	<b>0.460</b>	<b>0.843</b>	-	-
<b>Multiple Concepts</b>	Custom Diffusion [37]	0.282	0.813	<b>0.636</b>	0.380	-	-	-
	Our	<b>0.320</b>	<b>0.821</b>	0.604	<b>0.477</b>	-	-	-

### 4.3 ADDITIONAL EXPERIMENTS

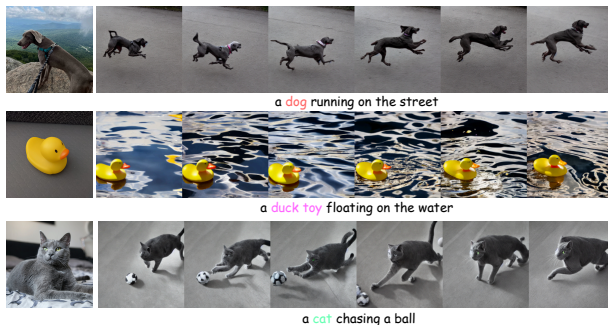
**User Study** We also performed user study to validate the effectiveness of our method. We used the same set of images generated in the Section 4.2 for user study, details of which are available in Appendix. I. The results of the user study are located in Tab. 1. The numbers in the table indicate at what percentage our method is considered by humans to be superior to the compared methods. The result of the user study shows that our method outperforms all methods in text similarity (> 50%). Although our method does not outperform all methods in image similarity, given the high CLIP-I scores, our method still produces images that are highly consistent with the given concept.

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499



500 Figure 8: Generation results for the prompt “a photo of a dog wearing a headphone” with different  
501 step counts and SPL weights. All results are generated using the same random seed.

502  
503  
504  
505  
506  
507  
508  
509  
510  
511



512  
513 Figure 9: Result of generated videos, showing good textual alignment and similarity of given concepts.

514  
515  
516  
517  
518  
519  
520  
521  
522

**Personalized Video Generation** We investigate the implementation of our method in personalized video generation. We utilized AnimateDiff V2 [26, 27] for video generation, configuring parameters to a resolution of  $512 \times 512$ , a guidance scale of 7.5, and 25 inference steps. The outcomes of the video generation process are illustrated in Fig. 9. Utilizing AnimateDiff, our technique produces videos that exhibit strong textual and conceptual coherence without the need for additional training. This demonstrates that our approach, which aligns personalized phrases with superclass-centric semantics, can generate engaging videos with dynamic generation capabilities stemming from pre-training, along with the ability to transition across corresponding domains.

523  
524  
525  
526  
527

**Abalation Experiments** We studied the influence of different weights of semantic preservation loss (SPL). The results show that higher SPL weights preserve combining ability better. At s SPL weight of 100, all steps successfully depict "wearing a headphone." However, at lower weights of 0.1 and 0.01, the headphone details diminish by 750 steps and disappear by 1500 steps. On the other hand, lower SPL weights restore specific concept features more effectively.

528

## 529 5 CONCLUSIONS

530  
531  
532  
533  
534  
535  
536  
537  
538  
539

In this work, we highlight the problem of weakened compositional ability due to individualized fine-tuning and provide an analysis of the causes of this problem from experimental observations and information-theoretic perspectives. We discovered that this weakening effect is primarily attributed to the semantic shift of the customized concepts throughout the fine-tuning process. As the model undergoes fine-tuning, the representations of these concepts gradually drift away from their original meanings, leading to a misalignment with the intended semantics. This semantic drift complicates the model’s ability to accurately sample from the joint conditional distribution, ultimately hindering its performance in generating or understanding the intended outcomes based on the fine-tuned concepts. We then introduce a new approach, termed ClassDiffusion, which mitigates the weakening of compositional ability by restoring the original semantic space. Finally, we present comprehensive experimental results showcasing the efficacy of ClassDiffusion and the fresh perspectives it offers on interconnected fields.

## REFERENCES

- [1] Yuval Alaluf, Elad Richardson, Gal Metzer, and Daniel Cohen-Or. A neural space-time representation for text-to-image personalization. *ACM Transactions on Graphics (TOG)*, 42(6):1–10, 2023.
- [2] Omri Avrahami, Thomas Hayes, Oran Gafni, Sonal Gupta, Yaniv Taigman, Devi Parikh, Dani Lischinski, Ohad Fried, and Xi Yin. Spatext: Spatio-textual representation for controllable image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18370–18380, 2023.
- [3] Yunpeng Bai, Xintao Wang, Yan-pei Cao, Yixiao Ge, Chun Yuan, and Ying Shan. Dreamdiffusion: Generating high-quality images from brain eeg signals. *arXiv preprint arXiv:2306.16934*, 2023.
- [4] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Qinsheng Zhang, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, et al. ediff-i: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022.
- [5] Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Universal guidance for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 843–852, 2023.
- [6] 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.
- [7] 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(4):1–10, 2023.
- [8] Kai Chen, Enze Xie, Zhe Chen, Lanqing Hong, Zhenguo Li, and Dit-Yan Yeung. Integrating geometric control into text-to-image diffusion models for high-quality detection data generation via text prompt. *arXiv preprint arXiv:2306.04607*, 2023.
- [9] Wenhu Chen, Hexiang Hu, Chitwan Saharia, and William W Cohen. Re-imagen: Retrieval-augmented text-to-image generator. *arXiv preprint arXiv:2209.14491*, 2022.
- [10] Zijiao Chen, Jiabin Qing, Tiange Xiang, Wan Lin Yue, and Juan Helen Zhou. Seeing beyond the brain: Conditional diffusion model with sparse masked modeling for vision decoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22710–22720, 2023.
- [11] Jiabin Cheng, Xiao Liang, Xingjian Shi, Tong He, Tianjun Xiao, and Mu Li. Layoutdiffuse: Adapting foundational diffusion models for layout-to-image generation. *arXiv preprint arXiv:2302.08908*, 2023.
- [12] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- [13] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. *Advances in Neural Information Processing Systems*, 34:19822–19835, 2021.
- [14] Ziyi Dong, Pengxu Wei, and Liang Lin. Dreamartist: Towards controllable one-shot text-to-image generation via positive-negative prompt-tuning. *arXiv preprint arXiv:2211.11337*, 2022.
- [15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [16] Yilun Du, Shuang Li, and Igor Mordatch. Compositional visual generation with energy based models. *Advances in Neural Information Processing Systems*, 33:6637–6647, 2020.
- [17] Yilun Du and Igor Mordatch. Implicit generation and modeling with energy based models. *Advances in Neural Information Processing Systems*, 32, 2019.



