
Enhancing Detail Preservation for Customized Text-to-Image Generation: A Regularization-Free Approach

Yufan Zhou¹, Ruiyi Zhang², Tong Sun², Jinhui Xu¹

¹ State University of New York at Buffalo ² Adobe Research
{yufanzho, jinhui}@buffalo.edu, {ruizhang, tsun}@adobe.com

Abstract

Recent text-to-image generation models have demonstrated impressive capability of generating text-aligned images with high fidelity. However, generating images of novel concept provided by the user input image is still a challenging task. To address this problem, researchers have been exploring various methods for customizing pre-trained text-to-image generation models. Currently, most existing methods for customizing pre-trained text-to-image generation models involve the use of regularization techniques to prevent over-fitting. While regularization will ease the challenge of customization and leads to successful content creation with respect to text guidance, it may restrict the model capability, resulting in the loss of detailed information and inferior performance. In this work, we propose a novel framework for customized text-to-image generation without the use of regularization. Specifically, our proposed framework consists of an encoder network and a novel sampling method which can tackle the over-fitting problem without the use of regularization. With the proposed framework, we are able to customize a large-scale text-to-image generation model within half a minute on single GPU, with only one image provided by the user. We demonstrate in experiments that our proposed framework outperforms existing methods, and preserves more fine-grained details.

1 Introduction

Text-to-image generation is a research topic that has been explored for years [33, 36, 38, 39, 41, 42], with remarkable progresses recently. Nowadays, researchers are able to perform zero-shot text-to-image generation with arbitrary text input by training large-scale models on web-scale datasets. Starting from DALL-E [21] and CogView [5], numerous methods have been proposed [3, 6, 7, 20, 22, 24, 37, 40], leading to impressive capability in generating text-aligned images of high resolution with exceptional fidelity. Besides text-to-image generation, these large-scale models also have huge impacts on many other applications including image manipulation [1, 10] and video generation [11, 29].

Although aforementioned large-scale text-to-image generation models are able to perform text-aligned and creative generation, they may face difficulties in generating novel and unique concepts [8] specified by users. Thus, researchers have exploited different methods in customizing pre-trained text-to-image generation models. For instance, [17, 23] propose to fine-tune the pre-trained generative models with few samples, where different regularization methods are applied to prevent over-fitting. [8, 9, 34] propose to encode the novel concept of user input image in a word embedding, which is obtained by an optimization method or from an encoder network. All these methods lead to customized generation for the novel concept, while satisfying additional requirements described in arbitrary user input text.



Figure 1: Customized generation with the proposed framework. Given only single testing image, we are able to perform customized generation which satisfies arbitrary specified requirements and preserves fine-grained details.

Despite these progresses, recent research also makes us suspect that the use of regularization may potentially restrict the capability of customized generation, leading to the information loss of fine-grained details. In this paper, we propose a novel framework called *ProFusion*, which consists of an encoder called *PromptNet* and a novel sampling method called *Fusion Sampling*. Different from previous methods, our ProFusion does not require any regularization, the potential over-fitting problem can be tackled by the proposed Fusion Sampling method at inference, which saves training time as there is no need to tune the hyper-parameters for regularization method. Our main contributions can be summarized as follows:

- We propose ProFusion, a novel framework for customized generation. Given single testing image containing a unique concept, the proposed framework can generate customized output for the unique concept and meets additional requirement specified in arbitrary text. Only about 30 seconds of fine-tuning on single GPU is required;
- The proposed framework does not require any regularization method to prevent over-fitting, which significantly reduces training time as there is no need to tune regularization hyper-parameters. The absence of regularization also allows the proposed framework to achieve enhanced preservation of fine-grained details;
- Extensive results, including qualitative, quantitative and human evaluation results, have demonstrated the effectiveness of the proposed ProFusion. Ablation studies are also conducted to better understand the components in the proposed framework;

2 Methodology

We now present our proposed ProFusion framework, which consists of a neural network called PromptNet and a novel sampling method called Fusion Sampling. Specifically, PromptNet is an encoder network which can generate word embedding S^* conditioned on input image x , inside the input embedding space of the text encoder from Stable Diffusion 2. The major benefit of mapping x into S^* is that S^* can be readily combined with arbitrary text to construct prompt for creative generation, e.g., " S^* from a superhero movie screenshot"; Meanwhile, the Fusion Sampling is a sampling method leads to promising generation which meets the specified text requirements while maintaining fine-grained details of the input image x .

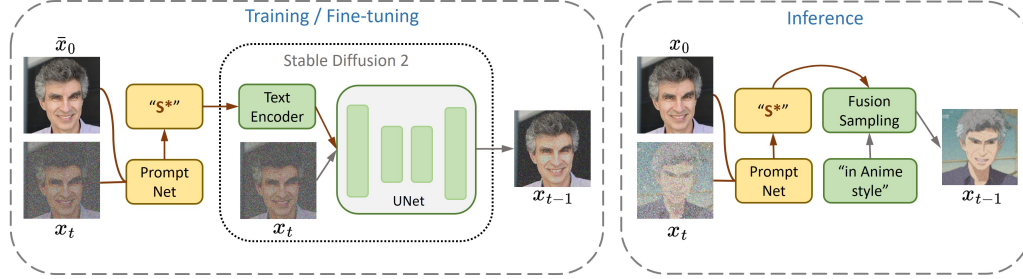


Figure 2: Illustration of the proposed framework.

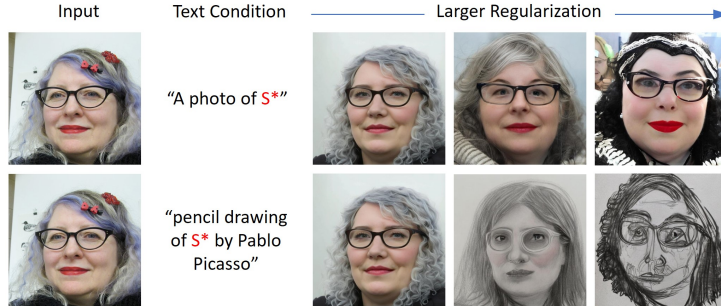


Figure 3: The performance of customized generation is impacted by the level of regularization.

