

JeDi: Joint-Image Diffusion Models for Finetuning-Free Personalized Text-to-Image Generation

Yu Zeng^{1,3} Vishal M. Patel¹ Haochen Wang² Xun Huang³
 Ting-Chun Wang³ Ming-Yu Liu³ Yogesh Balaji³
¹Johns Hopkins University ²TTI-Chicago ³NVIDIA Research

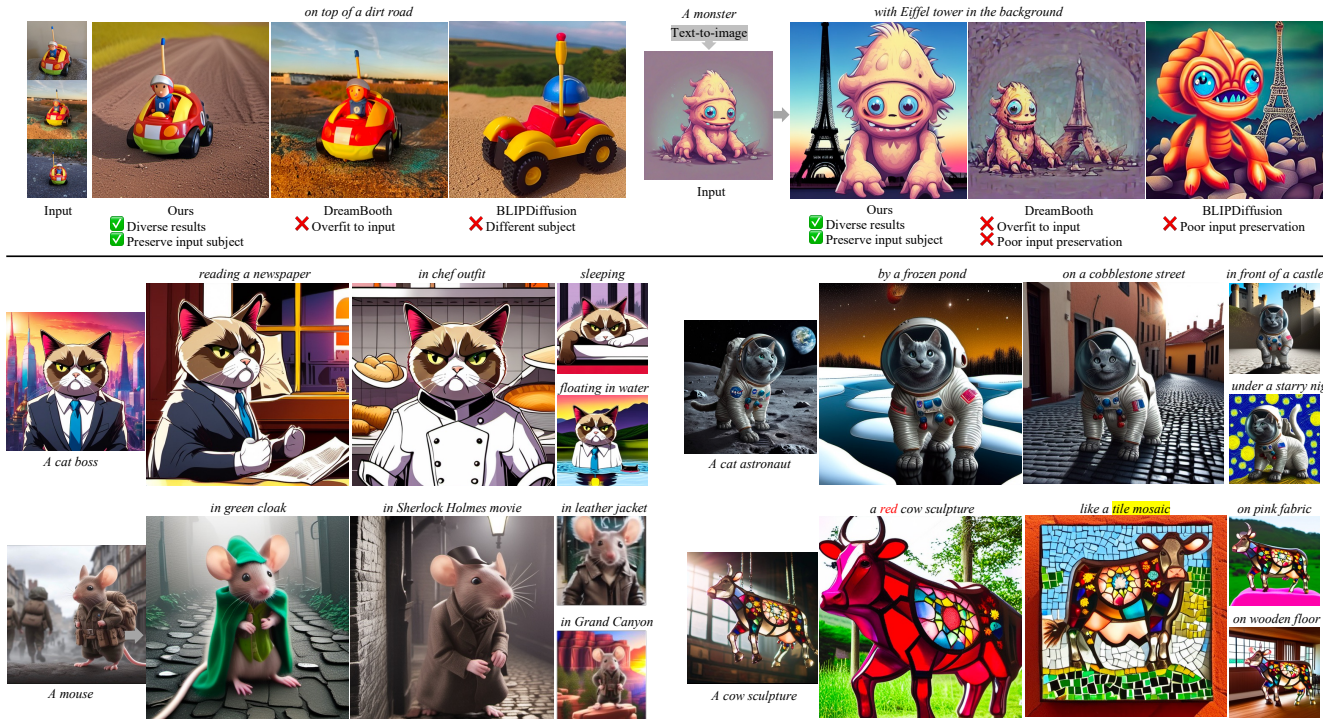


Figure 1. We present Joint-Image Diffusion (JeDi), a finetuning-free image personalization model that can operate on any number of reference images. JeDi is able to preserve the appearance of custom subjects while generating novel variations. As shown in the top row, JeDi does not suffer from the issues of overfitting and lack of diversity exhibited by the prior models. The examples in the bottom two rows demonstrate JeDi’s high-quality results on challenging personalization tasks.

Abstract

Personalized text-to-image generation models enable users to create images that depict their individual possessions in diverse scenes, finding applications in various domains. To achieve the personalization capability, existing methods rely on finetuning a text-to-image foundation model on a user’s custom dataset, which can be non-trivial for general users, resource-intensive, and time-consuming. Despite attempts to develop finetuning-free methods, their generation quality is much lower compared to their finetuning counterparts. In this paper, we propose Joint-Image

Diffusion (JeDi), an effective technique for learning a finetuning-free personalization model. Our key idea is to learn the joint distribution of multiple related text-image pairs that share a common subject. To facilitate learning, we propose a scalable synthetic dataset generation technique. Once trained, our model enables fast and easy personalization at test time by simply using reference images as input during the sampling process. Our approach does not require any expensive optimization process or additional modules and can faithfully preserve the identity represented by any number of reference images. Experimental results show that our model achieves state-of-the-

art generation quality, both quantitatively and qualitatively, significantly outperforming both the prior finetuning-based and finetuning-free personalization baselines. Project page <https://research.nvidia.com/labs/dir/jedi>

1. Introduction

The state-of-the-art in text-to-image generation has advanced significantly in the last two years, propelled by the emergence of large-scale diffusion models and paired image-text datasets [2, 3, 7, 23, 26, 27, 29]. Despite their superior capability in generating high-quality images well-aligned to the input text prompts, existing models cannot generate novel images depicting specific custom objects or styles that are only available as few reference images external to the training datasets. To address this important use case, various personalization methods have been developed.

The key challenge of personalized image generation is to produce distinct variations of a custom subject while preserving its visual appearance. Most existing approaches achieve this goal by finetuning a pre-trained model on a reference set of images to make it memorize the custom concept. Although these methods can yield good synthesis results, they require substantial resources and a long training time to fit the custom subject, and more than one reference image is needed to avoid overfitting. To overcome these challenges, there has been recent interest in developing *finetuning-free personalization* methods [17, 31, 33, 34]. These methods typically encode the reference images into a compact feature space, and condition the diffusion model on the encoded features. However, the encoding step results in information loss, leading to poor appearance preservation, especially for challenging unusual objects as seen in Fig. 1. Therefore, the performance of encoder-based personalization techniques is inferior to finetuning-based approaches.

In this paper, we present JeDi, a novel approach for finetuning-free personalized text-to-image generation that excels at preserving input reference content. Our core idea is to train a diffusion model to learn a joint distribution of multiple related text-image pairs that share a common subject. As illustrated in Fig. 2, this goal is achieved using two key ingredients: First, we construct a synthetic dataset of related images in which each sample contains a set of text-image pairs that share a common subject. We present a scalable approach for creating such a dataset using LLMs and pre-trained single-image diffusion models. Second, we modify the architecture of existing text-to-image diffusion models to encode relationships between multiple images in a sample set. Specifically, we adapt the self-attention layers of the diffusion U-Net so that the attention blocks corresponding to different input images are coupled. That is, the self-attention layer corresponding to each image co-attends to every other image in the sample set. The use of the cou-

Table 1. In contrast to prior work, JeDi does not require finetuning or the use of image encoders for image personalization.

Method	Finetuning-free	Encoder-free
DreamBooth [28]	✗	✓
CustomDiffusion [16]	✗	✓
ELITE [33]	✓	✗
BLIPDiffusion [17]	Optional	✗
JeDi	✓	✓

pled self-attentions at different levels of hierarchy in the U-Net provides a much stronger representation needed for good input preservation.

At test time, JeDi can take multiple text prompts as input and generate images of the same subject in different contexts. By simply substituting reference images as observed variables in the sampling process, JeDi can generate personalized images based on any number of reference images. We utilize guidance techniques [12] on reference images to further improve the image alignment. JeDi can achieve high-fidelity personalization results even in challenging cases involving unique subjects (Fig. 1, 7), using as few as a single reference image.

Our **key contributions** are summarized as follows:

- We propose a finetuning-free text-to-image generation method with a novel joint-image diffusion model.
- We present a simple and scalable data synthesis pipeline for generating a multi-image personalization dataset with images sharing the same subject.
- We design novel architecture and sampling techniques such as coupled self-attention and image guidance for achieving high-fidelity personalization.

2. Related Work

2.1. Text-to-Image Generation

Denosing diffusion models [9, 13, 32] formulate the image generation task as a series of progressive denoising steps. The denoising network can be trained conditioned on text embeddings to generate images from an input caption. DALL-E2 [26] achieves high-resolution text-to-image synthesis using two diffusion models: the first model transforms a CLIP text embedding to a CLIP image embedding, while the second model transforms the image embedding to an output image. Imagen [29] trains a cascaded diffusion model conditioned on T5 language embeddings [25]. eDiff-I [2] uses an ensemble of expert denoisers to increase the model capacity, with each expert specializing in a specific noise range. Latent diffusion models [7, 23, 27] train the diffusion model in a compact latent space of an autoencoder for efficient training and sampling.