- 594 [18] Weixi Feng, Xuehai He, Tsu-Jui Fu, Varun Jampani, Arjun Akula, Pradyumna Narayana, Sugato  
595 Basu, Xin Eric Wang, and William Yang Wang. Training-free structured diffusion guidance for  
596 compositional text-to-image synthesis. *arXiv preprint arXiv:2212.05032*, 2022.
- 597 [19] Honghao Fu, Zhiqi Shen, Jing Jih Chin, and Hao Wang. Brainvis: Exploring the bridge between  
598 brain and visual signals via image reconstruction. *arXiv preprint arXiv:2312.14871*, 2023.
- 599 [20] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and  
600 Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using  
601 textual inversion. In *The Eleventh International Conference on Learning Representations*, 2023.
- 602 [21] Rinon Gal, Moab Arar, Yuval Atzmon, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or.  
603 Designing an encoder for fast personalization of text-to-image models. *arXiv e-prints*, pages  
604 arXiv-2302, 2023.
- 605 [22] Songwei Ge, Taesung Park, Jun-Yan Zhu, and Jia-Bin Huang. Expressive text-to-image  
606 generation with rich text. In *Proceedings of the IEEE/CVF International Conference on*  
607 *Computer Vision*, pages 7545–7556, 2023.
- 608 [23] Xinyu Gong, Shiyu Chang, Yifan Jiang, and Zhangyang Wang. Autogan: Neural architecture  
609 search for generative adversarial networks. In *Proceedings of the IEEE/CVF International*  
610 *Conference on Computer Vision*, pages 3224–3234, 2019.
- 611 [24] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
612 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural*  
613 *information processing systems*, 27, 2014.
- 614 [25] Paul Grimal, Hervé Le Borgne, Olivier Ferret, and Julien Tourille. Tiam – a metric for evaluating  
615 alignment in text-to-image generation, 2024.
- 616 [26] Yuwei Guo, Ceyuan Yang, Anyi Rao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Sparsectrl:  
617 Adding sparse controls to text-to-video diffusion models. *arXiv preprint arXiv:2311.16933*,  
618 2023.
- 619 [27] Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh  
620 Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image  
621 diffusion models without specific tuning. *International Conference on Learning Representations*,  
622 2024.
- 623 [28] Cusuh Ham, James Hays, Jingwan Lu, Krishna Kumar Singh, Zhifei Zhang, and Tobias Hinz.  
624 Modulating pretrained diffusion models for multimodal image synthesis. In *ACM SIGGRAPH*  
625 *2023 Conference Proceedings*, pages 1–11, 2023.
- 626 [29] Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris Metaxas, and Feng Yang.  
627 Svdiff: Compact parameter space for diffusion fine-tuning. In *Proceedings of the IEEE/CVF*  
628 *International Conference on Computer Vision*, pages 7323–7334, 2023.
- 629 [30] Yutong He, Alexander Robey, Naoki Murata, Yiding Jiang, Joshua Williams, George J Pappas,  
630 Hamed Hassani, Yuki Mitsufuji, Ruslan Salakhutdinov, and J Zico Kolter. Automated black-box  
631 prompt engineering for personalized text-to-image generation. *arXiv preprint arXiv:2403.19103*,  
632 2024.
- 633 [31] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances*  
634 *in neural information processing systems*, 33:6840–6851, 2020.
- 635 [32] Hexiang Hu, Kelvin C. K. Chan, Yu-Chuan Su, Wenhui Chen, Yandong Li, Kihyuk Sohn, Yang  
636 Zhao, Xue Ben, Boqing Gong, William Cohen, Ming-Wei Chang, and Xuhui Jia. Instruct-  
637 imagen: Image generation with multi-modal instruction, 2024.
- 638 [33] Miao Hua, Jiawei Liu, Fei Ding, Wei Liu, Jie Wu, and Qian He. Dreamtuner: Single image is  
639 enough for subject-driven generation. *arXiv preprint arXiv:2312.13691*, 2023.
- 640 [34] Sangwon Jang, Jaehyeong Jo, Kimin Lee, and Sung Ju Hwang. Identity decoupling for multi-  
641 subject personalization of text-to-image models, 2024.
- 642 [35] Xuan Ju, Ailing Zeng, Chenchen Zhao, Jianan Wang, Lei Zhang, and Qiang Xu. Humansd:  
643 A native skeleton-guided diffusion model for human image generation. In *Proceedings of the*  
644 *IEEE/CVF International Conference on Computer Vision*, pages 15988–15998, 2023.
- 645 [36] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*  
646 *arXiv:1312.6114*, 2013.