Our core idea is presented in Figure 2. The proposed PromptNet infers S^* from an input image \mathbf{x}_0 and current noisy generation \mathbf{x}_t . Instead of using \mathbf{x}_0 , we can use $\bar{\mathbf{x}}_0$ during the training of PromptNet, which denotes a different view of \mathbf{x}_0 and can be obtained by data augmentation, *e.g.*, resizing, rotation. The PromptNet is trained with diffusion loss:

$$L_{\text{Diffusion}} = \mathbb{E}_{\mathbf{x}, \mathbf{y}(S^*), t, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{y}(S^*), t)\|_2^2], \quad (1)$$

where $\mathbf{y}(S^*)$ denotes the constructed prompt containing S^* , *e.g.* "A photo of S^* ".

Existing works [8, 9] use similar idea to obtain S^* . However, regularization are often applied in these works. For instance, E4T [9] proposes to use an encoder to generate S^* , which is optimized with

$$L = L_{\text{Diffusion}} + \lambda \|S^*\|_2^2, \quad (2)$$

where the L_2 norm of S^* is regularized. Similarly, Textual Inversion [8] proposes to directly obtain S^* by solving

$$S^* = \operatorname{argmin}_{S'} L_{\text{Diffusion}} + \lambda \|S' - S\|_2^2$$

with optimization method, where S denotes a coarse embedding*.

In this work, we argue that although the use of regularization will ease the challenge and enables successful content creation with respect to testing text. It also leads to the loss of detailed information, resulting in inferior performance. To verify this argument, we conduct a simple experiment on FFHQ dataset [15]. We train several encoders with different levels of regularization by selecting different λ in (2). After training, we test their capability by classifier-free sampling [13] with different prompts containing resulting S^* . The results are shown in Figure 3, from which we can find that smaller regularization leads to less information loss, which results in better preservation of details. However, the information could be too strong to prevent creative generation with respect to user input text. Meanwhile, large regularization leads to successful content creation, while fails to capture details of the input image, resulting in unsatisfactory results.

A consequent question is, **is it possible to perform successful customized generation using S^* obtained without regularization so that the details from original image can be well-preserved?** To answer this question, we propose a novel sampling method called Fusion Sampling.

let S^ be a target embedding for a specific human face image, S can be set to be the embedding of text "face".

2.1 Fusion Sampling

Given a PromptNet pre-trained without regularization which can map input image \mathbf{x}_0 into word embedding S^* , our goal is to successfully perform customized generation which preserves details of \mathbf{x}_0 , and meets the requirements specified in arbitrary prompt containing S^* .

The task can be formulated as a conditional generation task with conditions S^* and C , where C denotes arbitrary user input text. We start from the most commonly used classifier-free sampling [13]. To sample \mathbf{x}_{t-1} given current noisy sample \mathbf{x}_t and conditions $[S^*, C]$, the diffusion model first outputs the predictions of conditional noise $\epsilon_\theta(\mathbf{x}_t, S^*, C)$ and unconditional noise $\epsilon_\theta(\mathbf{x}_t)$. Then an updated prediction (with hyper-parameter ω)

$$\tilde{\epsilon}_\theta(\mathbf{x}_t, S^*, C) = (1 + \omega)\epsilon_\theta(\mathbf{x}_t, S^*, C) - \omega\epsilon_\theta(\mathbf{x}_t), \quad (3)$$

will be used in different sampling strategies [12, 14, 30, 31].

In customized generation, the reason that vanilla classifier-free sampling does not work without regularization is that, information from S^* can become too strong without regularization. As a result, $\epsilon_\theta(\mathbf{x}_t, S^*, C)$ will degenerate to $\epsilon_\theta(\mathbf{x}_t, S^*)$ and information of C will be lost. Thus, we need to propose a new sampling method, to produce a new prediction for $\tilde{\epsilon}_\theta(\mathbf{x}_t, S^*, C)$ which is enforced to be conditioned on both S^* and C .

Sampling with independent conditions We begin by assuming that S^* and C are independent. According to [13], we know that

$$\epsilon_\theta(\mathbf{x}_t, S^*, C) = -\sqrt{1 - \bar{\alpha}_t} \nabla \log p(\mathbf{x}_t | S^*, C), \quad (4)$$

where $\bar{\alpha}_t$ is a hyper-parameter as defined in [12]. By (4) and Bayes' Rule, we can re-write (3) as

$$\tilde{\epsilon}_\theta(\mathbf{x}_t, S^*, C) = \epsilon_\theta(\mathbf{x}_t) - (1 + \omega)\sqrt{1 - \bar{\alpha}_t} \nabla \log p(S^*, C | \mathbf{x}_t). \quad (5)$$

Since we assume that S^*, C are independent, we can further re-write the above as

$$\begin{aligned} \tilde{\epsilon}_\theta(\mathbf{x}_t, S^*, C) &= \epsilon_\theta(\mathbf{x}_t) - (1 + \omega)\sqrt{1 - \bar{\alpha}_t} \nabla \log p(S^* | \mathbf{x}_t) - (1 + \omega)\sqrt{1 - \bar{\alpha}_t} \nabla \log p(C | \mathbf{x}_t) \\ &= \epsilon_\theta(\mathbf{x}_t) + (1 + \omega)\{\epsilon_\theta(\mathbf{x}_t, S^*) - \epsilon_\theta(\mathbf{x}_t)\} + (1 + \omega)\{\epsilon_\theta(\mathbf{x}_t, C) - \epsilon_\theta(\mathbf{x}_t)\}. \end{aligned}$$

We re-write it as

$$\tilde{\epsilon}_\theta(\mathbf{x}_t, S^*, C) = \epsilon_\theta(\mathbf{x}_t) + (1 + \omega_1)\{\epsilon_\theta(\mathbf{x}_t, S^*) - \epsilon_\theta(\mathbf{x}_t)\} + (1 + \omega_2)\{\epsilon_\theta(\mathbf{x}_t, C) - \epsilon_\theta(\mathbf{x}_t)\} \quad (6)$$

for more flexibility. (6) can be readily extended to more complicated scenarios, where a list of conditions $\{S_1^*, S_2^*, \dots, S_k^*, C\}$ are provided. The corresponding $\tilde{\epsilon}_\theta(\mathbf{x}_t, \{S_i^*\}_{i=1}^k, C)$ is

$$\tilde{\epsilon}_\theta(\mathbf{x}_t, \{S_i^*\}_{i=1}^k, C) = \epsilon_\theta(\mathbf{x}_t) + \sum_{i=1}^k (1 + \omega_i)\{\epsilon_\theta(\mathbf{x}_t, S_i^*) - \epsilon_\theta(\mathbf{x}_t)\} + (1 + \omega_C)\{\epsilon_\theta(\mathbf{x}_t, C) - \epsilon_\theta(\mathbf{x}_t)\}.$$

Fusion Sampling with dependent conditions One major drawback of (6) is that the independence does not always hold in practice. As we will show in later experiment, assuming S^* and C to be independent can lead to inferior generation.

