# Efficient Personalization of Quantized Diffusion Model without Backpropagation

Hoigi Seo[1*]   Wongi Jeong[1*]   Kyungryeol Lee[1]   Se Young Chun[1,2†]

[1]Dept. of Electrical and Computer Engineering, [2]INMC & IPAI
Seoul National University, Republic of Korea

{seohoiki3215, wg7139, kr.lee, sychun}@snu.ac.kr

## Abstract

*Diffusion models have shown remarkable performance in image synthesis, but they demand extensive computational and memory resources for training, fine-tuning and inference. Although advanced quantization techniques have successfully minimized memory usage for inference, training and fine-tuning these quantized models still require large memory possibly due to dequantization for accurate computation of gradients and/or backpropagation for gradient-based algorithms. However, memory-efficient fine-tuning is particularly desirable for applications such as personalization that often must be run on edge devices like mobile phones with private data. In this work, we address this challenge by quantizing a diffusion model with personalization via Textual Inversion and by leveraging a zeroth-order optimization on personalization tokens without dequantization so that it does not require gradient and activation storage for backpropagation that consumes considerable memory. Since a gradient estimation using zeroth-order optimization is quite noisy for a single or a few images in personalization, we propose to denoise the estimated gradient by projecting it onto a subspace that is constructed with the past history of the tokens, dubbed Subspace Gradient. In addition, we investigated the influence of text embedding in image generation, leading to our proposed time steps sampling, dubbed Partial Uniform Timestep Sampling for sampling with effective diffusion timesteps. Our method achieves comparable performance to prior methods in image and text alignment scores for personalizing Stable Diffusion with only forward passes while reducing training memory demand up to 8.2×. Project page: https://ignoww.github.io/ZOODiP_project/*
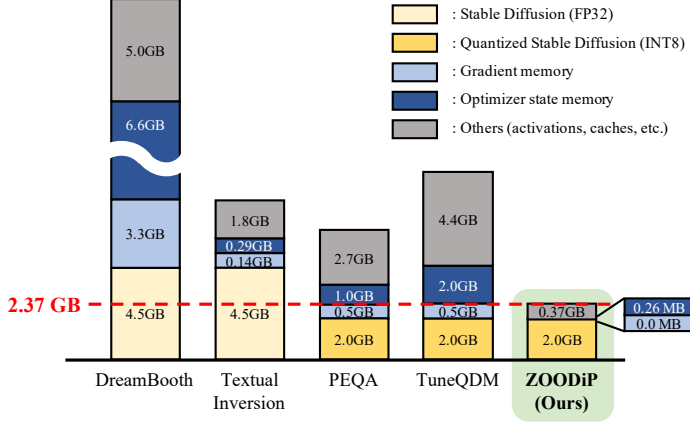
## 1. Introduction

Recent advances in diffusion models [21, 56, 58, 59, 64] have revolutionized generative AI, offering a powerful framework for high-quality, diverse data generation. Diffusion models have shown impressive capabilities across various applications, including image synthesis [21, 37, 45, 56–58, 64], text-guided image synthesis [9, 15, 46, 50, 51], video generation [4, 22, 23, 25, 38], and 3D content generation [31, 35, 47, 55, 63], outperforming traditional generative approaches. However, modern diffusion models face significant memory and computational costs during training, fine-tuning, and inference due to their size.
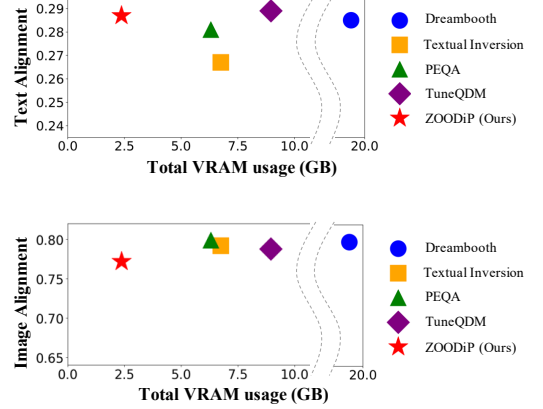
These overheads pose major challenges in personalization tasks, which require fine-tuning diffusion models with only a few user-provided images. Typically, personalization is achieved by either fine-tuning the denoising network [20, 28, 52, 53, 61, 65] or introducing new text tokens [1, 17, 28, 42, 62, 68] to represent desired subjects. However, both approaches require gradient-based optimization, leading to significant memory and computational overhead, which is especially problematic in memory-constrained on-device training for sensitive data. To tackle the overheads in personalization, efficient tuning techniques have emerged. Existing methods reduce trainable parameters [20, 28, 53, 61, 65], leverage quantized models [26, 54], or employ gradient-free optimization with evolution strategy [16]. However, these methods have limitations: 1) they mostly rely on backpropagation, unsuitable for most mobile processors, which are designed to accelerate inference [48]; 2) they still require significant memory for storing activations and gradients; 3) evolutionary algorithms can be unstable and inefficient in small-batch.

Here we propose Zeroth-Order Optimization for Diffusion model Personalization (ZOODiP) that can personalize Stable Diffusion with 2.37GB VRAM consumption using only forward passes without compromising the image quality. Our method relies on three key observations: First, Zeroth-Order (ZO) optimization effectively handles non-differentiable objectives [39] (*e.g.*, accuracy) during training. Second, tokens optimized via Textual Inversion [17] undergo significant changes in a low-dimensional subspace. Principal Component Analysis (PCA) on the token embeddings reveals that the initial and personalized tokens pri-

---

1

(a) GPU memory breakdown across various personalization methods.

(b) VRAM usage versus image and text alignment scores.

Figure 1. Analysis of memory consumption and performance of Stable Diffusion personalization methods. **(Left)** GPU memory breakdown for each method on a Stable Diffusion personalization with a batch size of 1. ZOODiP (Ours) shows significantly higher memory efficiency compared to other methods. **(Right)** Comparison of memory usage versus performance across methods. Performance is measured with text (CLIP-T) and image (CLIP-I) alignment scores. ZOODiP achieves comparable performance to other methods while using significantly less memory (up to $8.2\times$ less than DreamBooth). Memory usage was profiled using the PyTorch profiler and `nvidia-smi` command.

marily update within this subspace. Third, based on prior works [2, 7, 10, 12, 19, 30, 34, 44, 68], which argue that timesteps have distinct roles in diffusion models, we identified an effective timestep section for personalization.

Following the first observation, we trained token embeddings using ZO optimization in a quantized, non-differentiable model [3, 39, 66]. This method reduces memory usage for activations, weights, and computational costs on edge devices. Inspired by the second finding, we introduced Subspace Gradient (SG) to accelerate training by mitigating noisy gradients. SG projects out dimensions with noisy gradients based on parameter trajectory, improving performance. Based on the third observation, we applied Partial Uniform Timestep Sampling (PUTS) within targeted timestep sections, skipping less influential timesteps to maximize impact within fixed training iterations.

We comprehensively assessed ZOODiP using both qualitative and quantitative measures on DreamBooth [52] dataset. We measured performance across two key metrics: 1) text-image alignment and 2) reference-generated image alignment. Through GPU memory profiling during training, we demonstrate that ZOODiP utilizes up to $8.2\times$ **less memory** than existing methods, achieving similar image quality. This substantial reduction in memory footprint underscores ZOODiP's efficacy in enabling diffusion model personalization on memory-constrained devices.

Our contributions can be summarized as follows:

- We empirically identify that Textual Inversion tokens primarily optimize within a low-dimensional subspace, and that focusing on partial timesteps accelerates training.
- We introduce ZOODiP, a novel method that combines ZO optimization with a quantized model, subspace gra-

dient projection, and customized timestep sampling for efficient diffusion model personalization.
- We demonstrate that ZOODiP achieves competitive performance with significantly lower memory requirements on the DreamBooth dataset (see Fig. 1).