- 648 [37] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-  
649 concept customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF Conference*  
650 *on Computer Vision and Pattern Recognition*, pages 1931–1941, 2023.
- 651 [38] Gihyun Kwon, Simon Jenni, Dingzeyu Li, Joon-Young Lee, Jong Chul Ye, and Fabian Caba  
652 Heilbron. Concept weaver: Enabling multi-concept fusion in text-to-image models, 2024.
- 653 [39] Dongxu Li, Junnan Li, and Steven Hoi. Blip-diffusion: Pre-trained subject representation for  
654 controllable text-to-image generation and editing. *Advances in Neural Information Processing*  
655 *Systems*, 36, 2024.
- 656 [40] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image  
657 pre-training with frozen image encoders and large language models. In *International conference*  
658 *on machine learning*, pages 19730–19742. PMLR, 2023.
- 659 [41] Wei Li, Xue Xu, Jiachen Liu, and Xinyan Xiao. Unimo-g: Unified image generation through  
660 multimodal conditional diffusion, 2024.
- 661 [42] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
662 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer*  
663 *Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014,*  
664 *Proceedings, Part V 13*, pages 740–755. Springer, 2014.
- 665 [43] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional  
666 visual generation with composable diffusion models. In *European Conference on Computer*  
667 *Vision*, pages 423–439. Springer, 2022.
- 668 [44] Xihui Liu, Dong Huk Park, Samaneh Azadi, Gong Zhang, Arman Chopikyan, Yuxiao Hu,  
669 Humphrey Shi, Anna Rohrbach, and Trevor Darrell. More control for free! image synthesis  
670 with semantic diffusion guidance. In *Proceedings of the IEEE/CVF Winter Conference on*  
671 *Applications of Computer Vision*, pages 289–299, 2023.
- 672 [45] Yizhuo Lu, Changde Du, Qiongyi Zhou, Dianpeng Wang, and Huiguang He. Minddiffuser: Con-  
673 trolled image reconstruction from human brain activity with semantic and structural diffusion.  
674 In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 5899–5908,  
675 2023.
- 676 [46] Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models:  
677 Synthesizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*,  
678 2023.
- 679 [47] Yuanhuiyi Lyu, Xu Zheng, and Lin Wang. Image anything: Towards reasoning-coherent and  
680 training-free multi-modal image generation, 2024.
- 681 [48] Jian Ma, Junhao Liang, Chen Chen, and Haonan Lu. Subject-diffusion: Open domain personal-  
682 ized text-to-image generation without test-time fine-tuning. *arXiv preprint arXiv:2307.11410*,  
683 2023.
- 684 [49] Yiyang Ma, Huan Yang, Wenjing Wang, Jianlong Fu, and Jiaying Liu. Unified multi-modal  
685 latent diffusion for joint subject and text conditional image generation. *arXiv preprint*  
686 *arXiv:2303.09319*, 2023.
- 687 [50] David JC MacKay. *Information theory, inference and learning algorithms*. Cambridge university  
688 press, 2003.
- 689 [51] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint*  
690 *arXiv:1411.1784*, 2014.
- 691 [52] Pengyu Ni and Yifeng Zhang. Natural image reconstruction from fmri based on self-supervised  
692 representation learning and latent diffusion model. In *Proceedings of the 15th International*  
693 *Conference on Digital Image Processing*, pages 1–9, 2023.
- 694 [53] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,  
695 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing  
696 with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- 697 [54] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,  
698 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning  
699 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 700  
701

- 702 [55] Furkan Ozcelik and Rufin VanRullen. Natural scene reconstruction from fmri signals using  
703 generative latent diffusion. *Scientific Reports*, 13(1):15666, 2023.
- 704 [56] Jae Wan Park, Sang Hyun Park, Jun Young Koh, Junha Lee, and Min Song. Cat: Contrastive  
705 adapter training for personalized image generation, 2024.
- 706 [57] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe  
707 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image  
708 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- 709 [58] Can Qin, Ning Yu, Chen Xing, Shu Zhang, Zeyuan Chen, Stefano Ermon, Yun Fu, Caiming  
710 Xiong, and Ran Xu. Gluegen: Plug and play multi-modal encoders for x-to-image generation.  
711 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23085–  
712 23096, 2023.
- 713 [59] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
714 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
715 models from natural language supervision. In *International conference on machine learning*,  
716 pages 8748–8763. PMLR, 2021.
- 717 [60] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,  
718 Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified  
719 text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 720 [61] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical  
721 text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3,  
722 2022.
- 723 [62] Royi Rassin, Eran Hirsch, Daniel Glickman, Shauli Ravfogel, Yoav Goldberg, and Gal Chechik.  
724 Linguistic binding in diffusion models: Enhancing attribute correspondence through attention  
725 map alignment. *Advances in Neural Information Processing Systems*, 36, 2024.
- 726 [63] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
727 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF*  
728 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, June  
729 2022.
- 730 [64] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.  
731 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In  
732 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
733 22500–22510, 2023.
- 734 [65] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton,  
735 Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al.  
736 Photorealistic text-to-image diffusion models with deep language understanding. *Advances in*  
737 *neural information processing systems*, 35:36479–36494, 2022.
- 738 [66] Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion  
739 distillation, 2023.
- 740 [67] Mingkai Shing. Svdiff: Stochastic video diffusion for conditional video generation. <https://github.com/mkshing/svdiff-pytorch>, 2023.
- 741 [68] Jaskirat Singh and Liang Zheng. Divide, evaluate, and refine: Evaluating and improving  
742 text-to-image alignment with iterative vqa feedback. In *Thirty-seventh Conference on Neural*  
743 *Information Processing Systems*, 2023.
- 744 [69] Kihyuk Sohn, Nataniel Ruiz, Kimin Lee, Daniel Castro Chin, Irina Blok, Huiwen Chang, Jarred  
745 Barber, Lu Jiang, Glenn Entis, Yuanzhen Li, et al. Styledrop: Text-to-image generation in any  
746 style. *arXiv preprint arXiv:2306.00983*, 2023.
- 747 [70] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*  
748 *preprint arXiv:2010.02502*, 2020.
- 749 [71] Kunpeng Song, Yizhe Zhu, Bingchen Liu, Qing Yan, Ahmed Elgammal, and Xiao Yang. Moma:  
750 Multimodal llm adapter for fast personalized image generation, 2024.
- 751 [72] Yu Takagi and Shinji Nishimoto. High-resolution image reconstruction with latent diffusion  
752 models from human brain activity. In *Proceedings of the IEEE/CVF Conference on Computer*  
753 *Vision and Pattern Recognition*, pages 14453–14463, 2023.