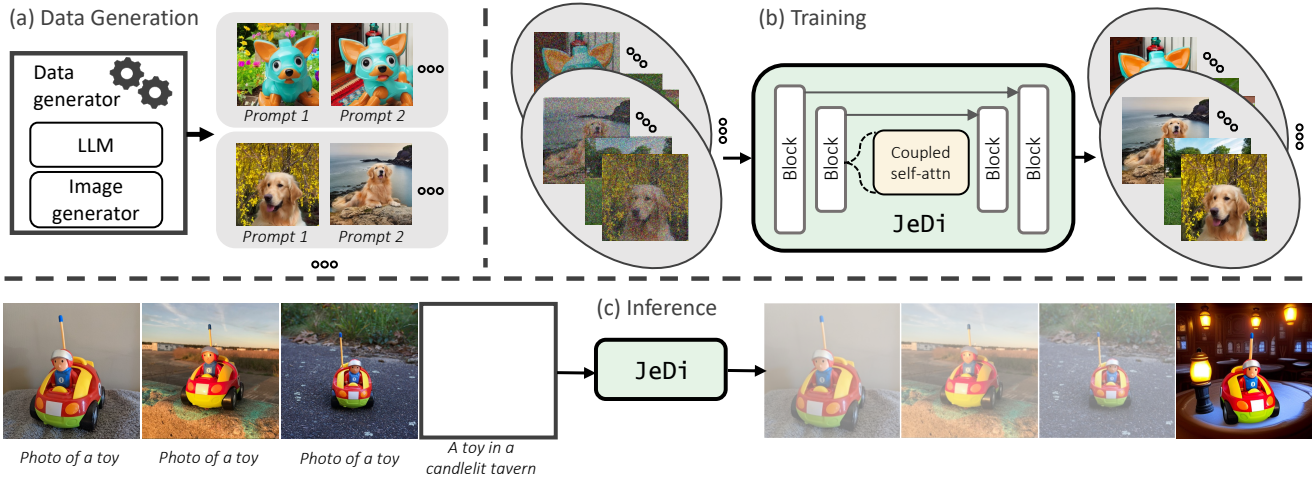


Figure 2. **Overall framework.** (a) We generate training data by using large language models and prompting pretrained single-image diffusion models. (b) During training, the JeDi model learns to denoise multiple same-subject images together, where each image attends to every image of the same subject set through coupled self-attention. (c) At inference, personalized generation is performed in an inpainting fashion where the goal is to generate the missing images of the joint-image set.

2.2. Personalized Text-to-Image Generation

Finetuning based methods. Most prior works achieve image personalization by finetuning the diffusion model on a custom dataset. Dreambooth [28] finetunes the entire model weights on the reference set, with a loss on images of similar concepts as regularization. CustomDiffusion [16] optimizes only a few parameters to enable fast tuning, and combines multiple finetuned models for multi-concept personalization. Textual Inversion [10] projects the reference images onto the text embedding space through an optimization process. SVDiff [11] finetunes only the singular values of the weight matrices to reduce the risk of overfitting. These finetuning-based methods require a substantial amount of resources and long training time, and often need multiple reference images per custom subject.

Finetuning-free methods. To improve the efficiency of image customization, there has been recent interest in developing finetuning-free methods. These approaches typically use an image encoder to encode a reference image onto a compact feature space, and train the diffusion model conditioned on this feature vector. BLIPdiffusion [17] uses BLIP-2 [18] encoder, while FastComposer [34] uses a CLIP [24] encoder for image encoding. ELITE [33] and InstantBooth [31] use a learnable image encoder trained jointly with the diffusion model. These encoder-based methods produce reasonable results for common subjects, but often fail to generate uncommon subjects and preserve fine-grained details due to the information loss in the encoding step. In contrast, our approach directly trains a joint-image diffusion model without an encoding step, resulting in better input preservation even for challenging objects.

Personalization dataset. To achieve finetuning-free personalization, a training dataset comprised of same-

subject image sets is required. Methods like [17, 31, 33, 34] rely on image augmentation and background removal to construct training data, which often does not provide sufficient variations for the same subject. To improve the diversity, we present a scalable data generation approach by prompting pre-trained single-image diffusion models to produce multi-image photo collages with good variation.

3. Method

3.1. Dataset Creation

Training a model to produce a joint distribution of multiple same-subject images requires a dataset where each sample is a set of images sharing a common subject. While there exist some same-subject datasets such as CustomConcept101 [16] and DreamBooth [28], they are small in scale and lack sufficient variations desired for diffusion model training. Therefore, we create a diverse large-scale dataset of image-text pairs containing the same subject, called the Synthetic Same-Subject (S^3) dataset, using large language models [22] and single-image diffusion models [23].

Fig. 3 illustrates our data generation process. We first start with a list of common objects and prompt ChatGPT to generate a text description for each object in the list. Then, we use the pre-trained SDXL [23] model to generate a dataset of same-subject photo collages by appending the text “photos of the same” to each of the text prompt generated in the previous step. We observed that by prompting the SDXL model this way, it can generate photo collages of the same subject with varying poses. However, the generated images usually contain a close-up view of an object in a simple background. To increase the data diversity, we employ a post-processing step that performs background augmentation on the generated objects.

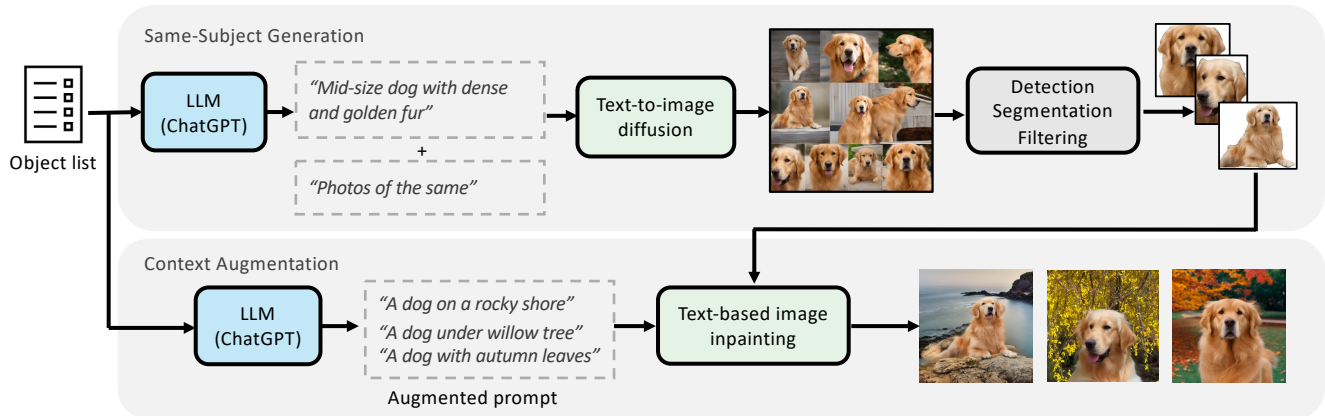


Figure 3. **Data generation process.** We construct Synthetic Same-Subject (S^3) dataset by first prompting the pretrained text-to-image diffusion models to generate same-subject photo collages, and then increasing the diversity using text-based background inpainting.



Figure 4. **Samples from the synthetic same-subject (S^3) dataset.** Each column denotes different images from one joint-data sample.

Given a generated photo collage, we first run object detection [19] and segmentation [15] to separate out object instances and extract foreground region. We discard pairs of instances with CLIP [24] image scores lower than 0.95 as they may not contain the same subject. We then paste the object at a random location in a blank image, and use the stable diffusion inpainting model to inpaint the background based on a new prompt related to the same object name. In addition, we use InstructPix2Pix [4] to stylize the generated samples with a probability of 0.5 to increase style variation using randomly selected style prompts. Fig. 4 shows some examples of the text-image data generated through this process. The generated samples have consistent subjects with good diversity and pose variations.

3.2. Joint-Image Diffusion

The goal of training a joint-image diffusion model is to generate multiple related images sharing the same subject. Conventional diffusion models [2, 21, 27, 29], however, can only generate individual images independently as the network architecture does not have any connections between different samples in a batch. We found that with a simple modification, a single-image diffusion model can be

adapted to a joint-image model which can generate images having related content (such as the same custom subject).

More specifically, given a set of same-subject noisy input images, we modify the attention layers of the U-Net to fuse together the self-attention features for different images in the same set. As illustrated in Fig. 6, a coupled self-attention layer has features at each spatial location attending to every other location across all images in the set. Since the U-Net architecture has attention layers at various resolutions, the use of coupled self-attentions at multiple resolutions makes the generated image set to have consistent high-level semantic features as well as low-level attributes. Fig. 5 visualizes the pixel-wise correspondences and attention heat maps of coupled self-attention layers. We observe that the co-attended regions across different images form the right correspondences across all resolutions.

After a coupled self-attention layer, the output is fed to a regular cross-attention layer, which aligns the visual feature of each image to the corresponding text prompt. The coupled self-attention layer can be implemented by simply adding two reshaping operations before and after a regular self-attention, thus enabling simple and easy adaptation from pre-trained single-image diffusion models.