## 2. Related Works

**Diffusion model personalization.** Diffusion model personalization aims to adapt a pre-trained model to generate images of new, user-defined concepts using a small set of provided images. Textual Inversion [17] tackles personalization by optimizing a single token embedding that represents the target concept. DreamBooth [52] tunes the denoising U-Net, while Custom Diffusion [28] personalizes multiple concepts by adopting new tokens and adjusting the key and value matrices in cross-attention. $\mathcal{P}+$ [62] assigns unique text tokens to each U-Net stage, and TextBoost [42] introduces an augmentation token and employs SNR-based sampling for single-image personalization. However, these methods often require significant computational resources and memory, limiting their applicability on resource-constrained devices. Inspired by the efficiency of single-token optimization in Textual Inversion, ZOODiP introduces a new token to represent the target concept and optimizes it efficiently.

**Efficient fine-tuning.** Fine-tuning large models requires significant memory, motivating the memory-efficient training methods. Many approaches aim to reduce the number of trainable parameters. For example, LoRA [24] applies low-rank matrices to linear layers, effectively decreas-
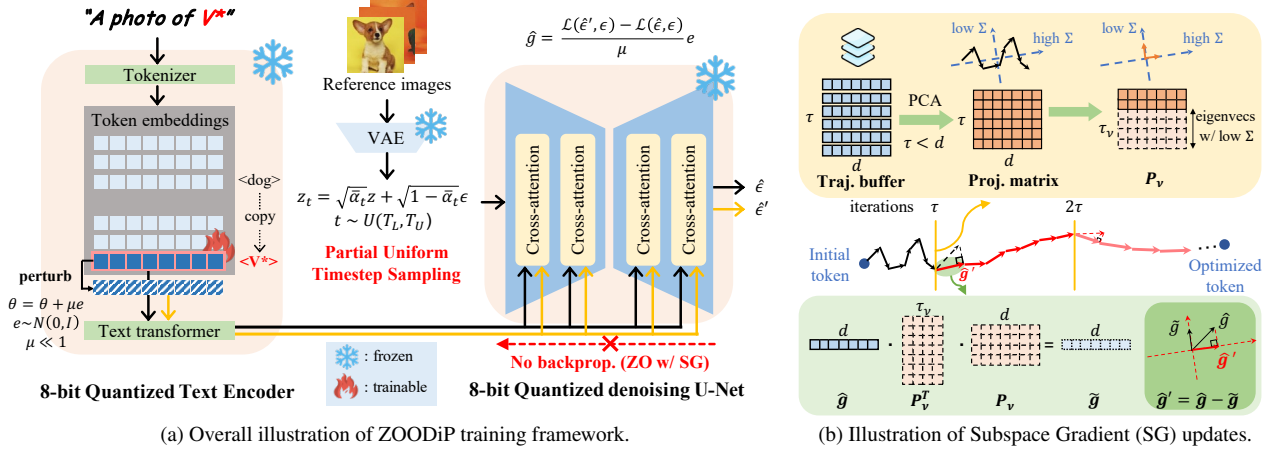
(a) Overall illustration of ZOODiP training framework.

(b) Illustration of Subspace Gradient (SG) updates.

Figure 2. **(a)** Illustration of overall ZOODiP framework. A target token is initialized and added to the prompt. Reference images are encoded, and Partial Uniform Timestep Sampling (PUTS)-sampled timestep noise is predicted. The loss is calculated with the original and perturbed token to estimate the gradient. **(b)** Illustration of Subspace Gradient (SG). Updated tokens from the previous $\tau$ iterations are stored. PCA identifies low-variance eigenvectors to project out noisy dimensions from the estimated gradient for the next $\tau$ iterations.

ing the number of parameters to update. QLoRA [13] further reduces memory usage by applying LoRA to quantized models. Another direction explores fine-tuning quantized models by adjusting quantization parameters. PEQA [26] only tunes the quantization scales to reduce the optimizer state memory. TuneQDM [54] personalizes quantized diffusion models by tuning the scales specific to each diffusion timestep set. Beyond parameter reduction and quantization, Gradient-Free Textual Inversion [16] employs an evolution strategy to optimize tokens, bypassing backpropagation.

However, these techniques can be limited by memory-intensive backpropagation or training instability. ZOODiP overcomes these limitations by utilizing zeroth-order optimization on a quantized model, eliminating backpropagation and its associated memory overhead.

**Zeroth-order optimization.** As model sizes increased, memory consumption during backpropagation became a major challenge for fine-tuning. Zeroth-order (ZO) optimization [5, 6, 14, 18, 36, 40, 60], which estimates gradients using only forward passes with random perturbations, has emerged as a promising solution. MeZO [39] demonstrated memory efficient tuning of large language models with forward-pass only, storing random seeds instead of perturbation vectors. Additionally, MeZO showed that ZO's convergence rate depends on the model's effective rank and that non-differentiable objectives can be optimized with ZO. Similarly, DeepZero [8] trained models from scratch using parallelized ZO, achieving performance comparable to first-order methods. Inspired by these advancements, we propose a personalization framework that employs ZO optimization on a quantized model, significantly reducing memory overhead for both activations and weights.

## 3. Method

In this section, we propose ZOODiP, a memory-efficient approach to personalize diffusion models. It leverages zeroth-order (ZO) optimization with a quantized model, eliminating the need for backpropagation, and thereby significantly reducing memory usage. Based on the observation that trained tokens in Textual Inversion primarily change within a low-dimensional subspace, we introduce Subspace Gradient (SG) to optimize tokens within this reduced space. Furthermore, we propose Partial Uniform Timestep Sampling (PUTS) for efficient timestep selection, capitalizing on the observation that text embeddings predominantly influence image generation at specific diffusion timesteps. By integrating ZO optimization with quantization, SG, and PUTS, ZOODiP enables efficient personalization on resource-constrained devices. The overall framework is illustrated in Fig. 2 and Algorithm. 1.

### 3.1. Formulation

ZOODiP builds upon Textual Inversion, aiming to learn a pseudo-token $v^*$ that represents a desired concept from a set of reference images $x_1, ..., x_n \in \mathcal{X}$ containing that concept. This pseudo-token is incorporated into the model's text embedding. We denote the condition for denoising network $y^*$ as the text prompt that includes $v^*$. Following the previous studies [17, 51], we formulate the objective as follows:

$$\mathcal{L}_{\text{LDM}} = \mathbb{E}_{z \sim \mathcal{E}(x), y^*, \epsilon \sim \mathcal{N}(0,I), t} \left[ ||\epsilon - \epsilon_\phi(z_t, t, c(y^*))||_2^2 \right],$$
$$z_t = \sqrt{\bar{\alpha}_t} z + \sqrt{1 - \bar{\alpha}_t} \epsilon$$
$$(1)$$

where $x$ is the input reference image, $\mathcal{E}$ is a Variational Autoencoder (VAE) encoder, $z$ is the latent encoded with the