- 756 [73] Yoad Tewel, Rinon Gal, Gal Chechik, and Yuval Atzmon. Key-locked rank one editing for  
757 text-to-image personalization. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–11,  
758 2023.
- 759 [74] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive mod-  
760 eling: Scalable image generation via next-scale prediction. *arXiv preprint arXiv:2404.02905*,  
761 2024.
- 762 [75] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif  
763 Rasul, Mishig Davaadorj, Dhruv Nair, Sayak Paul, William Berman, Yiyi Xu, Steven Liu,  
764 and Thomas Wolf. Diffusers: State-of-the-art diffusion models. [https://github.com/  
765 huggingface/diffusers](https://github.com/huggingface/diffusers), 2022.
- 766 [76] Anton Voronov, Mikhail Khoroshikh, Artem Babenko, and Max Ryabinin. Is this loss informa-  
767 tive? faster text-to-image customization by tracking objective dynamics. *Advances in Neural  
768 Information Processing Systems*, 36, 2024.
- 769 [77] Andrey Voynov, Kfir Aberman, and Daniel Cohen-Or. Sketch-guided text-to-image diffusion  
770 models. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–11, 2023.
- 771 [78] Andrey Voynov, Qinghao Chu, Daniel Cohen-Or, and Kfir Aberman.  $p+$ : Extended textual  
772 conditioning in text-to-image generation. *arXiv preprint arXiv:2303.09522*, 2023.
- 773 [79] Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. Elite:  
774 Encoding visual concepts into textual embeddings for customized text-to-image generation.  
775 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15943–  
776 15953, 2023.
- 777 [80] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony  
778 Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer,  
779 Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain  
780 Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-  
781 art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods  
782 in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020.  
783 Association for Computational Linguistics.
- 784 [81] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng  
785 Li. Human preference score v2: A solid benchmark for evaluating human preferences of  
786 text-to-image synthesis. *arXiv preprint arXiv:2306.09341*, 2023.
- 787 [82] You Wu, Kean Liu, Xiaoyue Mi, Fan Tang, Juan Cao, and Jintao Li. U-vap: User-specified visual  
788 appearance personalization via decoupled self augmentation. *arXiv preprint arXiv:2403.20231*,  
789 2024.
- 790 [83] Chendong Xiang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. A closer look at parameter-  
791 efficient tuning in diffusion models. *arXiv preprint arXiv:2303.18181*, 2023.
- 792 [84] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao  
793 Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation,  
794 2023.
- 795 [85] Zhongcong Xu, Jianfeng Zhang, Jun Hao Liew, Hanshu Yan, Jia-Wei Liu, Chenxu Zhang, Jiashi  
796 Feng, and Mike Zheng Shou. Magicanimate: Temporally consistent human image animation  
797 using diffusion model. *arXiv preprint arXiv:2311.16498*, 2023.
- 798 [86] Han Xue, Zhiwu Huang, Qianru Sun, Li Song, and Wenjun Zhang. Freestyle layout-to-image  
799 synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
800 Recognition*, pages 14256–14266, 2023.
- 801 [87] Yue Yang, Kaipeng Zhang, Yuying Ge, Wenqi Shao, Zeyue Xue, Yu Qiao, and Ping Luo. Align,  
802 adapt and inject: Sound-guided unified image generation. *arXiv preprint arXiv:2306.11504*,  
803 2023.
- 804 [88] Yuyang Yin, Dejia Xu, Chuangchuang Tan, Ping Liu, Yao Zhao, and Yunchao Wei. Cle diffusion:  
805 Controllable light enhancement diffusion model. In *Proceedings of the 31st ACM International  
806 Conference on Multimedia*, pages 8145–8156, 2023.
- 807 [89] Jiwen Yu, Yinhuai Wang, Chen Zhao, Bernard Ghanem, and Jian Zhang. Freedom: Training-  
808 free energy-guided conditional diffusion model. In *Proceedings of the IEEE/CVF International  
809 Conference on Computer Vision*, pages 23174–23184, 2023.