To solve this problem, we propose Fusion Sampling, which consists of two stages at each timestep t : a **fusion stage** which encodes information from both S^* and C into \mathbf{x}_t with an updated $\tilde{\mathbf{x}}_t$, and a **refinement stage** which predicts \mathbf{x}_{t-1} based on Equation (6). The proposed algorithm is presented in Algorithm 1. Sampling with independent conditions can be regarded as a special case of Fusion Sampling with $m = 0$. In practice, $m = 1$ works well, thus we set $m = 1$ in all our experiments.

The remaining challenge in Algorithm 1 is to sample $\tilde{\mathbf{x}}_{t-1} \sim q(\tilde{\mathbf{x}}_{t-1} | \tilde{\mathbf{x}}_t, \tilde{\mathbf{x}}_0)$ and $\tilde{\mathbf{x}}_t \sim q(\tilde{\mathbf{x}}_t | \tilde{\mathbf{x}}_{t-1}, \tilde{\mathbf{x}}_0)$. We take Denoising Diffusion Implicit Models (DDIM) [30] as an example, while the following derivation can be extended to other diffusion models. Let \mathbf{I} be the identity matrix, σ_t denotes a hyper-parameter controlling randomness. In DDIM, we have

$$q(\tilde{\mathbf{x}}_t | \tilde{\mathbf{x}}_0) = \mathcal{N}(\tilde{\mathbf{x}}_t; \sqrt{\bar{\alpha}_t} \tilde{\mathbf{x}}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \quad (7)$$

and

$$q(\tilde{\mathbf{x}}_{t-1} | \tilde{\mathbf{x}}_t, \tilde{\mathbf{x}}_0) = \mathcal{N}(\tilde{\mathbf{x}}_{t-1}; \sqrt{\bar{\alpha}_{t-1}} \tilde{\mathbf{x}}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \frac{\tilde{\mathbf{x}}_t - \sqrt{\bar{\alpha}_t} \tilde{\mathbf{x}}_0}{\sqrt{1 - \bar{\alpha}_t}}, \sigma_t^2 \mathbf{I}). \quad (8)$$

Algorithm 1 Fusion Sampling at Timestep t

- 1: **Require: Conditions** S^* **and** C , **a noisy sample** \mathbf{x}_t , **a pre-trained diffusion model** ϵ_θ , **hyper-parameters** $0 < \sigma_t, 0 \leq \gamma \leq 1$.
 - 2: Set $\tilde{\mathbf{x}}_t = \mathbf{x}_t$
 - 3: // Fusion Stage
 - 4: **for** $i = 1, \dots, m$ **do**
 - 5: Generate $\tilde{\epsilon}_\theta(\tilde{\mathbf{x}}_t, \gamma S^*, C)$ by (3).
 - 6: Generate predicted sample $\tilde{\mathbf{x}}_0 = \frac{\tilde{\mathbf{x}}_t - \sqrt{1 - \bar{\alpha}_t} \tilde{\epsilon}_\theta(\tilde{\mathbf{x}}_t, \gamma S^*, C)}{\sqrt{\bar{\alpha}_t}}$.
 - 7: Inject fused information into $\tilde{\mathbf{x}}_{t-1}$ by sampling $\tilde{\mathbf{x}}_{t-1} \sim q(\tilde{\mathbf{x}}_{t-1} | \tilde{\mathbf{x}}_t, \tilde{\mathbf{x}}_0)$.
 - 8: **if** Use refinement stage **then**
 - 9: Inject fused information into $\tilde{\mathbf{x}}_t$ by sampling $\tilde{\mathbf{x}}_t \sim q(\tilde{\mathbf{x}}_t | \tilde{\mathbf{x}}_{t-1}, \tilde{\mathbf{x}}_0)$.
 - 10: **else**
 - 11: Return $\mathbf{x}_{t-1} = \tilde{\mathbf{x}}_{t-1}$.
 - 12: **end if**
 - 13: **end for**
 - 14: // Refinement Stage
 - 15: **if** Use refinement stage **then**
 - 16: Generate $\tilde{\epsilon}_\theta(\tilde{\mathbf{x}}_t, S^*, C)$ by (6) and perform classifier-free sampling step. Return \mathbf{x}_{t-1} .
 - 17: **end if**
-

By the property of Gaussian distributions [2], we know that

$$q(\tilde{\mathbf{x}}_t | \tilde{\mathbf{x}}_{t-1}, \tilde{\mathbf{x}}_0) = \mathcal{N}(\tilde{\mathbf{x}}_t; \Sigma(A^T L(\tilde{\mathbf{x}}_{t-1} - b) + B\boldsymbol{\mu}), \Sigma) \quad (9)$$

where

$$\Sigma = \frac{(1 - \bar{\alpha}_t)\sigma_t^2}{1 - \bar{\alpha}_{t-1}} \mathbf{I}, \quad \boldsymbol{\mu} = \sqrt{\bar{\alpha}_t} \tilde{\mathbf{x}}_0, \quad b = \sqrt{\bar{\alpha}_{t-1}} \tilde{\mathbf{x}}_0 - \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1} - \sigma_t^2)}{\sqrt{1 - \bar{\alpha}_t}} \tilde{\mathbf{x}}_0$$

$$A = \frac{\sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2}}{\sqrt{1 - \bar{\alpha}_t}} \mathbf{I}, \quad L = \frac{1}{\sigma_t^2} \mathbf{I}, \quad B = \frac{1}{1 - \bar{\alpha}_t} \mathbf{I}$$

which leads to

$$\tilde{\mathbf{x}}_t = \frac{\sqrt{(1 - \bar{\alpha}_t)(1 - \bar{\alpha}_{t-1} - \sigma_t^2)}}{1 - \bar{\alpha}_{t-1}} \tilde{\mathbf{x}}_{t-1} + \frac{(1 - \bar{\alpha}_t)\sigma_t^2}{1 - \bar{\alpha}_{t-1}} \mathbf{z}$$

$$+ \frac{\tilde{\mathbf{x}}_0}{1 - \bar{\alpha}_{t-1}} \{ \sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1}) - \sqrt{\bar{\alpha}_{t-1}(1 - \bar{\alpha}_t)(1 - \bar{\alpha}_{t-1} - \sigma_t^2)} \}, \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (10)$$