3

**Algorithm 1** Fine-tuning algorithm of ZOODiP

**Require:** token embedding $\theta$, reference image set $\mathcal{X}$, text input $y^*$, text encoder $c$, VAE encoder $\mathcal{E}$, denoising network $\epsilon_\phi$, buffer size $\tau$, threshold $\nu$, number of estimation $n$, perturbation size $\mu$, total iteration $L$, trajectory buffer $B$, learning rate $\eta$
1: $P_\nu \leftarrow \mathbf{0}$
2: **for** $l = 1, \ldots, L$ **do**
3:      Sample image $x \in \mathcal{X}$
4:      $t \sim U(T_L, T_U)$            ▷ PUTS
5:      $\hat{g}_\theta \leftarrow \text{RGE}(x, y^*, n, t, \mathcal{E}, c, \epsilon_\phi, \mu)$ ▷ Eq. 1 and Eq. 3
6:      $\hat{g}'_\theta \leftarrow \hat{g}_\theta(I - P_\nu^\top P_\nu)$        ▷ SG update
7:      $\theta \leftarrow \text{optimizer}(\theta, \hat{g}'_\theta, \eta)$
8:      $B_{\text{mod}(l,\tau)} \leftarrow \theta$ ▷ Fill the buffer with updated token
9:      **if** $\text{mod}(l, \tau) = 0$ **then**     ▷ Update $P_\nu$ for SG
10:         $P_\nu \leftarrow \text{update} P_\nu(B, \nu)$      ▷ Algorithm. 2
11:      **end if**
12: **end for**

---

**Algorithm 2** Subspace generation

1: **function** UPDATE$P_\nu(B, \nu)$
2:      $\bar{B} \leftarrow \text{Normalize}(B)$         ▷ Eq. 4
3:      $\lambda, V \leftarrow \text{PCA}(\bar{B})$
4:      $i^* \leftarrow \arg\min_i \left\{ i \in \{1 : \tau\} \mid \frac{\sum_{j=1}^i \lambda_j}{\sum_{k=1}^\tau \lambda_k} > 1 - \nu \right\}$
5:      $P_\nu \leftarrow V_{i^*:\tau}^\top$
6:      **return** $P_\nu$
7: **end function**

---

encoder, $\epsilon$ is the noise from the forward process, $\epsilon_\phi$ is a denoising network, $t$ is the diffusion timestep, $c$ is a condition encoder which is text encoder in ZOODiP, and $\bar{\alpha}_t$ is cumulative product of $\alpha$ from 0 to $t$ as defined in DDPM [21].

### 3.2. ZO Optimization with Quantized Model

To optimize the $v^*$ with memory efficiency, we quantize the weights in the Linear and Convolution layer of all components in the diffusion model: U-Net, VAE, and the text encoder. The quantization is formulated as follows [27]:

$$\widetilde{W} = s \cdot \left( \text{clamp}\left( \left\lfloor \frac{W}{s} \right\rceil + o, Q_N, Q_P \right) - o \right) \quad (2)$$

where $W$ is a weight to be quantized, $\lfloor \cdot \rceil$, $s$, $o$ represent the rounding function, quantization scale, and zero-point, respectively. The minimum and maximum quantization bounds, $Q_N = -2^{N-1}$ and $Q_P = 2^{N-1} - 1$, are configured to ensure symmetric $N$-bit quantization. After the quantization of the network, $\mathcal{L}_{\text{LDM}}$ in Eq. 1 becomes nondifferentiable due to the existence of the rounding function $\lfloor \cdot \rceil$, which makes training a quantized model difficult [3].

To estimate the gradient on a quantized model, we leverage Random Gradient Estimation (RGE) [6, 8, 40, 60].
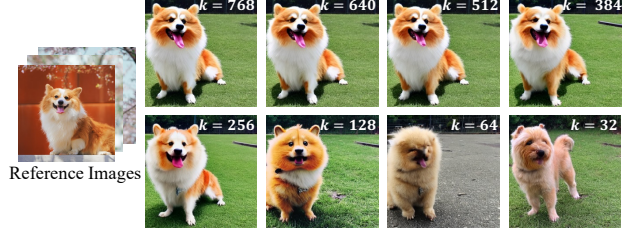


Figure 3. Sparse effective dimension in the token trained with Textual Inversion. Notably, the concept was preserved even when retaining only one-third of the optimized dimensions ($k = 256$).

RGE perturbs $\theta \in \mathbb{R}^d$, the token embedding for $v^*$, with random variables and estimates the gradient by calculating the non-differentiable objective $\mathcal{L}_{\text{LDM}}(\theta)$ as follows:

$$\hat{g}_\theta = \frac{1}{n} \sum_{i=1}^n \left[ \frac{\mathcal{L}_{\text{LDM}}(\theta + \mu e_i) - \mathcal{L}_{\text{LDM}}(\theta)}{\mu} e_i \right] \quad (3)$$

where $\hat{g}_\theta$ denotes an estimation of First-Order (FO) gradient $\nabla_\theta \mathcal{L}(\theta)$ with respect to $\theta$, $\mu$ is a perturbation size, $n$ is the number of random directions to estimate the gradient and $e_i$ a random vector sampled from $\mathcal{N}(0, I)$. With the ZO approach, we were able to estimate the gradient without backpropagation while effectively reducing memory usage.

### 3.3. Subspace Gradient

We investigated the optimization of the pseudo-token by Textual Inversion [17](TI) using Principal Component Analysis (PCA) of the token embeddings. Both the initial and optimized tokens were projected onto the principal components. The components with the top-$k$ largest changes in the optimized token were retained, while the rest were replaced with components from the initial token before back-projection. Fig. 3 shows that the personalized concept is preserved after back-projection, indicating that key changes from TI lie within this subspace.

ZO optimization suffers from slow training due to noisy gradient estimates. Inspired by the fact that token optimization mainly occurs in a lower-dimensional subspace, we propose Subspace Gradient (SG) to accelerate training. SG leverages token trajectories to eliminate noisy dimensions. Specifically, we store updated tokens for $\tau$ iterations in a trajectory buffer $B \in \mathbb{R}^{\tau \times d}$, where $d > \tau$. Once $B$ is full, we normalize it as follows:

$$\bar{B} = (B - \mu_B)/\sigma_B \quad (4)$$

where $\mu_B$ and $\sigma_B$ is the mean and standard deviation of $B$ along each feature dimension, respectively. Then we perform PCA on $\bar{B}$ with singular value decomposition to get eigenvalues and eigenvectors of the covariance matrix:

$$\bar{B} = U\Sigma V^\top, \quad \lambda_i = \Sigma_{ii}^2 \quad (5)$$

The square of $i$-th diagonal elements of $\Sigma$, denoted as $\lambda_i$, represent the variance explained by the corresponding eigenvectors. We calculate the ratio of the cumulative sum of $\lambda_i$ to the total sum:

$$i^* = \arg\min_i \left\{ i \in \{1, \ldots, \tau\} \,\middle|\, \frac{\sum_{j=1}^i \lambda_j}{\sum_{k=1}^\tau \lambda_k} > 1 - \nu \right\} \quad (6)$$

where $\nu$ is a hyperparameter that controls the amount of variance retained. This determines $i^*$, the smallest index $i^*$ for which the cumulative ratio exceeds $1 - \nu$. We then construct a projection matrix $P_\nu$ using the eigenvectors corresponding to the remaining variance:

$$P_\nu := V_{(i^*+1):\tau}^\top \quad (7)$$

This matrix $P_\nu$ represents the dimensions to be removed. We project the estimated row vector gradient $\hat{g}$ onto the subspace orthogonal to $P_\nu$ and subtract from $\hat{g}$ to get $\hat{g}'$:

$$\hat{g}' = \hat{g}(I - P_\nu^\top P_\nu) \quad (8)$$

This subtraction effectively eliminates noisy components from the $\hat{g}'$. After updating $P_\nu$, the buffer $B$ is cleared. Over the next $\tau$ iterations, the estimated gradient is projected out through $P_\nu$, the token trajectory is simultaneously accumulated in $B$. Fig. 2b, Algorithm. 1, and Algorithm. 2 represents the SG process as described above.