810 [90] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image  
811 diffusion models, 2023.  
812  
813 [91] Shiwei Zhang, Jiayu Wang, Yingya Zhang, Kang Zhao, Hangjie Yuan, Zhiwu Qin, Xiang Wang,  
814 Deli Zhao, and Jingren Zhou. I2vgen-xl: High-quality image-to-video synthesis via cascaded  
815 diffusion models. *arXiv preprint arXiv:2311.04145*, 2023.  
816 [92] Yuxin Zhang, Weiming Dong, Fan Tang, Nisha Huang, Haibin Huang, Chongyang Ma, Tong-  
817 Yee Lee, Oliver Deussen, and Changsheng Xu. Prospect: Prompt spectrum for attribute-aware  
818 personalization of diffusion models. *ACM Transactions on Graphics (TOG)*, 42(6):244:1–  
819 244:14, 2023.  
820 [93] Junbo Zhao, Michael Mathieu, and Yann LeCun. Energy-based generative adversarial network.  
821 *arXiv preprint arXiv:1609.03126*, 2016.  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

## A VISUALIZATION AFTER FINE TUNNING

In section 3, we present visualization results for the text space and the cross-attention layer of other methods, highlighting the semantic bias that emerges in the text space, leading to a decline in compositional ability. To further affirm the effectiveness of our method and our hypothesis, we also visualize the textual feature space and the cross-attention layer with our method in this section. The visualization results are depicted in Fig. 10. The ranking of the distance is 36 out of 71, as opposed to 26 out of 71 before fine-tuning and 67 out of 71 for other methods. A comparison with the visualizations in Fig. 3b and Fig. 3a reveals that our model effectively addresses the semantic drift in the text space.

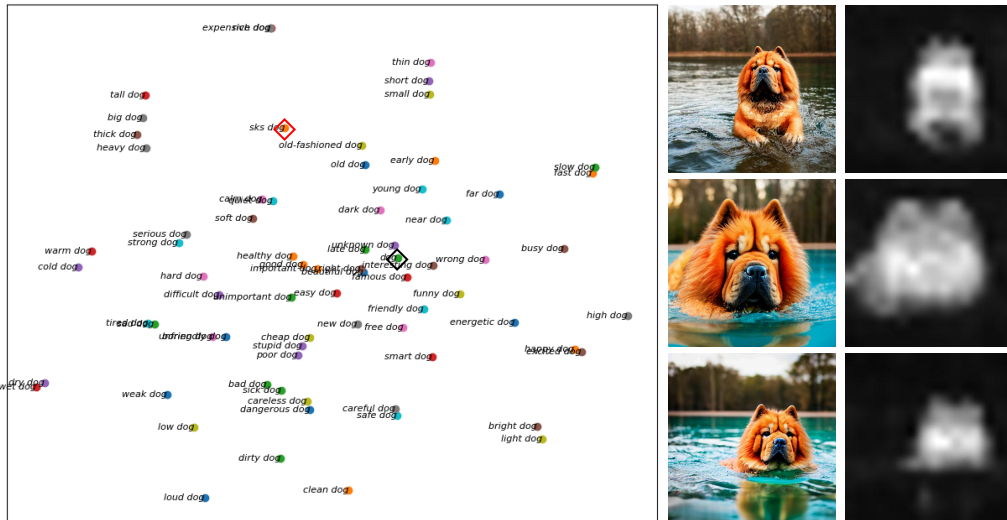


Figure 10: Visualization results after fine-tuning of our approach.

## B ENTROPY REDUCTION DURING THE FINE-TUNING

In section 3.3, we provide a solid theoretical analysis for the weakening of compositional ability due to the semantic drift from the perspective of information theory and probability distributions. In this section, we will discuss it in a more detailed way.

In the field of subject-driven personalization generation, two manifest phenomena are caused by overfitting: weakening of diversity in classes of given concepts and weakening of compositional ability. In addition to the calculations mentioned in the main body of the text, the entropy of combined conditional probability can also be calculated as conditional entropy:

$$H(X|c_1, c_2, \dots, c_i) \quad (1)$$

Make the given concept into a series of specific conditions:  $c_{s1}, \dots, c_{si}$ , each condition describes one of the features of the given concept. The entropy after fine-tuning will be:

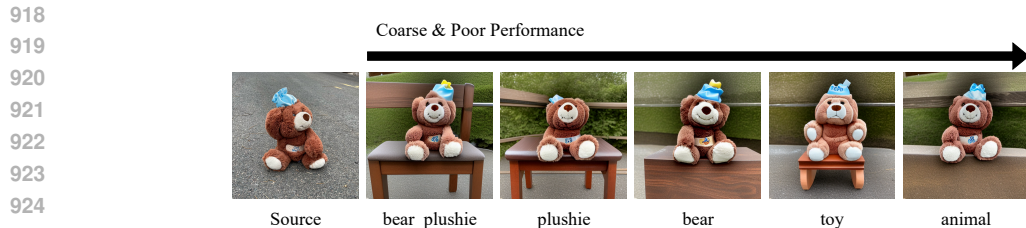
$$H(X|c_1, c_2, \dots, c_n, c_{s1}, \dots, c_{si}) = H(X|c_1 \dots, c_n) - I(X|c_1 \dots, c_n; c_{s1}, \dots, c_{si}) \quad (2)$$

Where  $I$  represent the mutual information, According to [50], we have:

$$I(X|c_1 \dots, c_n; c_{s1}, \dots, c_{si}) \geq 0 \quad (3)$$

$$H(X|c_1, c_2, \dots, c_n, c_{s1}, \dots, c_{si}) < H(X|c_1, c_2, \dots, c_i) \quad (4)$$





926 Figure 11: Visualization of generating “a <center word> sitting on the chair.”. This <center  
 927 word> is used in both prompt and class token. Experiments show that a fine-grained center word  
 928 will benefit our proposed method.