Fig. 2 (b) illustrates the training process. We start by creating noisy same-subject data by adding isotropic Gaussian noise, and train the joint-image diffusion model to denoise the data. Ideally, there is no limit on the size of each image set of the same subject. In our experiments, we randomly set the size to 2, 3, or 4 during training. The training loss of Jedi is very similar to that of a regular diffusion model. We use ϵ -prediction and a simplified training objective introduced in [13]. The loss function is as follows,

$$L = \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1), t \sim [1,T]} [\|\epsilon - \epsilon_{\theta}(\mathbf{x}_t)\|_2^2], \quad (1)$$

where ϵ_{θ} represents the network parameterized by θ , T is the number of diffusion steps, \mathbf{x}_t is the t -step noised image set of size N , i.e. $\mathbf{x} = [x^1, x^2, \dots, x^N]$. We omit the text and

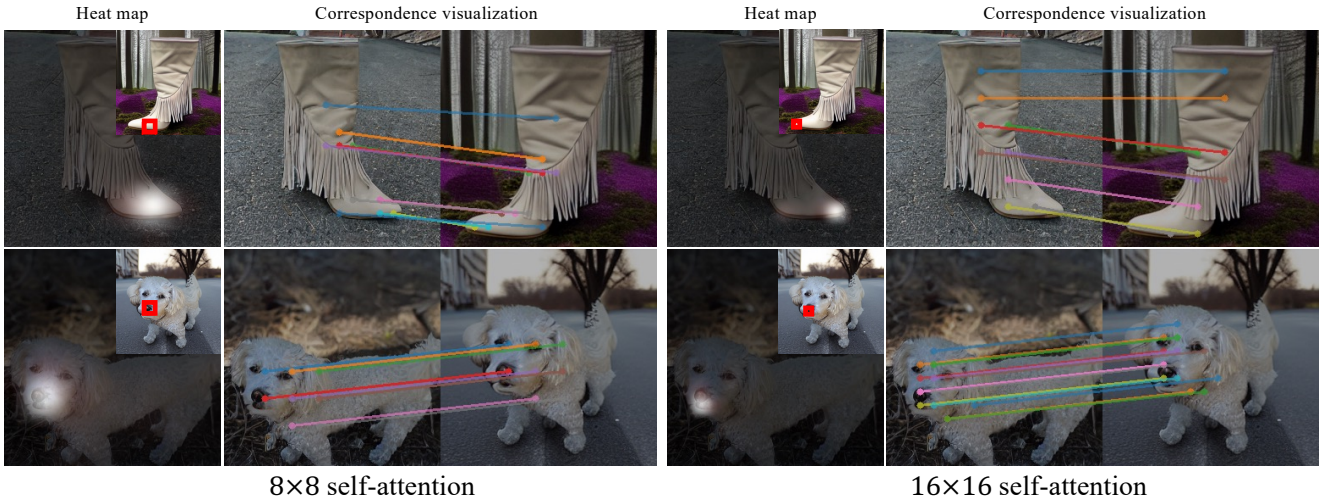


Figure 5. **Visualization of the coupled self-attentions.** For both scales (8x8 and 16x16), the correspondence map (Corr.) shows the connections with the highest weights between elements in the two images. The heatmap visualizes the distribution of the attention weights in an image for a specific element in another image (marked with a red box). We observe that similar regions in different images are co-attended in the coupled self-attention layers.

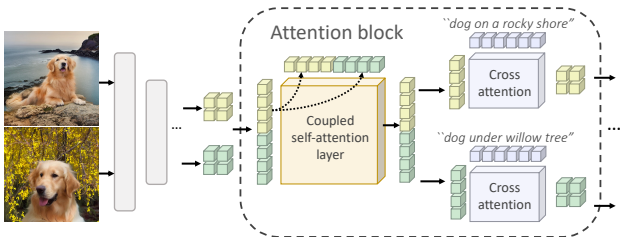


Figure 6. **Attention architecture.** In a coupled self-attention block, features corresponding to each spatial location attends to every location across all images in the image set. Following the coupled self-attention block, features of each individual image attends to their respective text embedding in cross-attention layers.

timestep conditioning in $\epsilon_\theta(\cdot)$ for simplicity. After training, the `Jedi` model can take multiple text prompts as input and generate images containing the same subject.

3.3. Personalized Text-to-Image Generation

Personalization as inpainting. While the joint-image diffusion model discussed in the previous section can generate same-subject images, it does not input a reference image that needs to be personalized. In this section, we propose to solve the input image personalization problem by casting it as an inpainting task. That is, given a few text-image pairs as reference, the task of generating a new personalized sample can be viewed as inpainting the missing images of a joint-image set containing reference images (Fig. 2 (c)).

We design the inpainting model by modifying the input layer of the diffusion U-Net so that it can be conditioned on reference images. More specifically, the input to the diffusion model is a concatenated list of noisy images, reference images and a binary mask indicating whether the reference image is used or not. When the binary mask is all 0's, the

reference image is used. On the other hand, when the binary mask is all 1's, the reference image is an empty black image (indicating the missing images that need to be generated). During training, for every image in the joint-image set, we use the reference image with a probability of 0.5 i.e., we assign the binary mask to 0 with probability 0.5. The training loss can then be written as follows,

$$L = \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1), t \sim [1,T], \mathbf{m} \sim \binom{N}{0.5}} [\|\epsilon - \epsilon_\theta(\mathbf{x}_t, \hat{\mathbf{x}}, \mathbf{M})\|_2^2], \quad (2)$$

where \mathbf{M} is the spatially repeated tensor of a binomial variable \mathbf{m} ; $\hat{\mathbf{x}} = \mathbf{x} \odot \mathbf{M}$ denotes the reference images, with the unknown elements set to zero.

During inference, we utilize the replacement trick [20, 32] in which the known part of the joint-image set is replaced with the forward diffusion response of the clean reference image. Let $\hat{\mathbf{x}} = [\hat{x}^1, \hat{x}^2, \dots, \hat{x}^n]$ be the reference images, which are the known elements in the image set $\mathbf{x} = [\hat{x}^1, \dots, \hat{x}^n, x^{n+1}, \dots, x^N]$ to be generated. During sampling, at each diffusion step t , we only keep the backward diffusion output for the unknown elements x^{n+1}, \dots, x^N while replacing the known part with the forward diffusion output, i.e. the noised real images $\hat{x}_t^1, \dots, \hat{x}_t^n$.

Image guidance. Classifier-free guidance is a popular technique used in single-image diffusion models to make the image generations more aligned to the input conditioning [12]. To improve the faithfulness of the generated samples to the input reference images, we use image guidance in addition to the text guidance during sampling. The score function with the use of image guidance is as follows:

$$\tilde{\epsilon}(\mathbf{x}_t, \hat{\mathbf{x}}, \mathbf{M}) = \epsilon^0 + \lambda[\epsilon_\theta(\mathbf{x}_t, \hat{\mathbf{x}}, \mathbf{M}) - \epsilon^0] \quad (3)$$

where $\epsilon^0 = \epsilon_\theta(\mathbf{x}_t, \mathbf{0}, \mathbf{M})$ represents the unconditional score when the text prompt and all reference images are



Figure 7. **Personalized text-to-image generations.** We show the generations obtained by our JeDi model on challenging uncommon input subjects shown in the left column. JeDi accurately preserves the reference image content while being faithful to the text prompt.

set to null; λ is the guidance scale. We find that the use of image guidance can significantly improve the fidelity to the input reference images.

4. Experiments

Dataset. We construct the Synthetic Same-Subject dataset (S^3 dataset) for training as described in Sec. 3.1. After CLIP filtering, we obtain 1.6M sets of images with each set containing 2-4 images. We also include the video frames from WebVid10M [1] and rendered multi-view images from Objaverse [8] during training, as the frames from the video and the rendered images from the asset usually have the same subject. We use the original video caption or asset caption as text prompts for all images obtained from the same video/asset. Additionally, we include the single-image data from LAION aesthetic dataset [30] and use a set size of 1 for these images. To evaluate our models, we use the test dataset proposed in DreamBooth [28]. DreamBooth test set contains 30 real-world subjects with 4-6 images and 25 prompts for each subject.

Evaluation metrics. The two main evaluation criteria

for personalized text-to-image generation include (1) alignment between generated images and the input text prompts, and (2) faithfulness of the generations to the input reference images. We use the CLIP image-text similarity (CLIP-T) between the generated images and the input captions for (1). For (2), we follow prior works [17, 28, 33] and use the cosine similarity of CLIP [24] image embedding (CLIP-I) and DINO [5] image embedding (DINO) between the generated images and the reference images. DINO is considered to be a preferred metric for measuring image similarity as it is sensitive to the appearance variations of different images in the same concept class. Additionally, we also report CLIP-I and DINO scores only on the foreground masked images, i.e., images with foreground objects cutout using object detection and segmentation [14, 19]. This helps remove the background variations when computing the image similarity scores to better reflect the faithfulness to the reference subject. We call these metrics MCLIP-I and MDINO.