### 3.4. Partial Uniform Timestep Sampling

To facilitate more efficient training, we introduce Partial Uniform Timestep Sampling (PUTS). The diffusion model can be interpreted as a mixture-of-experts based on the timesteps [2, 19, 30, 34], with each timestep playing a distinct role. In text-to-image diffusion models, several studies [7, 10, 12, 44, 68] have shown that text conditioning impact varies across timesteps. However, most works focus on inference or training from scratch, while using this insight for personalization remains underexplored.

TextBoost [42] empirically observed that text influence increases as the timestep nears noise, proposing SNR-based timestep sampling. However, our experiments showed text influence is negligible at certain timesteps (see Fig. 4). Specifically, sampling timestep $t$ from $U(0, 500)$ failed to effectively learn the reference image concept, whereas sampling from $U(500, 1000)$ led to successful learning.

Based on these observations, we propose PUTS, which focuses on uniformly sampling timesteps within a specific range of the diffusion process. This range is chosen to prioritize timesteps where the text embedding has the most significant influence. PUTS can be formulated as follows:

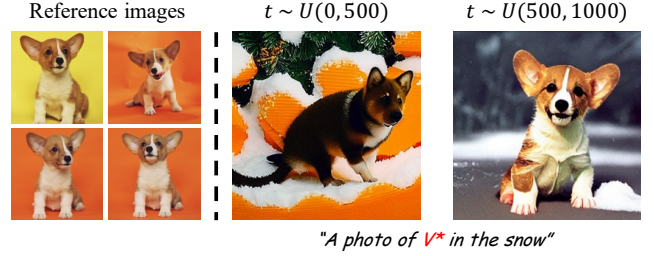$$z_t = \sqrt{\bar{\alpha}_t}z + \sqrt{1 - \bar{\alpha}_t}\epsilon, \quad t \sim U(T_L, T_U) \quad (9)$$



Reference images    $t \sim U(0, 500)$    $t \sim U(500, 1000)$

"A photo of V* in the snow"

Figure 4. Textual Inversion [17] with various timestep sampling. When the timestep $t$ for training is sampled from $U(0, 500)$, key conceptual features such as color and body shape of the reference image are not effectively trained. In contrast, sampling from $U(500, 1000)$ results in successful learning of these features.

where $\epsilon$ represents the noise sampled from Gaussian distribution, $\bar{\alpha}_t$ is the cumulative product of $\alpha$ values to sampled timestep $t$, and $T_L$ and $T_U$ denote the lower and upper bound of the partial timestep range. Our dense ablation study in Tab. 5 and Fig. 7 further display the effectiveness of the proposed method and validate the observations.

## 4. Experiments

To evaluate ZOODiP, we conducted quantitative and qualitative comparisons with methods based on DreamBooth (DB) [13, 26, 52, 54] and Textual Inversion (TI) [16, 17]. Personalization used the DB dataset and Stable Diffusion-v1.5 [51], with INT8 quantization applied to Linear and Conv2D layers. All experiments ran on a single Nvidia RTX 3090 GPU with batch size 1. ZOODiP was trained with $n = 2$, $\mu = 10^{-3}$, $\tau = 128$, $\nu = 10^{-3}$, $T_L = 500$, $T_U = 900$, $L = 30,000$, $\eta = 5 \times 10^{-3}$ using the ZOAdam [39] optimizer. Further results and experimental details are provided in the supplementary materials.

### 4.1. Quantitative Results

**Image and text alignment score.** To evaluate ZOODiP's personalization performance, we assessed text and image alignment against baselines. Using the DB dataset, we personalized 30 subjects with 25 prompts, generating 5 images per prompt for a total of 3,750 images. Image alignment was measured via cosine similarity with CLIP [49] (CLIP-I) and DINOv2 [41] (DINO) embeddings. Text alignment was measured using cosine similarity between CLIP text embeddings of prompts and generated image embeddings (CLIP-T). Tab. 1 shows ZOODiP achieves similar performance to prior methods, outperforming GF-TI by 43.0% on CLIP-I and 13.4% on CLIP-T with the same memory usage. Compared to TuneQDM, a state-of-the-art quantized diffusion personalization method, ZOODiP shows a 0.7% drop on CLIP-T and a 0.5% increase on DINO. Compared to our baseline, TI, performance differences were +0.7%, -0.8%,

| Base. | Method | Quant. | Grad. Free | Mem.↓ (GB) | Stor.↓ (MB) | CLIP-T↑ | CLIP-I↑ | DINO↑ |
|---|---|---|---|---|---|---|---|---|
| DB | DB [52] | ✗ | ✗ | 19.4 | 3438 | 0.281 | 0.782 | 0.592 |
|  | QLoRA [13] | ✓ | ✗ | 7.56 | 1.63 | 0.297 | <u>0.762</u> | 0.607 |
|  | PEQA [26] | ✓ | ✗ | 6.31 | 1.32 | <u>0.275</u> | 0.791 | 0.604 |
|  | TuneQDM [54] | ✓ | ✗ | 8.96 | 2.48 | 0.289 | 0.788 | <u>0.555</u> |
| TI | TI [17] | ✗ | ✗ | 6.75 | **0.003** | 0.285 | 0.778 | 0.559 |
|  | GF-TI [16] | ✓ | ✓ | **2.37** | **0.003** | <u>0.253</u> | <u>0.540</u> | <u>0.011</u> |
|  | **ZOODiP (Ours)** | ✓ | ✓ | **2.37** | **0.003** | 0.287 | 0.772 | 0.558 |

Table 1. Quantitative comparisons of DreamBooth [52] (DB), QLoRA [13] ($r = 2$), PEQA [26], TuneQDM [54], Textual Inversion [17] (TI), Gradient-Free Textual Inversion [16] (GF-TI), and Ours. ↑ / ↓ indicates higher / lower values are better. Performance was evaluated with CLIP-I and DINO for image alignment, CLIP-T for text-image alignment, and memory requirements of training (Mem.) and storage (Stor.). The <u>worst</u>-performance is double-underlined, and the <u>second worst</u> is single-underlined. ZOODiP achieves performance comparable to that of gradient-based methods with significantly less memory.

and -0.2% on CLIP-T, CLIP-I, and DINO, respectively.

**Memory efficiency.** To assess the training memory efficiency of each method, we tracked peak memory usage using the `nvidia-smi` command after quantization. As shown in Tab. 1, ZOODiP requires significantly less memory during training—87.8% reduction compared to DB and 64.9% reduction compared to TI, which serves as a baseline for ZOODiP. We also evaluated storage requirements for the optimized models in `safetensors` format. ZOODiP requires only 3KB of storage, similar to other TI-based methods, making it well-suited for edge devices.

**Training speed.** The training speed of personalization methods is shown in Tab. 2. Since ZOODiP relies solely on forward passes, avoiding costly backpropagation, it achieves significantly faster training speeds compared to other backpropagation-based methods. Specifically, ZOODiP is $2.2\times$ and $1.7\times$ faster than TI for $n = 1$ and $n = 2$, respectively, and $4.2\times$ and $3.3\times$ faster than TuneQDM. GF-TI is the slowest, with ZOODiP showing a $28\times$ and $22\times$ speed improvement due to GF-TI's requirement for 30 forward passes per iteration. Additionally, ZOODiP with FP16 precision, fully supported by the hardware, further enhances training speed.