With further derivation, we can summarize a single update in fusion stage as:

$$\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_t - \frac{\sigma_t^2 \sqrt{1 - \bar{\alpha}_t}}{1 - \bar{\alpha}_{t-1}} \tilde{\epsilon}_\theta(\tilde{\mathbf{x}}_t, \gamma S^*, C) + \frac{\sqrt{(1 - \bar{\alpha}_t)(2 - 2\bar{\alpha}_{t-1} - \sigma_t^2)}}{1 - \bar{\alpha}_{t-1}} \sigma_t \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (11)$$

Remark 1 Recall $\tilde{\epsilon}_\theta(\tilde{\mathbf{x}}_t, \gamma S^*, C) = -\sqrt{1 - \bar{\alpha}_t} \nabla \log \tilde{p}_\omega(\tilde{\mathbf{x}}_t | \gamma S^*, C)$ [13], we can re-write (11) as

$$\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_t + \frac{\sigma_t^2(1 - \bar{\alpha}_t)}{1 - \bar{\alpha}_{t-1}} \nabla \log \tilde{p}_\omega(\tilde{\mathbf{x}}_t | \gamma S^*, C) + \frac{\sqrt{(1 - \bar{\alpha}_t)(2 - 2\bar{\alpha}_{t-1} - \sigma_t^2)}}{1 - \bar{\alpha}_{t-1}} \sigma_t \mathbf{z}. \quad (12)$$

From (12), we can conclude that our fusion stage is actually an gradient-based optimization method similar to Langevin dynamics [35]. Compared to Langevin dynamics which is

$$\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_t + \lambda \nabla \log \tilde{p}_\omega(\tilde{\mathbf{x}}_t | \gamma S^*, C) + \sqrt{2\lambda} \mathbf{z}. \quad (13)$$

with λ being the step size, (12) has less randomness, because

$$\frac{(1 - \bar{\alpha}_t)(2 - 2\bar{\alpha}_{t-1} - \sigma_t^2)\sigma_t^2}{(1 - \bar{\alpha}_{t-1})^2} \leq \frac{2\sigma_t^2(1 - \bar{\alpha}_t)}{1 - \bar{\alpha}_{t-1}}.$$

Remark 2 If we set the DDIM hyper-parameter to be $\sigma_t = \sqrt{1 - \bar{\alpha}_{t-1}}$, then (11) becomes

$$\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_t - \sqrt{1 - \bar{\alpha}_t} \tilde{\epsilon}_\theta(\tilde{\mathbf{x}}_t, \gamma S^*, C) + \sqrt{1 - \bar{\alpha}_t} \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

which is equivalent to sampling $\tilde{\mathbf{x}}_t$ using (7) without sampling intermediate $\tilde{\mathbf{x}}_{t-1}$ in our Algorithm 1. Thus directly sampling $\tilde{\mathbf{x}}_t$ using (7) is a special case of our Fusion Sampling algorithm.

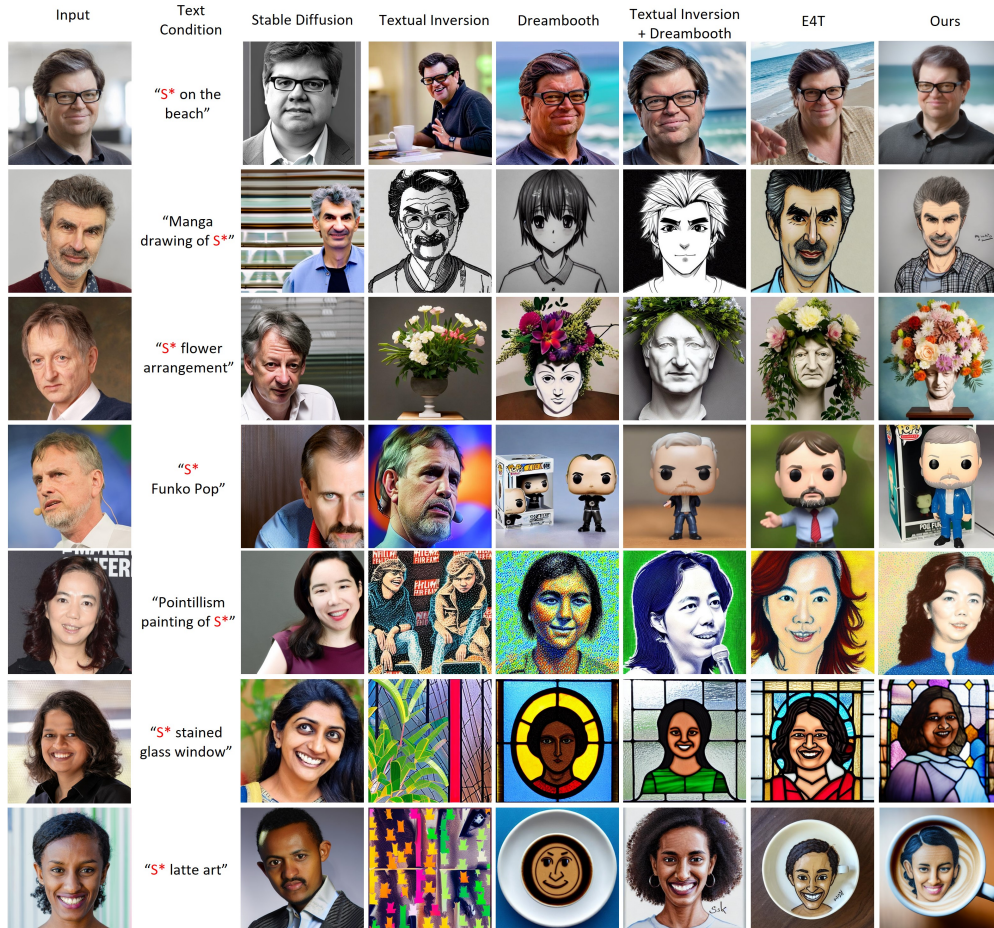


Figure 4: Comparison with baseline methods. Our proposed approach exhibits superior capability for preserving fine-grained details.