929

930 However, the entropy reduction here is different from the entropy reduction in the main text. The  
 931 entropy reduction here leads to a reduction in the diversity of the generated images. Specifically,  
 932 when the cue word is “a photo of a dog”, the image generated is closer to the given concept than to  
 933 the diversity of dogs.

### 934 C FINE-GRAINED EXPERIMENTS

935

936

937 At the core of our approach, we want the semantics of personalized phrases to be closer to the  
 938 category-centered words. In this section, we explore the effect of different category-centered words  
 939 on the results for the same given concept. Fig. 11 shows the training results using different category-  
 940 centered words. The results show that the use of different center words leads to significant differences  
 941 in performance, and a fine-grained center word benefits our method. Also, it indicates the importance  
 942 of recovering the semantical space of the customized concept.

### 943 D PROMPT USED IN THE VISUALIZATION OF CLIP TEXT SAMPLE SPACE

944

945

946 In this section, we provide a realistic visualization of the schematic in Fig. 3a, and discuss the prompt  
 947 to generate the 70 phrases that include adjectives and super-categories “dog”, and the whole 71  
 948 adjectives.

949 The realistic visualization of the schematic is:

950 The prompt we use is:

951 *Please help me generate some adjectives that can describe an attribute of a dog in a photo.*

952 The adjectives we use are:

953

954

955

<i>beautiful</i>	<i>happy</i>	<i>sad</i>	<i>tall</i>	<i>short</i>
<i>bright</i>	<i>dark</i>	<i>big</i>	<i>small</i>	<i>young</i>
<i>old</i>	<i>fast</i>	<i>slow</i>	<i>warm</i>	<i>cold</i>
<i>soft</i>	<i>hard</i>	<i>heavy</i>	<i>light</i>	<i>strong</i>
<i>weak</i>	<i>good</i>	<i>bad</i>	<i>rich</i>	<i>poor</i>
<i>thick</i>	<i>thin</i>	<i>expensive</i>	<i>cheap</i>	<i>quiet</i>
<i>loud</i>	<i>clean</i>	<i>dirty</i>	<i>smart</i>	<i>stupid</i>
<i>interesting</i>	<i>boring</i>	<i>new</i>	<i>old-fashioned</i>	<i>safe</i>
<i>dangerous</i>	<i>healthy</i>	<i>sick</i>	<i>easy</i>	<i>difficult</i>
<i>right</i>	<i>wrong</i>	<i>high</i>	<i>low</i>	<i>near</i>
<i>far</i>	<i>early</i>	<i>late</i>	<i>wet</i>	<i>dry</i>
<i>busy</i>	<i>free</i>	<i>careful</i>	<i>careless</i>	<i>friendly</i>
<i>unfriendly</i>	<i>important</i>	<i>unimportant</i>	<i>famous</i>	<i>unknown</i>
<i>excited</i>	<i>calm</i>	<i>serious</i>	<i>funny</i>	<i>tired</i>
<i>energetic</i>				

969

970

971

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

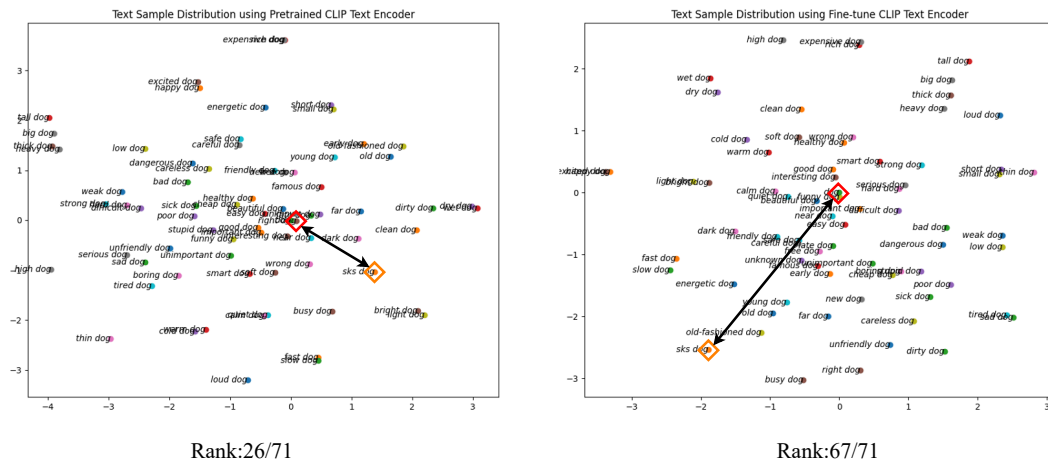


Figure 12: Visualization of the CLIP sample space. Using ChatGPT, we created 70 phrases containing adjectives related to the superclass of dogs. Subsequently, text features derived from these phrases were processed through the CLIP text encoder, downsampled, and their distance from the central point of the superclass (representing a dog image) was calculated. The comparison between the pre-trained model (illustrated in the left figure) and the fine-tuned model (depicted in the right figure) indicates that in the pre-trained model, phrases with special tokens ranked 26 out of 71, while in the fine-tuned model, they ranked 67 out of 71. Moreover, it is evident that phrases containing special tokens are situated further away from the central point of the superclass.