Implementation details. We implement our method based on StableDiffusion V1.4 to enable fair comparison with prior approaches. We train the model with batch size

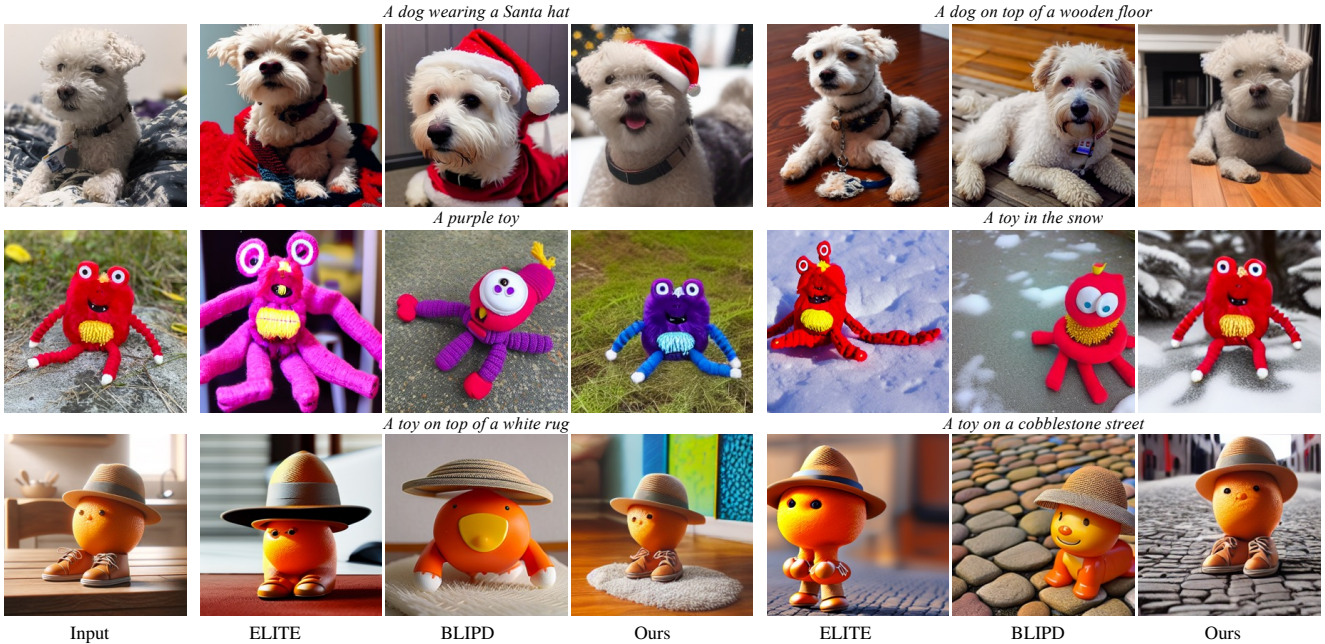


Figure 8. **Visual comparison with finetuning-free methods.** JeDi can faithfully preserve the details of input content even for challenging uncommon objects (row 2 and 3). ELITE and BLIPD fail in such cases and can get good results only in common object classes (row 1).



Figure 9. **Visual comparison with finetuning methods.** Finetuning methods memorize the input images (DB outputs) or result in poor input preservation (CD outputs) when the number of reference images is small. On the other hand, JeDi can generate images that are faithful both to the input images and the text prompt.

2048 and a learning rate $5e-5$. We initialize the weights using the pre-trained StableDiffusion model. For each batch, we randomly sample image sets from S^3 , WebVid10M, Objaverse and LAION datasets with equal probability. We randomly choose the image set size between 2-4 except for LAION, where the set size is always 1. Our model takes 36 hours to train on 32 A100 GPUs for 140K steps.

4.1. Comparison

We compare our method with the state-of-the-art finetuning-free methods - BLIPD [17] and ELITE [33], along with two finetuning-based methods - DreamBooth

(DB) [28] and CustomDiffusion (CD) [16]. For each subject, we use one input reference image for finetuning-free approaches and three for finetuning-based approaches.

Comparison to finetuning-free methods. Fig. 8 shows the visual comparison of our results to BLIPD and ELITE. It can be seen that our method can faithfully capture the visual features of the input reference image, including both semantic attributes and low-level details. However, the generations from BLIPD and ELITE can only roughly resemble the color patterns and semantic features of the input. We also observe that BLIPD and ELITE can produce reasonable results for common object classes such as dogs. This is

Table 2. **Quantitative comparisons.** All finetuning-based methods use 3 input images, while the finetuning-free methods use 1 input image. JeDi model with 1 input image outperforms all finetuning-free baselines, while the JeDi model with 3 input images outperforms all finetuning baselines. JeDi obtains a much higher masked DINO scores, which suggests that we achieve stronger input subject preservation compared to other baselines.

	Method	CLIP-T (↑)	CLIP-I (↑)	MCLIP-I (↑)	DINO (↑)	MDINO (↑)
Finetuning-based	DreamBooth [28]	0.2812	0.8135	0.8683	0.6341	0.7115
	Custom Diffusion [16]	0.3015	0.7952	0.8640	0.6343	0.7109
Finetuning-free	BLIP Diffusion [17]	0.2934	0.7899	0.8620	0.5855	0.6692
	ELITE [33]	0.2961	0.7924	0.8615	0.5922	0.6805
	JeDi (1 input)	0.3040	0.7818	0.8764	0.6190	0.7510
	JeDi (3 inputs)	0.2932	0.8139	0.9011	0.6791	0.8037

Table 3. **Effect of varying the S³ dataset size.**

	Dataset size			
	0	0.6M	1.2M	1.6M
DINO (↑)	0.8895	0.7501	0.7443	0.7467
MDINO (↑)	0.8931	0.8639	0.8572	0.8636
CLIP-T (↑)	0.2524	0.3020	0.3017	0.3015

Table 4. **Ablation of different joint-image diffusion designs.**

	CLIP baseline	JeDi	(+) CLIP emb	w/o IG
DINO (↑)	0.3411	0.7501	0.7432	0.4652
MDINO (↑)	0.4394	0.8639	0.8617	0.5922
CLIP-T (↑)	0.3325	0.3020	0.3041	0.3259

because their encoder can easily recognize the common object categories (such as the dog breed) which makes the personalized generation easier. However, for unique uncommon objects, their results tend to be much worse, *e.g.* the toy in the second row and the image in the third row. Note that even for the common classes such as the dog example in the first row, the generations from BLIPD and ELITE can miss some input features (different haircuts despite being the same breed). In contrast, our method eliminates information loss caused by the encoder and results in much better preservation of the custom concept. This is also reflected in the quantitative comparison Table 2, where our method outperforms BLIPD and ELITE by a large margin.

Comparison to finetuning-based methods. Fig. 9 shows the visual comparison of our approach with DreamBooth (DB) and CustomDiffusion (CD). When the number of reference images is limited, it is challenging for finetuning-based methods to avoid overfitting and generate novel variations of the input subject. From Fig. 9, we see that DB often directly copies the input image due to overfitting, while the images generated by CD do not faithfully preserve the features of the input subject. In contrast, our method creates proper variations of the reference subjects without changing its key visual features. Even without any expensive finetuning, our method outperforms DB and CD in quantitative comparisons when we provide the same

number of reference images (JeDi-3), as shown in Table 2.

4.2. Ablation Study

Size of S³ dataset. Table 3 reports the results of training our model using different numbers of synthetic images. Dataset size 0 refers to training the model on only videos and multi-view images. This setting yields the best image alignment (DINO and MDINO) and the worst text alignment as the model learns a shortcut to ignore the text prompts and copy the input images. The performance of columns 2-4 are roughly similar, which shows that we do not obtain much gains by increasing the size of the synthetic dataset.

Joint-image diffusion model. In Table. 4, we report the ablation study of different design choices in the training of JeDi. The first column shows the results of a CLIP encoder baseline, which is a single-image diffusion model conditioned on the CLIP image features of the reference image. Our JeDi model yields a much better image alignment than the CLIP encoder baseline, which demonstrates the advantages of using the joint-image model over the image encoders. We also find that adding CLIP image embedding as extra conditional input to JeDi (+ CLIP emb) does not improve the performance as shown in column 3. This implies that the joint-image model already captures the information extracted in the image embedding. The last column reports the results without image guidance (w/o IG). By comparing the second and the last columns, we can see that image guidance is crucial to obtaining good personalization results.

5. Conclusion

This paper presents JeDi, a novel approach for finetuning-free personalized text-to-image generation using a joint-image diffusion model. We show how a single image diffusion U-Net can be adapted to learn a joint image distribution using coupled self-attention layers. To train the joint-image diffusion model, we construct a synthetic dataset called S³, in which each sample contains a set of images sharing the same subject. The experimental results show that the proposed JeDi model outperforms the previous approaches both quantitatively and qualitatively in benchmark datasets.

JeDi: Joint-Image Diffusion Models for Finetuning-Free Personalized Text-to-Image Generation

Supplementary Material

6. Implementation Details

In this section, we describe the key modifications based on StableDiffusion v1.4¹ to implement the proposed method. We point to the original location in StableDiffusion code and highlight the modified lines in each code snippet.