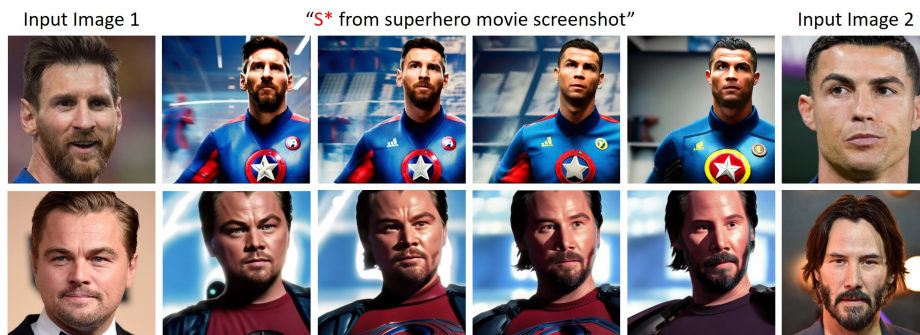


Figure 5: The proposed framework enables generation conditioned on multiple input images and text. Creative interpolation can be performed.

3 Experiments

We conduct extensive experiments to evaluate the proposed framework. Specifically, we first pre-train a PromptNet on FFHQ dataset [15] on 8 NVIDIA A100 GPUs for 80,000 iterations with a batch size of 64, without any data augmentation. Given a testing image, the PromptNet and all attention layers of the pre-trained Stable Diffusion 2 are fine-tuned for 50 steps with a batch size of 8. Only half a minute and a single GPU is required in fine-tuning such a customized generative model, indicating the efficiency of the proposed method, especially considering the impressive results we could obtain. Some more implementation details are provided in the Appendix. Our code and pre-trained models will be publicly available at <https://github.com/drboog/ProFusion>.

Method	Pre-trained CLIP Models								
	ViT-B/32	ViT-B/16	ViT-L/14	ViT-L/14@336px	RN101	RN50	RN50×4	RN50×16	RN50×64
Stable Diffusion 2	0.271	0.256	0.196	0.196	0.428	0.202	0.355	0.254	0.181
Textual Inversion	0.257	0.251	0.197	0.201	0.426	0.195	0.350	0.247	0.173
DreamBooth	0.283	0.267	0.205	0.210	0.434	0.209	0.363	0.260	0.187
E4T	0.277	0.264	0.203	0.213	0.429	0.206	0.358	0.260	0.191
ProFusion (Ours)	0.293	0.283	0.225	0.229	0.446	0.223	0.374	0.279	0.202

Table 1: Similarity (\uparrow) between generated example and input text.

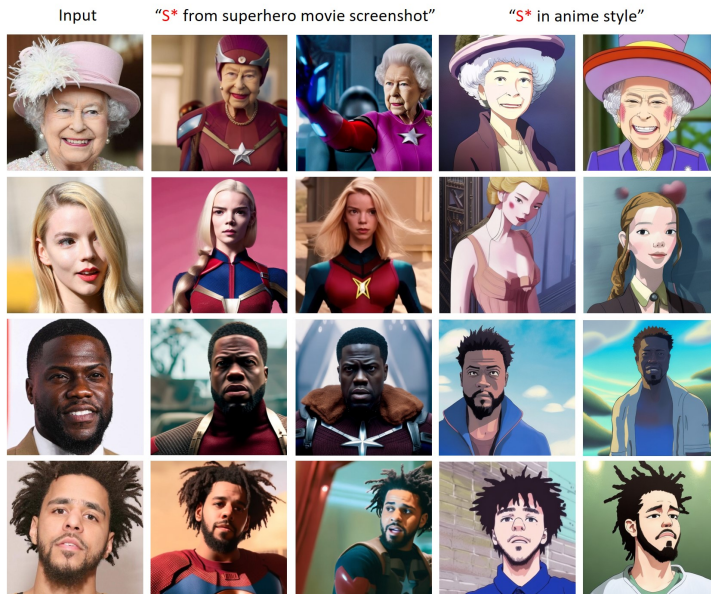


Figure 6: Some results of customized generation with the proposed framework.

3.1 Qualitative Results

Our main results are shown in Figure 1 and Figure 6. From the results, we can see that the proposed framework effectively achieves customized generation which meets the specified text requirements while maintaining fine-grained details of the input image. More results are provided in the Appendix. As mentioned previously, our proposed framework is also able to perform generation conditioned on multiple images. We also provide these generated examples in Figure 5.

Following [9], we then compare proposed framework with several baseline methods including Stable Diffusion[†] [22], Textual Inversion [8], DreamBooth [23], E4T [9]. The qualitative results are presented in Figure 4, where the results of related methods are directly taken from [9]. From the comparison we can see that our framework results in better preservation of fine-grained details.

3.2 Quantitative Results

We also evaluate our methods and baseline methods quantitatively. Specifically, we utilize different pre-trained CLIP models [19] to calculate the image-prompt similarity between the generated image and input text. The results are shown in Table 1, our ProFusion obtains higher image-prompt similarity on all CLIP models, indicating better prompt-adherence and edit-ability..

We then calculate the identity similarity between the generated image and input image, which is cosine similarity computed using features extracted by pre-trained face recognition models. The identity similarity is also evaluated across different pre-trained models [4, 16, 18, 25, 26, 27, 28, 32]. The results are shown in Table 2. In general, our ProFusion obtains higher similarity, indicating better identity preservation.

[†]The results of Stable Diffusion is obtained by directly feeding corresponding researcher’s name and text requirements into the pre-trained text-to-image generation model.

Method	Pre-trained Face Recognition Models							
	VGG-Face	Facenet	Facenet512	OpenFace	DeepFace	ArcFace	SFace	AdaFace
Stable Diffusion 2	0.530	0.334	0.323	0.497	0.641	0.144	0.191	0.093
Textual Inversion	0.516	0.410	0.372	0.566	0.651	0.248	0.231	0.210
DreamBooth	0.518	0.483	0.415	0.516	0.643	0.379	0.304	0.307
E4T	0.677	0.596	0.621	0.660	0.732	0.454	0.398	0.426
ProFusion (Ours)	0.720	0.616	0.597	0.681	0.774	0.459	0.443	0.432

Table 2: Similarity (\uparrow) between generated example and input image.