## E 2D TEXT SPACE DISTANCE CALCULATION

In this section, we discuss how we visualized CLIP’s text sample space in Fig. 3a. First, we collect 70 phrases that combine adjectives with words representing the class of given concepts (e.g., "a happy dog" or "a cool dog"). Using the CLIP text encoder, we extract text embeddings for these phrases. To visualize the semantic space and intuitively track modifications within it, we use t-SNE to reduce these high-dimensional embedding vectors to 2D. An overview of this result are shown in Fig.3(a), meanwhile we show the real experimental results in Fig.12. Formally, we use the adjectives generated which are described in Section. D as the initialize set  $S$ , and use the following pseudocode to get a 2D point set  $T$ :

---

### Algorithm 1 Algorithm to Convert Character Set to 2D Point Set

---

- 1: **Input:** Initial character set  $S$
  - 2: **Output:** 2D point set  $T$ , Distance set  $Dis$
  - 3:  $E \leftarrow \text{CLIP text encoder encoding}(S)$  ▷ Encode the character set to an encoding set
  - 4:  $T \leftarrow \text{TSNE}(E)$  ▷ Dimensionality reduction of the encoding set to a 2D point set
  - 5:  $Dis \leftarrow \{ \|T_i - T_{class}\| \mid i = 1, 2, \dots, |T| \}$  ▷ Calculate the 2D distance to  $T_{class}$  for each point in  $T$
  - 6: **return**  $T, Dis$
- 

## F MULTI-CONCEPTS EXPERIMENTS

In this section, we demonstrate the ability of ClassDiffusion to generate multiple concepts, specifically 3 concepts in one model. Fig 13 shows the result of this experiment. The experiments demonstrate that ClassDiffusion generates high-quality results when combining multiple concepts, validating the effectiveness of our proposed methods.

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

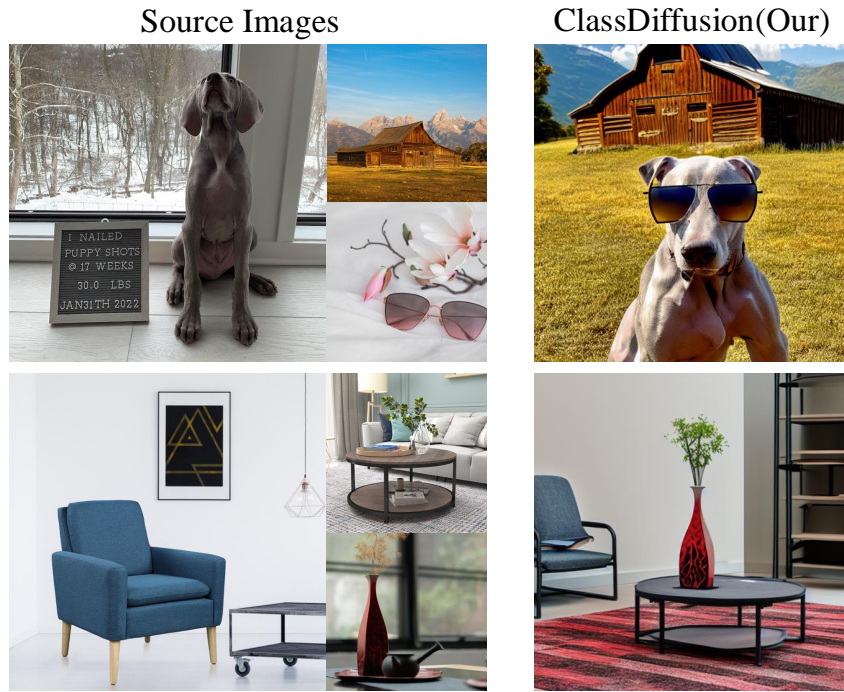


Figure 13: Qualitative result of generating three concepts.

## G MORE BASELINE COMPARISON

To further evaluate the performance of our proposed method, we further introduce a new baseline Prospect [92] which aims at solving similar problems we observe. The result of the experiment is shown in Tab. 2. The quantitative results below show that our model achieves superior performance.

Models	CLIP-T $\uparrow$	CLIP-I $\uparrow$	DINO-I $\uparrow$
Prospect	0.294	0.815	0.588
ClassDiffusion(Ours)	<b>0.300</b>	<b>0.828</b>	<b>0.673</b>

Table 2: Quantitative results comparing with Prospect

## H QUANTITATIVE ABLATION OF SPL WEIGHT

In this section, we conduct a quantitative ablation experiment on the choose of SPL weight. Tab. 3 shows the result of the experiment result. This result shows that CLIP-T becomes higher (increase in the ability to follow prompts) and DINO-I decreases(decrease in the ability to customize concepts) with increasing SPL weight, which is consistent with our expectations. Meanwhile, we find that SPL is a loss function that is insensitive to weight. Combined with the qualitative observation in Fig. 8, we prefer to choose 1 as the SPL weight.

## I DETAILS OF USER STUDY

We offer users a comprehensive user study guide that includes user selection criteria which is shown in Fig. 14. Additionally, to maintain fairness, we positioned our method alongside the baseline method randomly to prevent users from showing bias towards either method. Our percentages in Tab. 1 are obtained by calculating the number that chose our model better as a percentage of the overall number that made a preference.