Data loading. We store all images and captions that belong to an image set in a single image file and text file. Images are vertically concatenated and the captions corresponding to different images are separated by a special token `<|split|>`, as illustrated in Fig. 10. This enables us to easily reuse existing single-image dataloaders in PyTorch. We only use square images in training and obtain

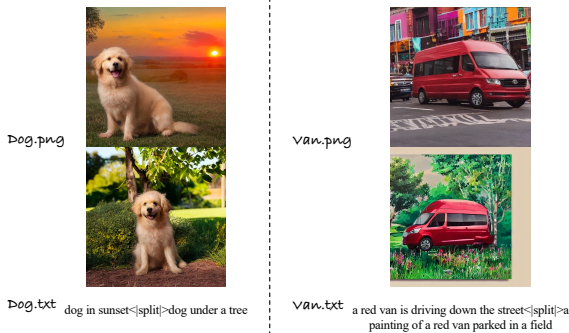


Figure 10. **Data format.** An image file is the concatenation of all images in the set. The text file contains corresponding captions separated by a special token.

the image set size by dividing the image height with the image weight. Given a batch of data, we extract individual images and text, and obtain the size of the image set `ng` as follows,

Joint-image diffusion models. The proposed joint-image diffusion model can be easily implemented based on a single-image diffusion model with a few simple modifications as follows,

Personalization as inpainting. As described in Sec. 3.3 of the paper, we cast the personalized generation problem into an inpainting task. In training, the inpainting masks are generated randomly. First, a binary random vector `ng_mask` is sampled with every bit set to zero or one with equal probability. Then a mask of the same spatial size as the input image is constructed by replicating `ng_mask` alone the height and width dimensions. The actual input fed into the U-Net is the concatenation of the mask `mask`, noisy image `x`, and masked clean image `x0*(1-mask)`.

¹We compare with previous approaches using the model based on v1.4 while the images in Fig. 1 are generated by the model based on SDXL.

```
1 # ldm/models/diffusion/ddpm.py#L865-L868
2 def get_group_data(self, img, txt):
3     b, c, h, w = img.shape
4     ng = h//w
5     assert h==w*ng
6     img = img.view(b, c, ng, w, w)
7     img = img.transpose(1, 2).reshape(b*ng, c, w, w)
8     txt = reduce(lambda
9     ↪ a, b: a+b, [t.split("<|split|>") for t in
10    ↪ txt])
11    return img, txt, ng
12 def shared_step(self, batch, **kwargs):
13    k1, kc = self.first_stage_key, self.cond_stage_key
14    assert kc=='txt'
15    batch[k1], batch[kc], ng =
16    ↪ self.get_group_data(batch[k1])
17    x, c = self.get_input(batch,
18    ↪ self.first_stage_key)
19    loss = self(x, c)
20    return loss
```

```
1 # ldm/modules/attention.py#L211C15-L215
2 # ng: size of the image set
3 def _forward(self, x, context=None, ng=None):
4     b, l, c = x.shape
5     if ng is not None:
6         x = x.view(-1, ng*l, c)
7     x = self.attn1(self.norm1(x)) + x
8     if ng is not None:
9         x = x.view(b, l, c)
10    x = self.attn2(self.norm2(x), context=context)
11    ↪ + x
12    x = self.ff(self.norm3(x)) + x
13    return x
```

The key implementation can be found in the following code snippet.

```
1 # ldm/modules/diffusionmodules/openaimodel.py
2 #L730-L730
3 bs=x.size(0)
4 if ng_mask is None:
5     ng_mask =
6     ↪ th.empty(bs, device=x.device).bernoulli_(1)
7 mask = ng_mask[:, None, None, None]
8 x0 = x0*(1-mask)
9 mask = mask.expand(-1, -1, *x0.shape[-2:])
10 x = th.cat((x, x0, mask), 1)
11 h = x.type(self.dtype)
```

At test time, `x0` is the concatenation of all input images and an all-zero image. The corresponding elements in `ng_mask` are set to 0 for the input images and 1 otherwise.

Synthetic same-subject (S³) dataset. Algorithm 1 describes the training data generation (Fig. 3 in the paper) in details. Please note that we use the term *instance segmentation* for simplicity; however, in our implementation, we combine an object detection model [19] and a segmentation model [15] to separate object instances rather than using an actual instance segmentation model. The *object-centric* prompts (GPT(*l*, *object*)) in Algorithm 1 are generated by instructing ChatGPT to generate details of an object *l*, and

Table 5. Quantitative comparisons on the unique subject test set.

Method	CLIP-T (↑)	CLIP-I (↑)	MCLIP-I (↑)	DINO (↑)	MDINO (↑)
BLIP Diffusion [17]	0.2851	0.8107	0.8234	0.6091	0.6018
ELITE [33]	0.2193	0.6082	0.6430	0.1862	0.2156
JeDi	0.2856	0.8697	0.8838	0.7934	0.7926

the *scene-centric* prompts ($GPT(l, scene)$) in Algorithm 1) are generated by instructing it to describe a scene involving the object l . The list of object names used in our implementation and a random subset of the generated training samples can be found in the attached file. We found that the initial text prompts generated by ChatGPT lack variations and therefore pair the images with the captions obtained from BLIPv2 [17] in training samples.

Algorithm 1: Generating the S^3 dataset.

```

1: Input: A list of object names  $L$ ;
2:   A text-to-image model  $G$ ;
3:   A text-based inpainting model  $G_I$ ;
4:   ChatGPT GPT;
5:   An instance segmentation model  $S$ ;
6:   CLIP image model CLIP;
7: Output: a database  $\mathcal{X}$  of image sets
    $\mathcal{X} = \{X_1, X_2, \dots, X_P\}$  where  $X_p = \{x_1, x_2, \dots, x_{N_p}\}$ 
   is a set of images share a common subject;
8:  $\mathcal{X} \leftarrow \phi$ ;
9: for  $l$  in  $L$  do
10:  Generate an object-centric prompt
    $t_o \leftarrow GPT(l, object)$ ;
11:  Generate an image  $x_0 \leftarrow G(t_o)$ ;
12:  Extract object instances from  $x_0$ :
    $\{\hat{x}_1, \hat{x}_2, \dots, \hat{x}_K\} \leftarrow S(x_0)$ ;
13:  Construct an affinity matrix
    $A, A_{ij} \leftarrow CLIP(\hat{x}_i, \hat{x}_j)$ ;
14:  Construct an adjacent matrix  $M, M_{ij} \leftarrow 1$  if
    $A_{ij} > 0.95$  else 0;
15:  Find connected components  $\mathcal{I} \leftarrow \{I_1, I_2, \dots, I_P\}$  of
   the graph represented by  $M$ ;
16:  for  $I_p$  in  $\mathcal{I}$  do
17:    $X \leftarrow \phi$ ;
18:   for  $i$  in  $I_p$  do
19:    Generate a scene-centric prompt involving
    object name  $l$ :  $t_s \leftarrow GPT(l, scene)$ ;
20:    Paste the object instance image  $\hat{x}_i$  at a random
    location in an empty image  $x$ ;
21:    Inpaint the unknown area in  $x$ :  $x \leftarrow G_I(x, t_s)$ ;
22:     $X \leftarrow X \cup \{x\}$ ;
23:   end for
24:    $\mathcal{X} \leftarrow \mathcal{X} \cup \{X\}$ 
25: end for
26: end for

```

Evaluation metric. We use the CLIP ViT-B/32 to compute CLIP-T, CLIP-I and MCLIP-I. We use DINO ViT-S/16

to compute DINO and MDINO. We use one input image per subject for comparison with finetuning-free methods, and average the pair-wise scores to all real images of the same subject and over all possible choices of the input image. For comparison with finetuning-based methods, we randomly select three input images. In ablation studies we use one randomly selected input image for each subject and only compute the scores using input/output pairs by default unless stated otherwise.

7. Additional Experiments

7.1. Additional Results

Fig. 13 shows additional personalized generation results on real-world human and object images. Our method can generate high-quality images with diverse content while preserving the key visual features of the subjects in input images. Although the model is not trained on human-specific data, it can still generate reasonable results for human subjects, as shown in the second and third row of Fig. 13.

7.2. Comparison with State-of-the-Art Methods

Fig. 14 provides additional visual comparisons with finetuning-based methods DreamBooth (DB) and CustomDiffusion (CD). Finetuning-based methods suffer from the overfitting issue and might fail to preserve the subject identity. For common subjects, they tend to extensively copy from the reference images, adding only minor adjustments to match the given text, *e.g.* in the first example, for the prompt *a backpack in the snow*, DreamBooth nearly replicate a reference image with slight snow patterns added in the bottom. For unique subjects, the finetuning-based methods often fail to preserve the distinctive features, *e.g.* the cartoon character in the fourth row. This is because these methods use the loss on retrieved or generated images of similar subjects as regularization during finetuning. For unique and rare objects, these images can be visually distinct from the reference images and interfere the model from memorizing the custom concept.