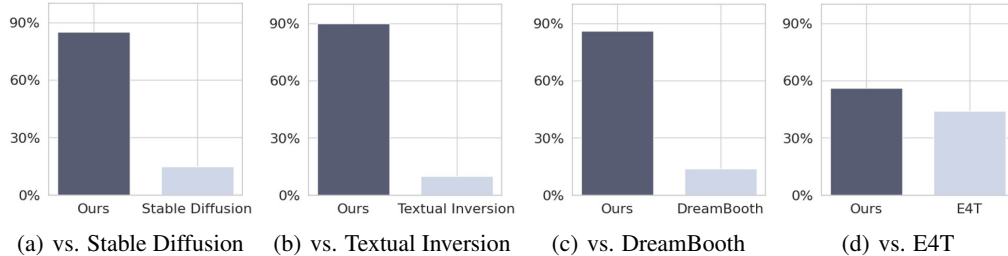


Figure 7: Results of human evaluation.

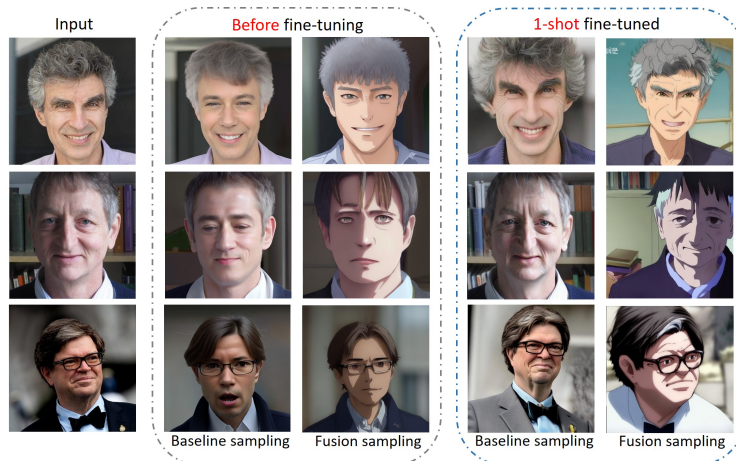


Figure 8: Examples with prompt " S^* in anime style", Fusion Sampling outperforms baseline.

3.3 Human Evaluation

We then conduct human evaluation on Amazon Mechanical Turk (MTurk). The workers are presented with two generated images from different methods along with original image and text requirements. They are then tasked with indicating their preferred choice. More details are provided in the Appendix. The results are shown in Figure 7, where we can find that our method obtains a higher preference rate compared to all other methods, indicating the effectiveness of our proposed framework.

3.4 Ablation Study

We conduct several ablation studies to further investigate the proposed ProFusion.

Fusion Sampling First of all, we apply the proposed Fusion Sampling with both pre-trained and fine-tuned PromptNet. As shown in Figure 8, Fusion Sampling obtains better results on both pre-trained and fine-tuned models compared to baseline classifier-free sampling. We then investigate the effects of removing fusion stage or refinement stage in the proposed Fusion Sampling. As we can see from Figure 10, removing refinement stage leads to the loss in detailed information, while removing fusion stage leads to a generated image with disorganized structure. Intuitively, S^* , which is the output of PromptNet, tries to generate a human face image following the structural information from the original image, while the text "is wearing superman costume" aims to generate a half-length photo. The conflicting nature of these two conditions results in an undesirable generation with a disorganized structure after we remove the fusion stage.

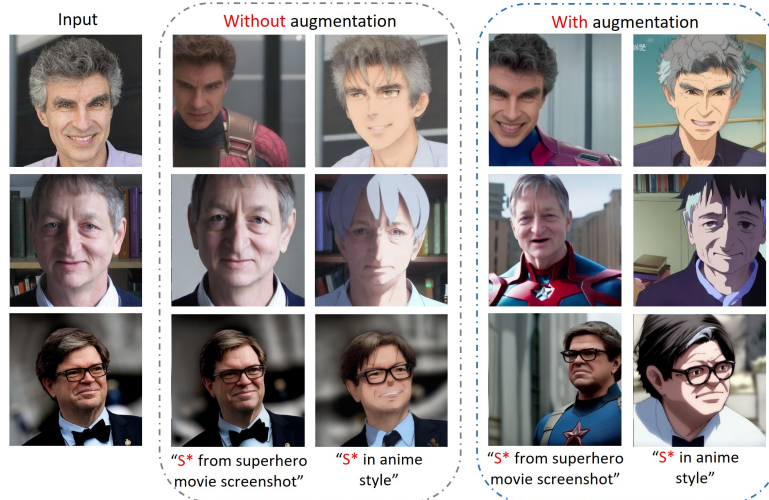


Figure 9: Data augmentation in fine-tuning stage leads to performance improvement.



Figure 10: Generated examples of ablation study, with prompt " S^* is wearing superman costume".

Data Augmentation We then analyze the effects of data augmentation. In particular, we conduct separate fine-tuning experiments: one with data augmentation and one without, both models are tested with Fusion Sampling after fine-tuning. The results are shown in Figure 9, we observe an improvement in performance as a result of employing data augmentation. Our data augmentation strategy is presented in the Appendix.

4 Discussion

Although the proposed framework has demonstrated remarkable capability in achieving high-quality customized generation, there are areas that can be improved. For instance, although ProFusion can reduce the training time by only requiring a single training without the need of tuning regularization hyperparameters, the proposed Fusion Sampling actually results in an increased inference time. This is due to the division of each sampling step into two stages. In the future, we would like to explore ways to improve the efficiency of Fusion Sampling.

Similar to other related works, our framework utilizing large-scale text-to-image generation models can raise ethical implications, both positive and negative. On the one hand, customized generation can create images with sensitive information and spread misinformation; On the other hand, it also holds the potential to minimize model biases as discussed in [8, 9]. Thus it is crucial to exercise proper supervision when implementing these methods in real-world applications.

5 Conclusion

In this paper, we present ProFusion, a novel framework for customized generation. Different from related methods which employs regularization, ProFusion successfully performs customized generation without any regularization, thus exhibits superior capability for preserving fine-grained details with less training time. Extensive experiments have demonstrated the effectiveness of the proposed ProFusion.