## 4.2. Qualitative Results

Fig. 5 shows qualitative comparisons between ZOODiP and other methods. ZOODiP generates images highly faithful to both prompts and reference images, with minimal distinction from other advanced techniques. In contrast, GF-TI, under similar memory constraints, struggles to maintain fidelity, resulting in noticeably inaccurate depictions. Fig. 6 shows the results of ZOODiP's style personalization. Using several photos with consistent styles as references, ZOODiP demonstrates accurate reflection of the learned style for the

| Method | Prec. | Speed (iter/sec) ↑ |
|---|---|---|
| TI | FP32 | 9.42 |
| Ours ($n = 1$) | FP16 | **22.3** |
| Ours ($n = 2$) |  | 18.4 |
| TuneQDM |  | 4.94 |
| GF-TI | INT8 | 0.74 |
| Ours ($n = 1$) |  | **20.7** |
| Ours ($n = 2$) |  | 16.1 |

Table 2. Training speed comparisons of prior works and ZOODiP. ZOODiP achieves the fastest training speed by estimating gradients with only forward passes, bypassing backpropagation.

given prompts, indicating that ZOODiP has successfully inherited the style personalization capabilities of TI.

## 4.3. Ablation Study

To analyze the contribution of each component in ZOODiP and assess the impact of hyperparameter variations, we conducted an ablation study. Experiments in this section were performed with <dog6> and <shiny_sneaker> data from the DreamBooth [52] dataset.

**Effectiveness of SG and PUTS.** To measure the impact of ZOODiP's key elements—SG and PUTS—we evaluated its performance by incrementally adding each component, as shown in Tab. 3. The results demonstrate that both SG and PUTS significantly enhance performance, outperforming naïve Textual Inversion with ZO optimization.

**Hyperparameter study.** We conducted exhaustive experiments to analyze the impact of two key hyperparameters used in SG: $\tau$ and $\nu$. As shown in Tab. 4, ZOODiP shows robust performance unless $\tau$ or $\nu$ is too large or small. When $\tau$ is small, frequent $P_\nu$ updates may occur, introduce com-
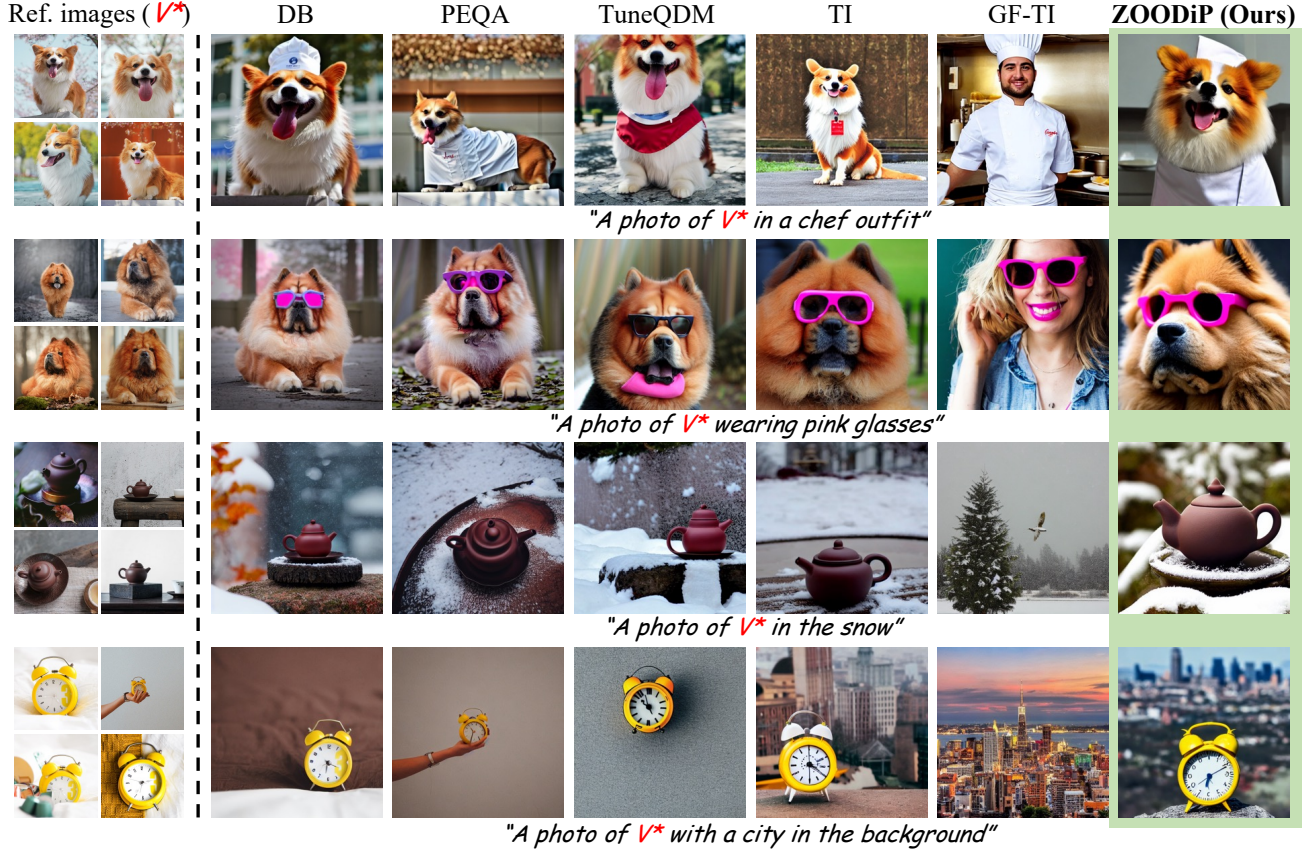
Figure 5. **Qualitative comparison of image and text alignment.** This figure shows how well each method generates images that match the input text prompt while preserving the identity of the personalized subject. ZOODiP generates images that faithfully reflect the prompt while maintaining the concept of the reference image, demonstrating strong image-text alignment.
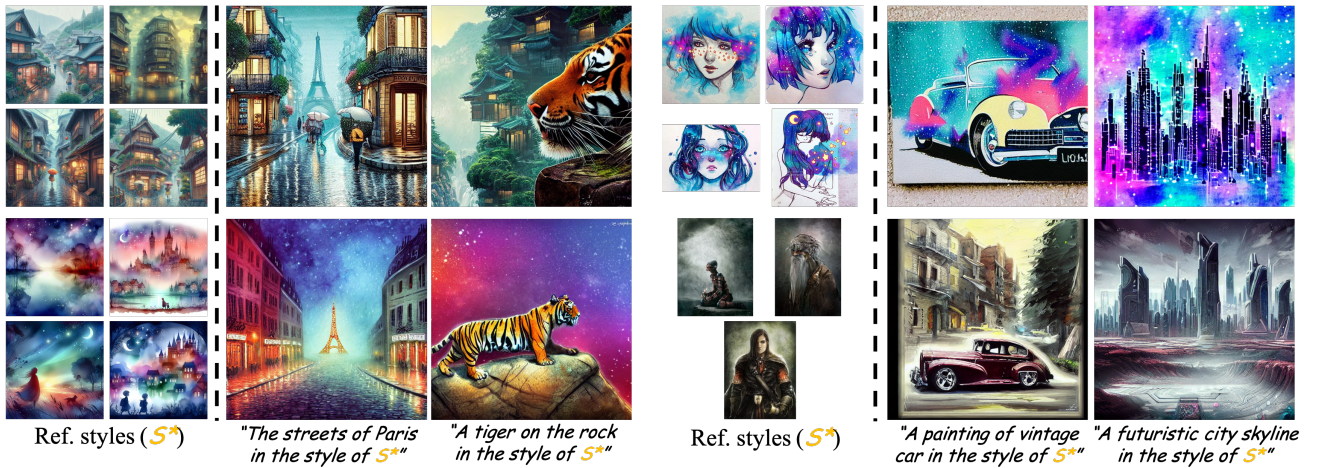