To further demonstrate the advantage of our methods for challenging cases, we collect a new test set containing unique subjects with only single input image for each subject. Most the input images are from Reddit AI Art channel². Fig. 15 visualize the input images. The first five rows in Fig. 16 compare the results of our method to state-of-the-art finetuning-free methods BLIP-Diffusion (BLIPD) and ELITE on the unique subject test set. We also include the

²<https://www.reddit.com/r/aiArt/>

results on common subjects from DreamBooth test set in the last two rows for comparison. It can be seen that BLIPD and ELITE can produce reasonable results for common subjects such as the dog of a typical breed and a common stuffed animal (row 6-7). However, for unique subjects, their results hardly resemble the subject from the reference image (row 1-5). In contrast, our method can faithfully capture the key visual features of the subject. The advantage of our method is also clearly reflected in the quantitative results in Table 5, where our method outperforms ELITE and BLIP-Diffusion by a large margin.

7.3. Additional Analysis

Image guidance. As we have discussed in the paper, the use of image guidance can significantly improve the faithfulness to the input images. This is also supported by the visual comparison in Fig. 12 (column 4-5).

In our main experiments in the paper, we use a simple strategy for image guidance where both the image and text input are set to null for unconditional inference. Here we discuss a more flexible guidance strategy to model trade-off between image alignment and text alignment. The score function with flexible image guidance is as follows,

$$\tilde{\epsilon}(\mathbf{x}_t, \hat{\mathbf{x}}, \mathbf{M}) = \epsilon^0 + \lambda_1(\epsilon^1 - \epsilon^0) + \lambda_2[\epsilon_\theta(\mathbf{x}_t, \hat{\mathbf{x}}, \mathbf{M}) - \epsilon^1], \quad (4)$$

where $\epsilon^0 = \epsilon_\theta(\mathbf{x}_t, \mathbf{0}, \mathbf{M})$ represents the unconditional score when the text prompt and all reference images are set to null; ϵ^1 represents the partially conditional score when either the text prompt or reference images are kept. We can compute the partial conditional score ϵ^1 using image conditioning to emphasize text alignment, or text conditioning to emphasize image alignment. We call these two options *text first* and *image first* strategies, respectively. Table 6 reports the quantitative results based on different strategies averaged over a varying guidance scale in [1.5, 10]. It can be seen that the *text first* strategy yields higher DINO and MDINO scores, indicating better image alignment. *Image first* strategy yields a higher CLIP-T score, which indicates better text alignment.

Table 6. Quantitative results with different guidance strategies averaged over a varying guidance scale in [1.5, 10].

Strategy	Text only	Joint	Image first	Text first
DINO	0.4652	0.7268	0.6558	0.7508
MDINO	0.5922	0.8384	0.7863	0.8527
CLIP-T	0.3259	0.3013	0.3156	0.2853

We can also adjust the ratio between the guidance scale of image condition and text condition for more flexible personalized generation. Fig. 11 visualize the change of DINO and CLIP-T scores with the varying ratio. We can see that the use of image guidance is important. Using only text

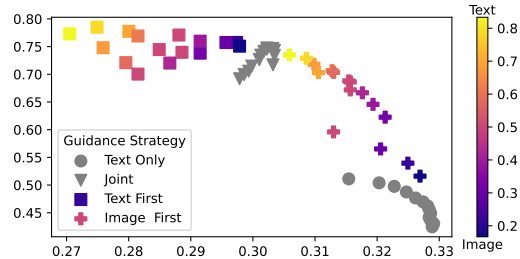


Figure 11. Effect of varying the ratio between the guidance scales for image and text guidance. X-axis: CLIP-T; Y-axis: DINO.

guidance (Text only in Fig. 11) yields low DINO score, indicating low resemblance to the custom subject. By varying the ratio of text and image guidance scales we can balance between subject identity preservation and text alignment. We found that the simple joint guidance strategy (Joint in Fig. 11) usually gives good balanced results.

Table 7. Quantitative comparison with three baseline models.

Model	CLIP baseline	Concate. baseline	Learned baseline	Jedi
DINO	0.3411	0.3379	0.7065	0.7501
MDINO	0.4394	0.4292	0.7740	0.8639
CLIP-T	0.3325	0.3247	0.3015	0.3020

Comparison with baseline models. We have discussed a CLIP-encoder baseline in the paper (CLIP baseline). Here we include the comparison to another two baseline models. Concate baseline: the input images are concatenated to the noisy images to be fed into the UNet. Learned baseline: similar to SuTI [6] where a learnable encoder is used to extract a feature vector from the input images. All models are trained on the same training data described in Sec. 4 of the paper and are based on the same StableDiffusion v1.4 backbone for fair comparison. Table 7 reports the quantitative comparison results and the visual comparison can be found in Fig. 12 (column 1,2,3,5). The results indicate that our method significantly outperforms all baseline models in terms of image alignment, as evidenced by considerably higher DINO and MDINO scores.

Training data identity similarity. Table 8 reports the average CLIP and DINO scores on training samples over the 1,000 ImageNet categories, which indicates a high overall identity similarity for a wide array of categories. For context, we also show the scores on real images from the DreamBooth test set, which has a slightly higher identity similarity but covers much less categories than our dataset.

	Subjects	Categories	CLIP-I (↑)	DINO (↑)
S ³ dataset	1.6M	~ 2K	0.849	0.751
Real images	30	15	0.885	0.774

Table 8. Training data statistics.

Quantitative results with more input images. Although

our method does not have an inherent constraint on the number of inputs, for simplicity, we only use 1-3 input images in the current implementation. We find that our method still outperforms DB and CD, even when they are finetuned with the maximum available reference images in the test set (4-6 images), as shown in Table 9. We will add the experiments with more reference images, e.g. 10, in the revised version.

Table 9. Comparison to DB and CD with the maximum available reference images.

	CLIP-T	CLIP-I (↑)	MCLIP-I (↑)	DINO (↑)	MDINO (↑)
DB	0.2971	0.8025	0.8736	0.6226	0.7175
CD	0.3071	0.7864	0.8586	0.6198	0.7011
Ours (1 input)	0.3040	0.7818	0.8764	0.6190	0.7510
Ours (3 inputs)	0.2932	0.8139	0.9011	0.6791	0.8037

Inference cost. The inference cost is comparable to other methods when N is small (reported in the table below). When N is substantially larger, e.g. a database, we can reduce the inference cost by first finetuning the model on the database, and then retrieving the few images closest to the text prompt to be the actual test time input (please refer to the future work section).

Table 10. Inference time for one diffusion step on one A100 GPU.

Method	BLIPD	ELITE	Ours
Time (second) ↓	0.0492	0.0719	0.0564

8. Limitations and Future Work

A limitation of J_{eDi} is that it needs to process all reference images at inference time. This enables finetuning-free personalization but leads to efficiency drop when the number of reference images increases. Therefore, J_{eDi} is more suitable for subject image generation given a few reference images, and are less efficient in adapting to a new domain given a large database of reference images. A potential solution is to combine J_{eDi} with finetuning-based methods. When a large database of reference images are available, we can first finetune J_{eDi} on the database. Then at inference time, given a text prompt, we retrieve the most relevant images from the database to use as the test-time inputs to J_{eDi} . Another limitation is that the current implementation cannot be directly applied for multi-subject image generation. There are two possible ways to extend J_{eDi} for multi-subject generation: (1) generate multiple subjects sequentially through inpainting, and (2) construct a multi-subject S^3 dataset by combining multiple sets of subjects. We will explore these directions in future work.

References

- [1] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *International Conference on Computer Vision*, 2021. 6
- [2] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 2, 4
- [3] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, Wesam Manassra, Prafulla Dhariwal, Casey Chu, Yunxin Jiao, and Aditya Ramesh. Improving image generation with better captions. <https://api.semanticscholar.org/CorpusID:264403242>. 2
- [4] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2023. 4
- [5] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *International Conference on Computer Vision*, 2021. 6
- [6] Wenhui Chen, Hexiang Hu, Yandong Li, Nataniel Ruiz, Xuhui Jia, Ming-Wei Chang, and William W Cohen. Subject-driven text-to-image generation via apprenticeship learning. *Conference on Neural Information Processing Systems*, 36, 2024. 4
- [7] Xiaoliang Dai, Ji Hou, Chih-Yao Ma, Sam Tsai, Jialiang Wang, Rui Wang, Peizhao Zhang, Simon Vandenhende, Xiao-fang Wang, Abhimanyu Dubey, et al. Emu: Enhancing image generation models using photogenic needles in a haystack. *arXiv preprint arXiv:2309.15807*, 2023. 2
- [8] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2023. 6
- [9] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *Conference on Neural Information Processing Systems*, 2021. 2
- [10] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *International Conference on Learning and Representation*, 2022. 3
- [11] Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris Metaxas, and Feng Yang. Svdiff: Compact parameter space for diffusion fine-tuning. *arXiv preprint arXiv:2303.11305*, 2023. 3
- [12] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS Workshop on Deep Generative Models and Downstream Applications*, 2021. 2, 5
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Conference on Neural Information Processing Systems*, 2020. 2, 4