References

- [1] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18208–18218, 2022.
- [2] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- [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.
- [4] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4690–4699, 2019.
- [5] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, and Jie Tang. Cogview: Mastering text-to-image generation via transformers, 2021.
- [6] Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. Cogview2: Faster and better text-to-image generation via hierarchical transformers. *arXiv preprint arXiv:2204.14217*, 2022.
- [7] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. *arXiv preprint arXiv:2203.13131*, 2022.
- [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. *arXiv preprint arXiv:2208.01618*, 2022.
- [9] Rinon Gal, Moab Arar, Yuval Atzmon, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. Designing an encoder for fast personalization of text-to-image models. *arXiv preprint arXiv:2302.12228*, 2023.
- [10] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022.
- [11] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [13] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.
- [14] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *Advances in Neural Information Processing Systems*.
- [15] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019.
- [16] Minchul Kim, Anil K Jain, and Xiaoming Liu. Adaface: Quality adaptive margin for face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18750–18759, 2022.
- [17] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. *arXiv preprint arXiv:2212.04488*, 2022.
- [18] Omkar M Parkhi, Andrea Vedaldi, and Andrew Zisserman. Deep face recognition. 2015.
- [19] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021.
- [20] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.

- [21] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8821–8831. PMLR, 18–24 Jul 2021.
- [22] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. *arXiv preprint arXiv:2112.10752*, 2021.
- [23] 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. 2022.
- [24] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.
- [25] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015.
- [26] Sefik Ilkin Serengil and Alper Ozpinar. Lightface: A hybrid deep face recognition framework. In *2020 Innovations in Intelligent Systems and Applications Conference (ASYU)*, pages 23–27. IEEE, 2020.
- [27] Sefik Ilkin Serengil and Alper Ozpinar. Hyperextended lightface: A facial attribute analysis framework. In *2021 International Conference on Engineering and Emerging Technologies (ICEET)*, pages 1–4. IEEE, 2021.
- [28] Sefik Ilkin Serengil and Alper Ozpinar. An evaluation of sql and nosql databases for facial recognition pipelines. <https://www.cambridge.org/engage/coe/article-details/63f3e5541d2d184063d4f569>, 2023. Preprint.
- [29] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022.
- [30] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2020.
- [31] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*.
- [32] Yaniv Taigman, Ming Yang, Marc’Aurelio Ranzato, and Lior Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1701–1708, 2014.
- [33] Ming Tao, Hao Tang, Songsong Wu, Nicu Sebe, Xiao-Yuan Jing, Fei Wu, and Bingkun Bao. Df-gan: Deep fusion generative adversarial networks for text-to-image synthesis, 2021.
- [34] Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. Elite: Encoding visual concepts into textual embeddings for customized text-to-image generation. *arXiv preprint arXiv:2302.13848*, 2023.
- [35] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 681–688, 2011.
- [36] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1316–1324, 2018.
- [37] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022.
- [38] Han Zhang, Jing Yu Koh, Jason Baldridge, Honglak Lee, and Yinfei Yang. Cross-modal contrastive learning for text-to-image generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 833–842, 2021.

- [39] Yufan Zhou, Chunyuan Li, Changyou Chen, Jianfeng Gao, and Jinhui Xu. Lafite2: Few-shot text-to-image generation. *ArXiv*, abs/2210.14124, 2022.
- [40] Yufan Zhou, Bingchen Liu, Yizhe Zhu, Xiao Yang, Changyou Chen, and Jinhui Xu. Shifted diffusion for text-to-image generation. *arXiv preprint arXiv:2211.15388*, 2022.
- [41] Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong Yu, Jiuxiang Gu, Jinhui Xu, and Tong Sun. Lafite: Towards language-free training for text-to-image generation. *arXiv preprint arXiv:2111.13792*, 2021.
- [42] Minfeng Zhu, Pingbo Pan, Wei Chen, and Yi Yang. Dm-gan: Dynamic memory generative adversarial networks for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5802–5810, 2019.

A More Generated examples

Some more generated examples are provided in Figure 11, Figure 12 and Figure 13.



Figure 11: More results of customized generation with the proposed framework.