Figure 6. **Qualitative results on style personalization.** This figure showcases the results of style personalization achieved through ZOODiP, using few reference images with a consistent style. The outcome highlights ability of ZOODiP to personalize not only the subject but also the style with a high degree of accuracy. This demonstrates the versatility and extensive personalization capabilities of ZOODiP, effectively adapting both stylistic elements and subject details to match the reference images.

| SG | PUTS | CLIP-T↑ | CLIP-I↑ | DINO↑ |
|----|------|---------|---------|-------|
| ✗ | ✗ | 0.273 | 0.736 | 0.505 |
| ✓ | ✗ | 0.265 | 0.747 | 0.527 |
| ✗ | ✓ | **0.277** | 0.744 | 0.562 |
| ✓ | ✓ | 0.266 | **0.759** | **0.569** |

Table 3. Ablation study on ZOODiP components with `<shiny_sneaker>`. ✓ denotes the component is applied, while ✗ means it is not. Without PUTS, timesteps are sampled uniformly.

| $\tau$ \ $\nu$ | $10^{-1}$ | $10^{-2}$ | $10^{-3}$ | $10^{-4}$ |
|------|-----------|-----------|-----------|-----------|
| 32 | 0.704 | 0.686 | 0.716 | 0.729 |
| 64 | 0.721 | 0.736 | 0.712 | 0.724 |
| 128 | 0.736 | 0.735 | **0.759** | 0.716 |
| 256 | 0.739 | 0.727 | 0.738 | 0.716 |
| 512 | 0.705 | 0.704 | 0.676 | 0.707 |

Table 4. Ablation study on hyperparameters $\tau$ and $\nu$ incorporated with SG. We optimized the pseudo-token with various $\tau$ and $\nu$ and measured the performance with the CLIP-I score. Experiments were done with `<shiny_sneaker>` dataset.

| Method | CLIP-T↑ | CLIP-I↑ | DINO↑ |
|--------|---------|---------|-------|
| Uniform | 0.265 | 0.747 | 0.527 |
| SNR-based | **0.271** | 0.719 | 0.545 |
| **PUTS (Ours)** | 0.266 | **0.759** | **0.569** |

Table 5. Ablation study on various diffusion timestep sampling method with `<shiny_sneaker>`. PUTS outperforms in image alignment score among all sample methods with minor degradation in text alignment score compared to SNR-based sampling.

putational overhead, so we choose values of $\tau = 128$ and $\nu = 10^{-3}$ which show the best performance.

**Diffusion timestep sampling.** To assess the effectiveness of PUTS on performance, we conducted experiments with different diffusion timestep sampling strategies: uniform sampling, SNR-based sampling [42], and PUTS. The results in Tab. 5 confirm that PUTS outperforms the other sampling methods, validating its effectiveness.

Furthermore, we conducted exhaustive experiments on various combinations of $T_U$ and $T_L$. As shown in Fig. 7, we divided the diffusion timesteps into 10 units and tested all possible combinations. The results indicate that training with timesteps closer to the noise improves both image and text fidelity, supporting our choice of $T_L$ and $T_U$ values.
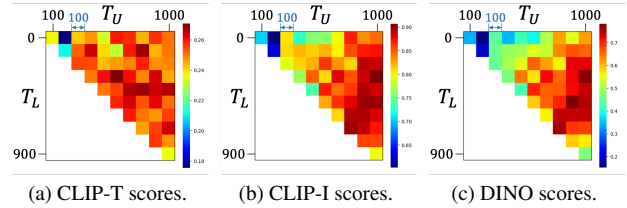


(a) CLIP-T scores. (b) CLIP-I scores. (c) DINO scores.

Figure 7. Heatmap of CLIP-T, CLIP-I and DINO scores across varying $T_L$ and $T_U$ on the `<dog6>` dataset. $x$-axis is the $T_U$ and $y$-axis is the $T_L$ applied to the sampling distribution.
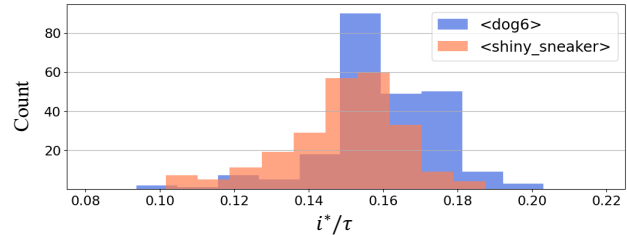


Figure 8. Histogram of $i^*/\tau$ ratios for `<dog6>` and `<shiny_sneaker>` dataset with hyperparameter $\tau = 128$, $\nu = 10^{-3}$ during training with SG. Despite the small $\nu$, a significant portion ($> 80\%$) of dimensions are projected out.

## 5. Discussion

ZO optimization's convergence rate is influenced by the effective dimension of the optimized parameters [33, 39, 67]. We conjecture that ZOODiP's success arises from the low dimensionality of the optimized token and its even smaller effective dimension (Sec. 3.3). The concentrated variance of the token trajectory allows SG to effectively project out over two-thirds of the dimensions in the trajectory buffer, even with a small $\nu$ value. It aligns with Fig. 3 and the analysis in Sec. 3.3. We selected $\nu = 10^{-3}$ to concentrate variance ratios within one-third of the total (Fig. 8).

## 6. Conclusion

We propose ZOODiP, a framework that enables the personalization of diffusion models even in highly memory-constrained environments. ZOODiP leverages quantized models during training to reduce memory footprint. We then propose zeroth-order optimization to personalize these models. To overcome the limitations of zeroth-order optimization, we introduce Subspace Gradient (SG), which utilizes the trajectory of the optimizing token to eliminate noisy gradient dimensions and enhance performance based on the empirical analysis of the effective dimension of personalized tokens. Additionally, we systematically analyzed the impact of diffusion timesteps on personalization and identified the most effective timesteps for efficient training. Based on this analysis, we proposed Partial Uniform

Timestep Sampling (PUTS), which samples only these relevant timesteps. Through extensive qualitative and quantitative experiments, we demonstrate that ZOODiP achieves performance comparable to prior arts while requiring significantly less memory of only 2.37GB during training.

## Acknowledgements

## Appendix

## S1. Additional Experimental Details

### S1.1. Dataset

All quantitative experiments were conducted using the DreamBooth [52] (DB) dataset, which includes 30 distinct subjects, which are denoted in Tab. S1, each paired with 25 unique prompts. Among these subjects, nine are living entities—specifically dogs and cats—while the remaining 21 represent non-living objects. Details of the subjects and their respective prompts are provided in the Tab. S2 (living) and Tab. S3 (non-living). Notably, the Textual Inversion [17]-based method did not incorporate class tokens in the prompts used for image generation, ensuring a fair comparison by excluding scenarios where class tokens are utilized during the unique token learning process. For Textual Inversion [17] (TI) and Gradient-Free Textual Inversion [16] (GF-TI), we employed the ImageNet-based template proposed in the original TI paper for training. ZOODiP followed the DCO [29], utilizing Comprehensive Captioning (CC) for the text prompt.

### S1.2. Metrics

To compute the CLIP-I score, we use the `openai/clip-vit-base-patch32` model [49] from Huggingface to extract image features for both the reference and generated images. We then calculate the cosine similarity for all possible pairs and average these values. For the CLIP-T score, we leverage the same model's text encoder to extract text features from the input prompt. We also calculate image features for the generated images and measure the pairwise cosine similarity between the text and image features, averaging the results to obtain the final score. To compute the DINO score, we utilize the `facebook/dinov2-base` model from Huggingface to extract DINOv2 [41] embeddings for both the reference and generated images. The score is determined by averaging the pairwise cosine similarities between these embeddings.