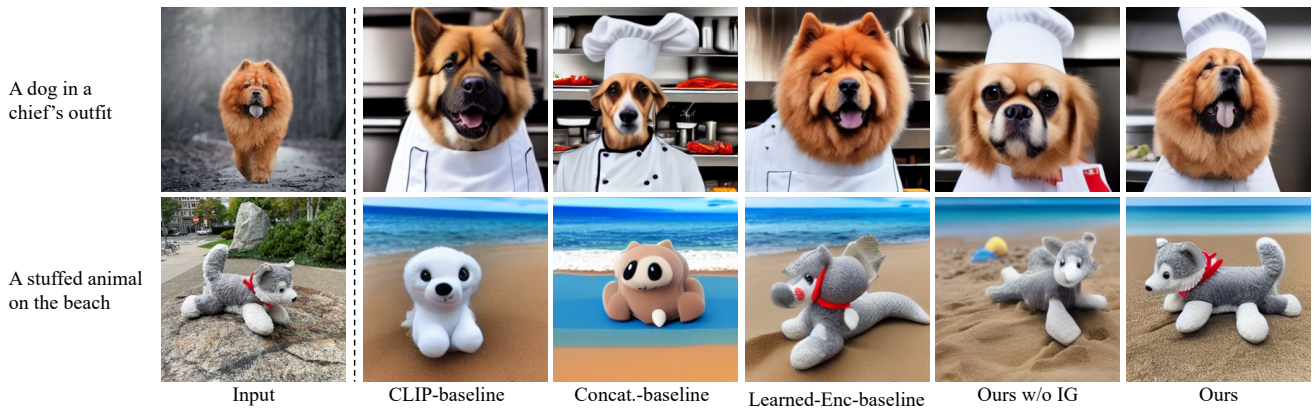


Figure 12. Visual comparisons to baseline models.

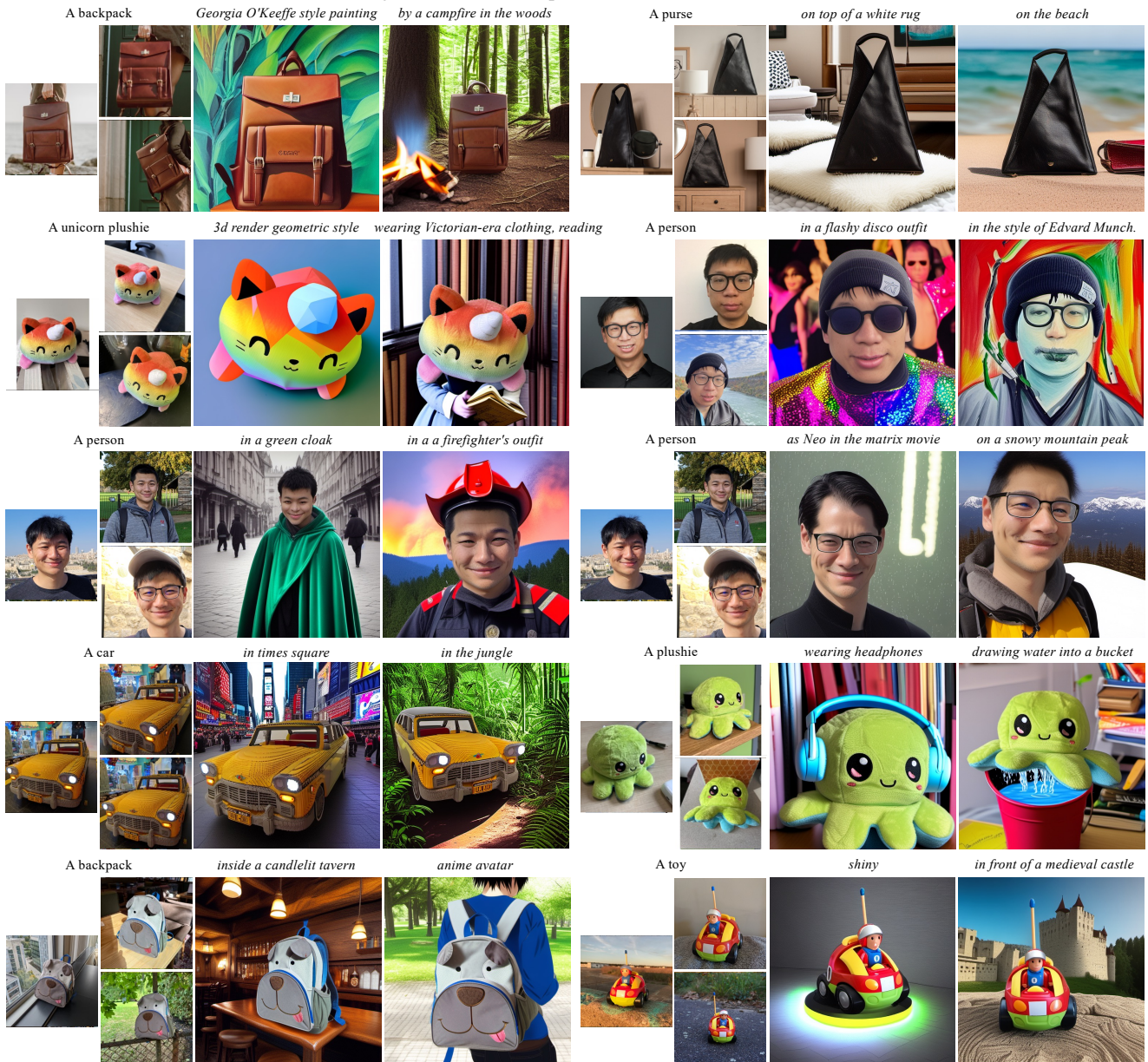


Figure 13. Personalized generation results for human and real-world objects.

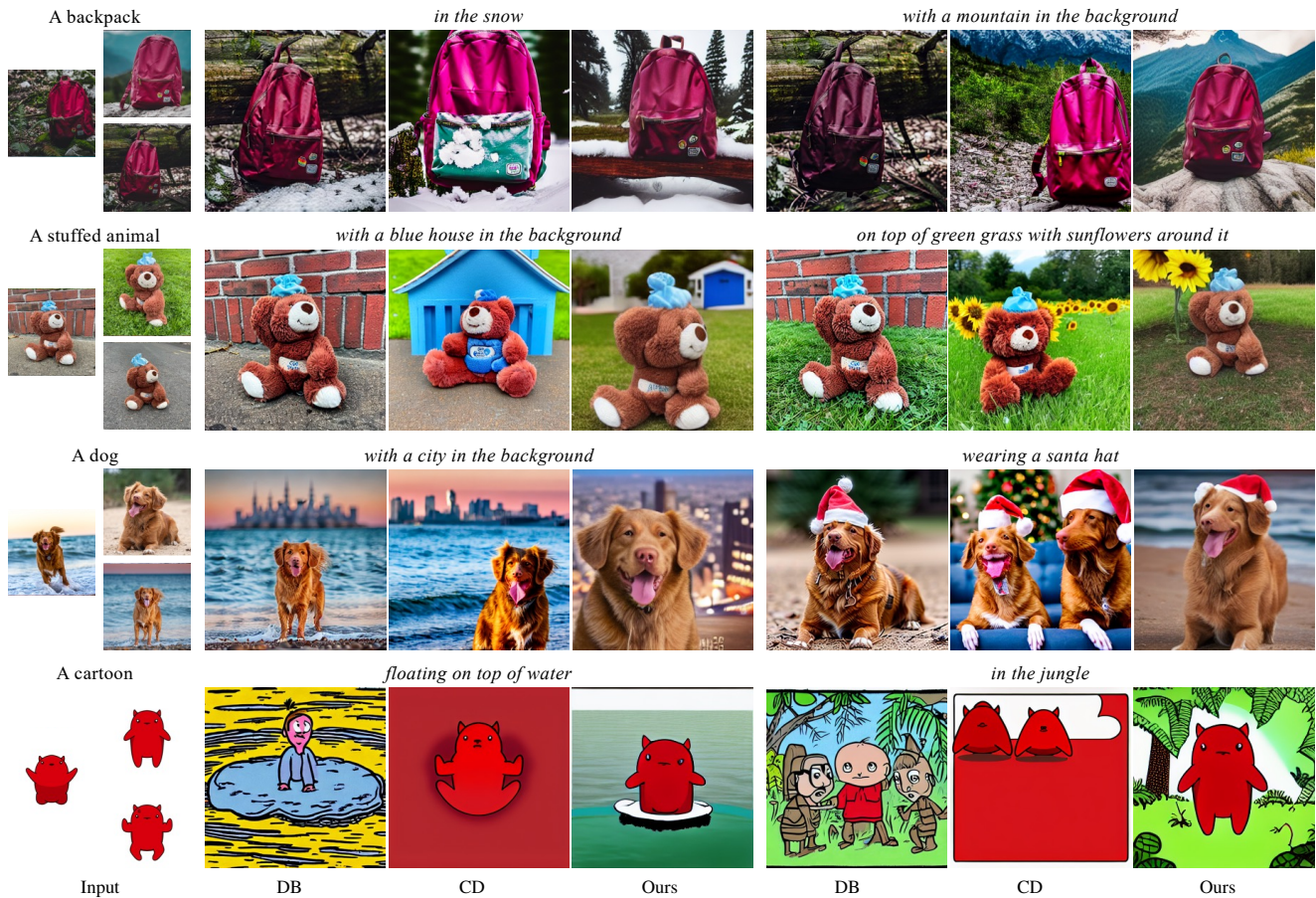


Figure 14. Visual comparison with finetuning-based methods on DreamBooth test set. DreamBooth (DB) and CustomDiffusion (CD) tend to overfit for common subjects and fail to capture the visual features of unique subjects.