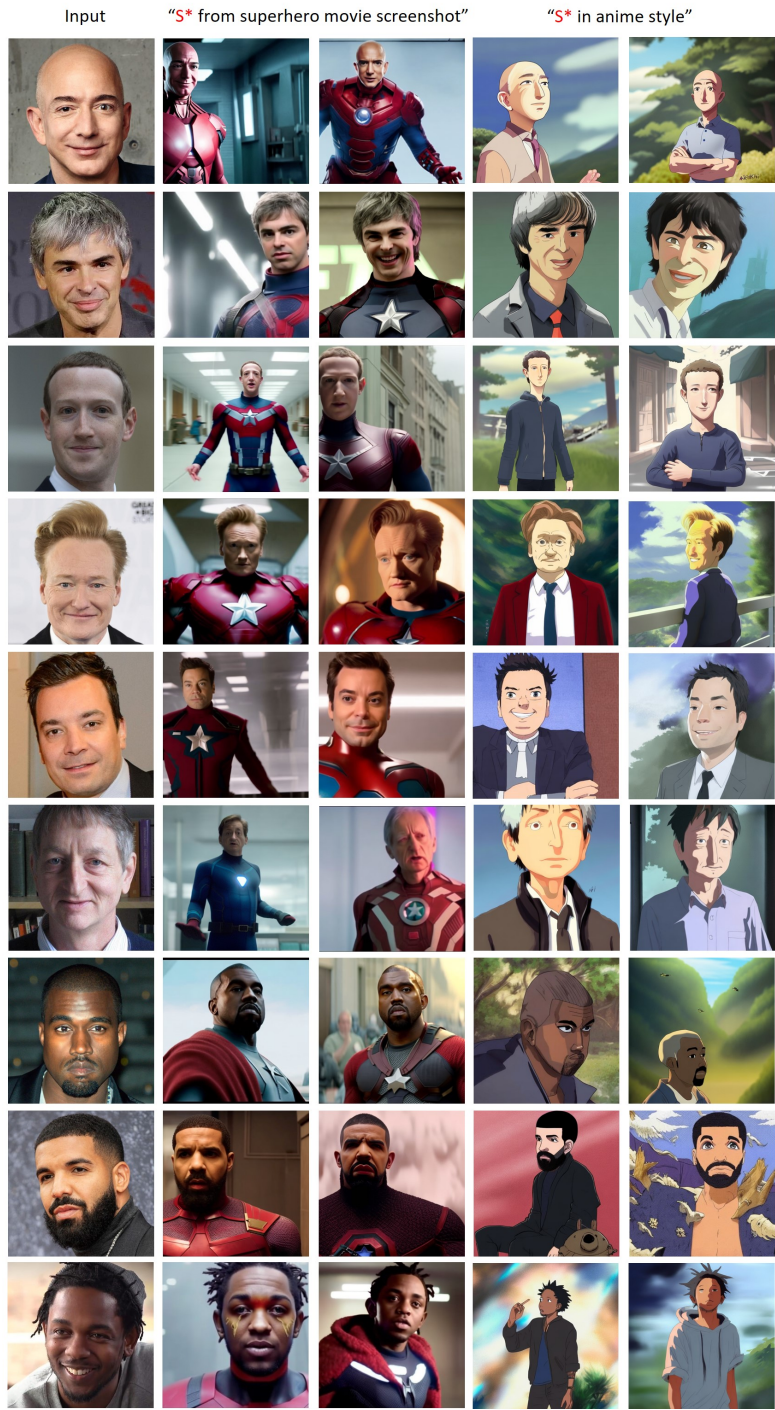


Figure 12: More results of customized generation with the proposed framework.

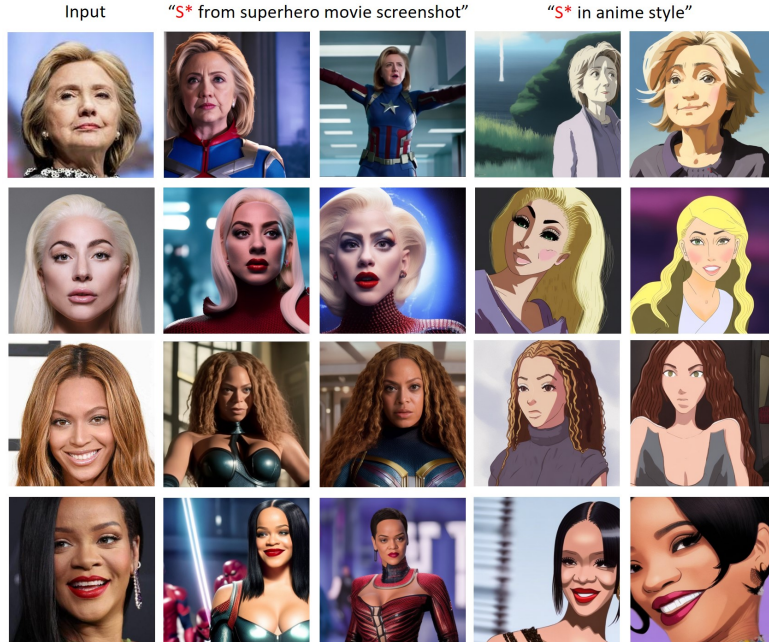


Figure 13: More results of customized generation with the proposed framework.

B Experiment Details

We provide some experiment details in this section.

Data augmentation We implement data augmentation at the fine-tuning stage, which is illustrated in Figure 14. Given a testing image, we first create a masked image where only the target face/object is unmasked. The masked image will be fed into a pre-trained Stable Diffusion inpainting model after random resize and rotation. The inpainting model will generate multiple augmented images, with different background. We use a positive prompt "a photo of a man/woman, highly detailed, soft natural lighting, photo realism, professional portrait, ultra-detailed, 4k resolution, wallpaper, hd wallpaper." and a negative prompt "magazine, frame, tiled, repeated, multiple people, multiple faces, group of people, split frame, multiple panel, split image, watermark, boarder, diptych, triptych" with a classifier-free guidance of 7.5 during inpainting.

PromptNet Our PromptNet is an encoder contains 5 encoder blocks, which are similar to the downsize and middle blocks in Stable Diffusion 2. The parameters are initialized with value from pre-trained Stable Diffusion 2 when applicable. Different from the blocks in Stable Diffusion 2, we use image embeddings from pre-trained CLIP ViT-H/14 instead of text embeddings as the contents for cross attention layers. The inputs \bar{x}_0 and x_t are first processed by different convolution layers, whose outputs are summed to serve as the input for the following blocks.

Human Evaluation Due to the fact that we do not have official implementation and pre-trained models of E4T [9], we directly take some generated examples from their paper for fair comparison. Then we use corresponding prompts in our framework to generate images to be compared. Specifically, there are 39 source image and prompt pairs for five different methods and each generated image is evaluated by five different workers with expertise. These workers are all from the US and required to have performed at least 10,000 approved assignments with an approval rate $\geq 98\%$. The human evaluation user interface is shown in Figure 7.

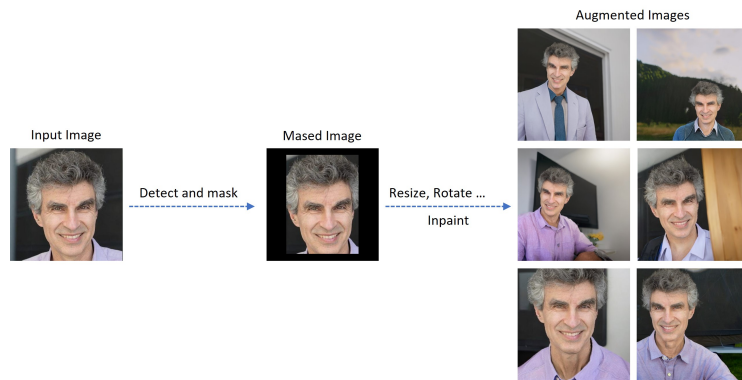



Figure 14: Illustration of data augmentation in our fine-tuning stage.



Read the text below and use the sliders below indicate your choice.

Input Image:



Text: Yoshua Bengio in a colorful graffiti

Generated image (1): **Generated image (2):**

Which image (1 or 2) is a better editing given the text?

Please consider the following factors when you choose the better result.

By "identity", whether the two persons are similar in terms of detailed features (expression, demeanor, and action) regardless of the styles.

By "generation quality", identify how well the it matches with the input texts (Please search the keywords for better understanding). Visually appealing is a minor factor.

Figure 15: Human Evaluation User Interface.