### S1.3. Quantization

We utilized `optimum-quanto`, the PyTorch [43] quantization backend provided by Huggingface, to enable accelerated matrix multiplications on CUDA devices. For the 8-bit weight quantization described in the main paper, we employed `optimum-quanto` to quantize the weights of all Linear and Conv2D layers in the VAE, U-Net, and text encoder modules to 8-bit integer. Other layers' parameters including text embedding, layer normalization, and batch normalization were not quantized. With this process, 96.4% of the whole Stable Diffusion [51] pipeline (U-Net, VAE, and text encoder) parameters are quantized to INT8, while the other 3.6% parameters remain FP16.

### S1.4. Memory usage

We primarily monitored memory usage using the `nvidia-smi` command to observe real-time memory consumption during training and employed the PyTorch profiler for detailed memory breakdowns. Notably, the memory utilized by the CUDA context can vary based on the CUDA and PyTorch versions, so the reported memory usage may differ depending on the experimental environment. Our measurements were conducted under the following settings: CUDA Toolkit 11.8, Torch 2.4.1, Torchvision 0.19.1, and Python 3.8.10. Since PyTorch loads all CUDA libraries by default, including unused ones, the actual memory usage could potentially be reduced by unloading unnecessary libraries. However, as our focus was not on such optimizations, we ensured a fair comparison by evaluating all training methods under the same experimental environment. We also observed that Variational Auto-Encoder (VAE) encoding introduces significant memory overhead. Given that all methods do not require backpropagation through the VAE, we standardized the process by blocking gradients in the VAE encoding step across all methods. This adjustment ensured that each personalization method used the minimal memory required, facilitating a fair evaluation of memory usage.

### S1.5. Baseline training configuration

For our experiments, we utilized the unofficial implementations provided by Huggingface for Textual Inversion [17] and DreamBooth [52] without prior-preservation loss which is suitable to memory-constrained environment. For PEQA [26] and TuneQDM [54], we conducted experiments using reproduced code based on the Huggingface implementation of DreamBooth. QLoRA [13] was imple-

**Subjects in DreamBooth [52] dataset**

backpack, backpack_dog, bear_plushie, berry_bowl, can, candle, cat, cat2, clock, colorful_sneaker, dog, dog2, dog3, dog5, dog6, dog7, dog8, duck_toy, fancy_boot, grey_sloth_plushie, monster_toy, pink_sunglasses, poop_emoji, rc_car, red_cartoon, robot_toy, shiny_sneaker, teapot, vase, wolf_plushie

Table S1. Full subjects name of DreamBooth dataset. The dataset names in the main paper are based on the corresponding subject datasets.

**Full prompt used in evaluation (living)**

```
1   'a {0} {1} in the jungle'.format(unique_token, class_token)
2   'a {0} {1} in the snow'.format(unique_token, class_token)
3   'a {0} {1} on the beach'.format(unique_token, class_token)
4   'a {0} {1} on a cobblestone street'.format(unique_token, class_token)
5   'a {0} {1} on top of pink fabric'.format(unique_token, class_token)
6   'a {0} {1} on top of a wooden floor'.format(unique_token, class_token)
7   'a {0} {1} with a city in the background'.format(unique_token, class_token)
8   'a {0} {1} with a mountain in the background'.format(unique_token, class_token)
9   'a {0} {1} with a blue house in the background'.format(unique_token, class_token)
10  'a {0} {1} on top of a purple rug in a forest'.format(unique_token, class_token)
11  'a {0} {1} wearing a red hat'.format(unique_token, class_token)
12  'a {0} {1} wearing a santa hat'.format(unique_token, class_token)
13  'a {0} {1} wearing a rainbow scarf'.format(unique_token, class_token)
14  'a {0} {1} wearing a black top hat and a monocle'.format(unique_token, class_token)
15  'a {0} {1} in a chef outfit'.format(unique_token, class_token)
16  'a {0} {1} in a firefighter outfit'.format(unique_token, class_token)
17  'a {0} {1} in a police outfit'.format(unique_token, class_token)
18  'a {0} {1} wearing pink glasses'.format(unique_token, class_token)
19  'a {0} {1} wearing a yellow shirt'.format(unique_token, class_token)
20  'a {0} {1} in a purple wizard outfit'.format(unique_token, class_token)
21  'a red {0} {1}'.format(unique_token, class_token)
22  'a purple {0} {1}'.format(unique_token, class_token)
23  'a shiny {0} {1}'.format(unique_token, class_token)
24  'a wet {0} {1}'.format(unique_token, class_token)
25  'a cube shaped {0} {1}'.format(unique_token, class_token)
```

Table S2. Full prompts used in evaluation of living category objects. `unique_token` represents the special token corresponds to object which aims to personalize, and `class_token` denotes the class that `unique_token` is in.

**Full prompt used in evaluation (non-living)**

```
1   'a {0} {1} in the jungle'.format(unique_token, class_token)
2   'a {0} {1} in the snow'.format(unique_token, class_token)
3   'a {0} {1} on the beach'.format(unique_token, class_token)
4   'a {0} {1} on a cobblestone street'.format(unique_token, class_token)
5   'a {0} {1} on top of pink fabric'.format(unique_token, class_token)
6   'a {0} {1} on top of a wooden floor'.format(unique_token, class_token)
7   'a {0} {1} with a city in the background'.format(unique_token, class_token)
8   'a {0} {1} with a mountain in the background'.format(unique_token, class_token)
9   'a {0} {1} with a blue house in the background'.format(unique_token, class_token)
10  'a {0} {1} on top of a purple rug in a forest'.format(unique_token, class_token)
11  'a {0} {1} with a wheat field in the background'.format(unique_token, class_token)
12  'a {0} {1} with a tree and autumn leaves in the background'.format(unique_token, class_token)
13  'a {0} {1} with the Eiffel Tower in the background'.format(unique_token, class_token)
14  'a {0} {1} floating on top of water.format(unique_token, class_token)
15  'a {0} {1} floating in an ocean of milk'.format(unique_token, class_token)
16  'a {0} {1} on top of green grass with sunflowers around it'.format(unique_token, class_token)
17  'a {0} {1} on top of a mirror'.format(unique_token, class_token)
18  'a {0} {1} on top of the sidewalk in a crowded street'.format(unique_token, class_token)
19  'a {0} {1} on top of a dirt road'.format(unique_token, class_token)
20  'a {0} {1} on top of a white rug'.format(unique_token, class_token)
21  'a red {0} {1}'.format(unique_token, class_token)
22  'a purple {0} {1}'.format(unique_token, class_token)
23  'a shiny {0} {1}'.format(unique_token, class_token)
24  'a wet {0} {1}'.format(unique_token, class_token)
25  'a cube shaped {0} {1}'.format(unique_token, class_token)
```

Table S3. Full prompts used in evaluation of non-living category objects. `unique_token` represents the special token corresponds to object which aims to personalize, and `class_token` denotes the class that `unique_token` is in.

| Method | CLIP-T↑ | CLIP-I↑ | DINO↑ |
|--------|---------|---------|-------|
| RGE | 0.266 | 0.759 | 0.569 |
| SPSA | 0.277 | 0.732 | 0.506 |
| One-point | 0.296 | 0.703 | 0.393 |

Table S4. Personalization performance across different gradient estimation methods with $n = 2$. The results show that RGE outperforms other two different gradient estimations methods. Notably, RGE is more efficient in terms of computational cost, requiring fewer forward passes than SPSA when $n > 1$. Due to this efficiency and performance advantages, we adopted RGE as the gradient estimation method in ZOODiP.