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

## User Study

The website of User Study is: xxxx-xxxx.com

After entering the correct token, you will see an interaction similar to the following:

a bowl in the jungle



Left is better Right is better Equally good

The prompt of images is shown on the top, you should make decision by:

1. If one of the image is aligned with the prompt and another isn't, choose the align one.
2. If none or both of the images is aligned with the prompt, choose the one that is more aligned with the prompt.
3. If you can't make a choice (i.e. you think both diagrams are equally good/bad for the prompt) choose Equally good.

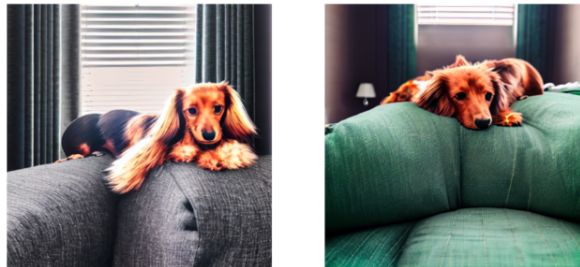
For the example image, you should choose "left is better". This is because although both the left and right images have bowls, the left image has a jungle and the right does not.

Or you will see an interaction similar to the following:

Reference Image



Generate Image



Left is better Right is better Equally good

You should choose the one that are more similar to the reference image or choose Equally good if you can not make a decision.

Figure 14: User Study Guide, which describes the user selection criteria and provides an example for reference.



1134  
1135  
1136  
1137  
1138  
1139  
1140

SPL weight	CLIP-T $\uparrow$	DINO-I $\uparrow$
0.01	0.299	0.677
0.1	0.300	0.674
1	0.300	0.673
10	0.300	0.665
100	0.301	0.661

1141  
1142  
1143

Table 3: Performance metrics for different SPL weights.

1144  
1145

## J MORE QUALITATIVE RESULT

1146  
1147  
1148

In Fig. 6, one generated image is provided for each prompt. Fig. 15 presents more images generated from the same prompts, thereby reinforcing the efficacy of our approach.

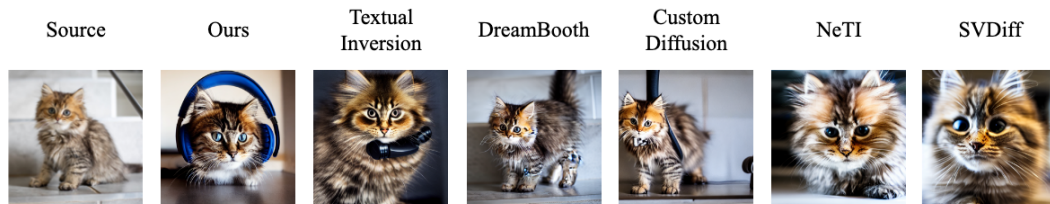
1149  
1150  
1151  
1152  
1153  
1154  
11551156  
1157  
1158  
1159  
1160  
1161A photo of a **cat** wearing a headphone1162  
1163  
1164  
1165  
1166A photo of a **dog** wearing chef's hat1167  
1168  
1169  
1170  
1171  
1172A photo of a flower in a **vase**1173  
1174  
1175  
1176  
1177  
1178A **teapot** on the hill of a mountain1179  
1180  
1181A **duck\_toy** wearing firefighter's outfit1182  
1183

Figure 15: Qualitative result of same prompts the main text.

1184  
1185

## K LIMITATION

1186  
1187

In this section, we discuss the limitations of our work. Our work are able to generate images that are aligned with the given prompt while keeping the features of the given concept. However, there are two major limitations in our work:



- 1188 • Considering that reconstruction of the human face is fine-grained, and the phrase “a photo of  
1189 a human” or “a photo of a human face” can not include extensive information about humans.  
1190 Whether our work can transfer to human-driven personalized generation remains explored.  
1191
- 1192 • For objects that have a combination of categories, choosing an appropriate center word  
1193 requires some experimentation.

## 1194 L SOCIAL IMPACT

1196 The advancements in text-to-image customization through fine-tuning diffusion models, as evidenced  
1197 by our work on ClassDiffusion, have significant social implications. By enhancing the compositional  
1198 capabilities of these models, our approach can contribute to a variety of fields, including digital  
1199 content creation, and education. In the realm of digital content creation, ClassDiffusion enables artists,  
1200 designers, and marketers to generate more precise and complex images based on textual descriptions.  
1201 This improvement reduces the time and effort required to produce customized visual content, fostering  
1202 creativity and innovation. It also allows for the seamless incorporation of personalized elements  
1203 into digital artworks, advertising materials, and user-generated content, thereby enhancing user  
1204 engagement and satisfaction. ClassDiffusion can also be a powerful tool in educational settings.  
1205 Educators can use this technology to create illustrative materials that are tailored to specific learning  
1206 objectives. For instance, teachers could generate images that accurately depict historical events,  
1207 scientific concepts, or literary scenes, making learning more interactive and engaging for students.  
1208 Furthermore, this technology can aid in the development of educational content for diverse learning  
1209 needs, including materials for students with disabilities.

1210 While the advancements in text-to-image generation hold promise, it is essential to address the ethical  
1211 considerations associated with their use. Ensuring that these models are free from biases and do not  
1212 perpetuate harmful stereotypes is crucial. Our work on ClassDiffusion includes measures to mitigate  
1213 semantic drift, which helps maintain the integrity and accuracy of generated content. Continuous  
1214 evaluation and updates are necessary to uphold these standards and ensure the technology benefits  
1215 society as a whole.

1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241