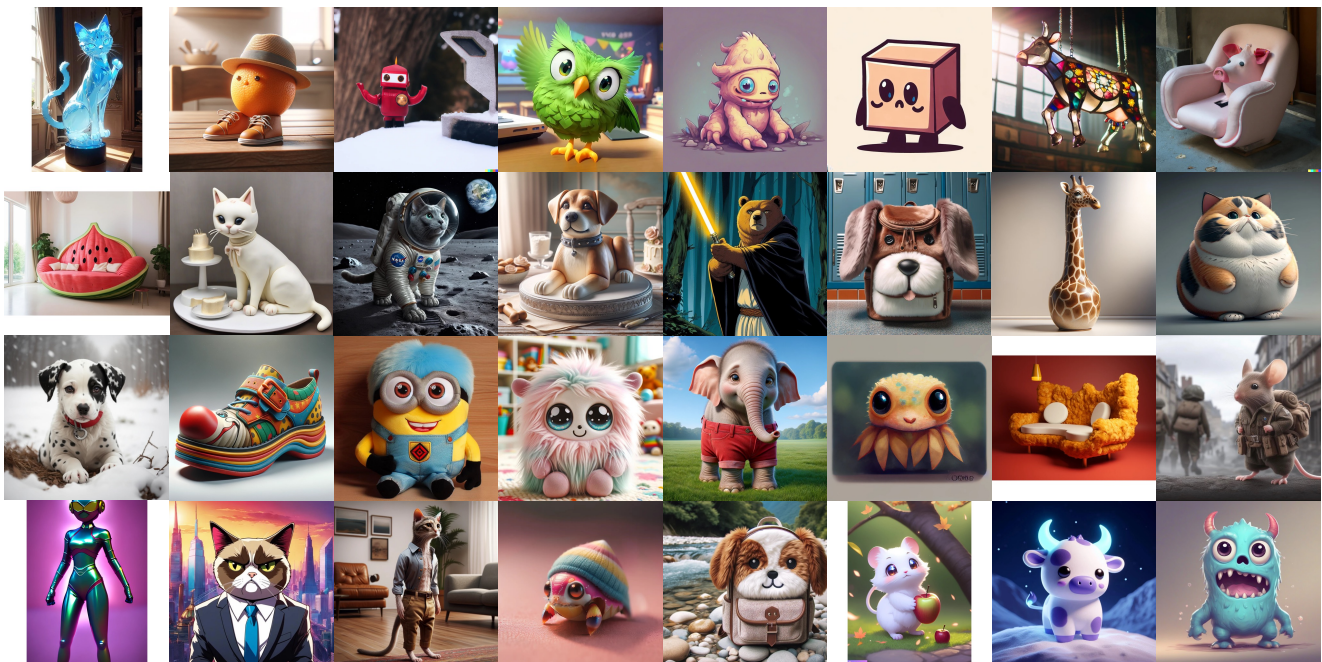


Figure 15. Input images in the unique subject test set.

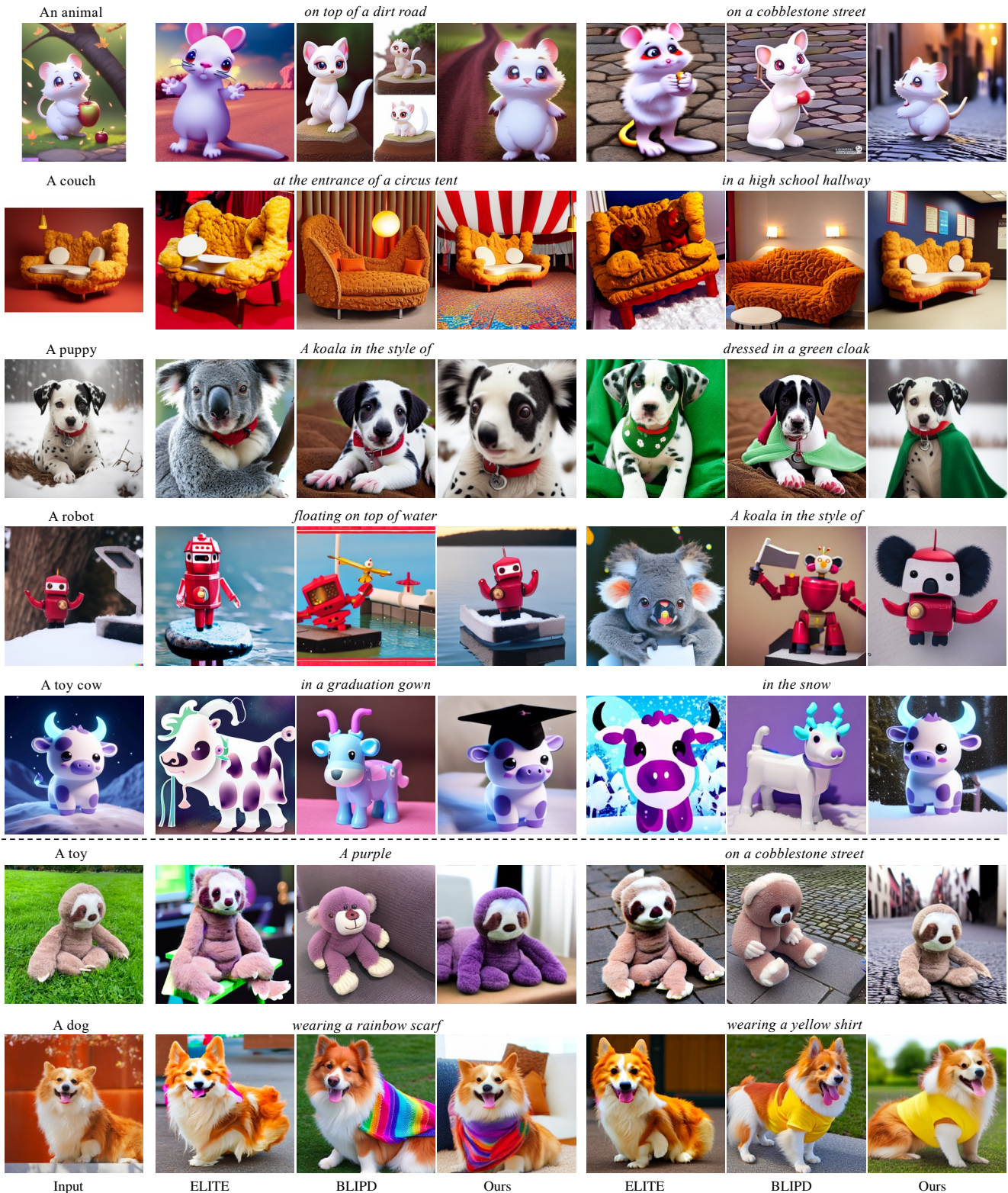


Figure 16. Visual comparison with finetuning-free methods on the unique subject test set (row 1-5) and on Dreambooth test set (row 6-7). BLIP Diffusion and ELITE can generate reasonable results for common subjects but often fail in challenging cases involving unique subjects. In contrast, our method can handle challenging cases and generate personalized images with well-preserved details.

- [14] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023. 6
- [15] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023. 4, 2
- [16] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2023. 2, 3, 7, 1
- [17] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 2, 3, 6, 7, 1
- [18] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023. 3
- [19] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 4, 6, 2
- [20] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 5
- [21] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 4
- [22] OpenAI. Chatgpt. <https://openai.com/blog/chatgpt>. Accessed: 2023. 3
- [23] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 2, 3
- [24] 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. In *International Conference on Machine Learning*, 2021. 3, 4, 6
- [25] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020. 2
- [26] 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. 2
- [27] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 2, 4
- [28] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2023. 2, 3, 6, 7, 1
- [29] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *Conference on Neural Information Processing Systems*, 2022. 2, 4
- [30] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. 6
- [31] Jing Shi, Wei Xiong, Zhe Lin, and Hyun Joon Jung. Instant-booth: Personalized text-to-image generation without test-time finetuning. *arXiv preprint arXiv:2304.03411*, 2023. 2, 3
- [32] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *Proc. ICLR*, 2021. 2, 5
- [33] 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. 2, 3, 6, 7, 1
- [34] Guangxuan Xiao, Tianwei Yin, William T Freeman, Frédo Durand, and Song Han. Fastcomposer: Tuning-free multi-subject image generation with localized attention. *arXiv preprint arXiv:2305.10431*, 2023. 2, 3