| $\mu$ | CLIP-T↑ | CLIP-I↑ | DINO↑ |
|-------|---------|---------|-------|
| $1 \times 10^{-2}$ | 0.281 | 0.778 | 0.599 |
| $\mathbf{1 \times 10^{-3}}$ | 0.277 | **0.797** | **0.613** |
| $1 \times 10^{-4}$ | **0.296** | 0.724 | 0.470 |

Table S5. Quantitative results for different $\mu$ values. Optimal performance is observed at $\mu = 10^{-3}$ with varying $\mu$, and we have set the value of $\mu$ accordingly for ZOODiP.

| $n$ | CLIP-T↑ | CLIP-I↑ | DINO↑ | Speed↑ (iter/s) |
|-----|---------|---------|-------|-----------------|
| 1 | 0.298 | 0.736 | 0.495 | 20.7 |
| 2 | 0.277 | 0.796 | 0.613 | 16.1 |
| 4 | 0.282 | 0.784 | 0.584 | 9.78 |
| 8 | 0.282 | 0.798 | 0.627 | 6.20 |

Table S6. Quantitative results for various $n$, the number of gradient estimations. $n = 2$ is the promising choice between performance and computation efficiency.

mented using the `BitsandBytes` library. For Gradient-Free Textual Inversion [16], we utilized the official implementation for our experiments. We configured the training settings for each method to be as consistent as possible, adjusting only the batch size to 1. For DreamBooth, we set $\eta = 5 \times 10^{-6}$, $L = 400$, for Textual Inversion, we set $\eta = 5 \times 10^{-3}$, $L = 5,000$, for QLoRA, we set $\eta = 1 \times 10^{-4}$, $L = 500$, for PEQA, we set $\eta = 3 \times 10^{-6}$, $L = 400$, for TuneQDM, we set $\eta = 3 \times 10^{-6}$, $L = 400$, for Gradient-Free Textual Inversion, $\eta = 5 \times 10^{6}$, $L = 13,000$, intrinsic dimension $d_i = 256$, $\sigma = 1$, and $\alpha = 1$.

## S2. Additional Ablation Study

We conducted an extensive ablation study to analyze the impact of each component of ZOODiP. Beyond evaluating RGE, we examined the effectiveness of alternative gradient estimation methods, the influence of varying perturbation scales on performance, the impact of changing the number of gradient estimations, and the resulting time overhead associated with these changes. Unless otherwise noted, all additional experiments were performed on two subjects: `<dog6>` as a representative of living objects and `<shiny_sneaker>` as a representative of non-living objects. The results from these experiments were averaged to ensure a balanced and comprehensive evaluation.

### S2.1. Gradient estimation method

There are various methods for estimating gradients, each with unique advantages and trade-offs. In ZOODiP, we opted for Random Gradient Estimation (RGE, Eq. S10) due to its efficiency in estimating gradients using Monte Carlo gradient estimation. While many existing works utilize Simultaneous Perturbation Stochastic Approximation (SPSA, Eq. S11), which perturbs the gradient in two directions, or a one-point estimation method (Eq. S12) that calculates the gradient directly at the perturbed point, we evaluated all three methods for gradient estimation as seen in Tab. S4.

$$\hat{g}_\theta = \frac{1}{n} \sum_{i=1}^{n} \left[ \frac{\mathcal{L}_{\text{LDM}}(\theta + \mu e_i) - \mathcal{L}_{\text{LDM}}(\theta)}{\mu} e_i \right] \quad (S10)$$

$$\hat{g}_\theta = \frac{1}{n} \sum_{i=1}^{n} \left[ \frac{\mathcal{L}_{\text{LDM}}(\theta + \mu e_i) - \mathcal{L}_{\text{LDM}}(\theta - \mu e_i)}{2\mu} e_i \right] \quad (S11)$$

$$\hat{g}_\theta = \frac{1}{n} \sum_{i=1}^{n} \left[ \frac{\mathcal{L}_{\text{LDM}}(\theta + \mu e_{i+1}) - \mathcal{L}_{\text{LDM}}(\theta + \mu e_i)}{\mu} e_i \right] \quad (S12)$$

Our results indicate that SPSA and RGE exhibit similar performance, with RGE occasionally outperforming SPSA. In contrast, the one-point estimation method performed poorly. From a computational perspective, SPSA requires $2n$ forward passes for $n$ gradient estimation steps, whereas RGE needs only $n + 1$ forward passes, making it more efficient. Given its favorable balance of computational efficiency and performance, we selected RGE as the gradient estimation method for ZOODiP.

### S2.2. Perturbation scale

We analyzed the impact of varying $\mu$, the scaling parameter (*a.k.a* smoothing parameter) for perturbations in Eq. S10, on personalization performance. In general, the optimal value of $\mu$ can vary depending on the configuration of the model to be tuned as seen in Sec. S3.2 and Sec. S3.4. The results in Tab. S5 demonstrate that the chosen value of $10^{-3}$ is the most suitable perturbation scale for our setup.

| Method | PUTS | CLIP-T↑ | CLIP-I↑ | DINO↑ |
|--------|------|---------|---------|-------|
| DB | ✗ | 0.266 | 0.853 | 0.751 |
|    | ✓ | **0.272** | **0.854** | **0.758** |
| TI | ✗ | 0.241 | 0.798 | 0.584 |
|    | ✓ | **0.247** | **0.825** | **0.661** |

Table S7. Quantitative results from applying PUTS to Dream-Booth and Textual Inversion. The results confirm that PUTS enhances the performance of gradient-based personalization methods by up to 13.2%. The improvement was particularly pronounced in Textual Inversion, which is highly influenced by the text encoder, highlighting the significant impact of PUTS in this context.

| U-Net Precision | CLIP-T↑ | CLIP-I↑ | DINO↑ |
|-----------------|---------|---------|-------|
| INT8 | **0.288** | 0.834 | **0.657** |
| INT4 | 0.212 | **0.835** | 0.647 |

Table S8. The performance comparison between INT8 and INT4 for U-Net precision on the `<dog6>` subset of the DB dataset.

## S2.3. Number of gradient estimation

According to MeZO [39], increasing $n$ when fine-tuning large language models provides only marginal performance gains compared to the proportional increase in the number of forward passes, making $n = 1$ the most efficient choice. Tab. S6 illustrates the personalization performance across different $n$ values with RGE. While increasing $n$ slightly improves performance due to more accurate gradient estimation, the associated increase in training time becomes a limiting factor. Thus, we set $n = 2$ as a compromise between performance and efficiency.

## S3. Additional Experiments

### S3.1. PUTS on gradient-based methods

Partial Uniform Timestep Sampling (PUTS), one of the key components of ZOODiP, not only enhances ZOODiP's performance but also proves effective in gradient-based approaches. To validate this, we conducted experiments using two representative gradient-based personalization methods on Stable Diffusion v1.5: Textual Inversion and DreamBooth. These experiments were performed with the same hyperparameter settings for $T_L$ and $T_U$ as those used in ZOODiP. The results in Tab. S7 demonstrate that applying Partial Uniform Timestep Sampling (PUTS) to gradient-based methods improves performance by up to 11.6%. PUTS enhances efficiency by prioritizing informative timesteps, showcasing its versatility in optimizing both zeroth-order and gradient-based approaches.



Figure S1. Generated images with the prompt *"a photo of a dog"* with various weight precision. While INT8 precision produces results nearly equivalent to full-precision performance, INT4 precision exhibits noticeable degradation in image quality, highlighting the trade-off between lower precision and fidelity.
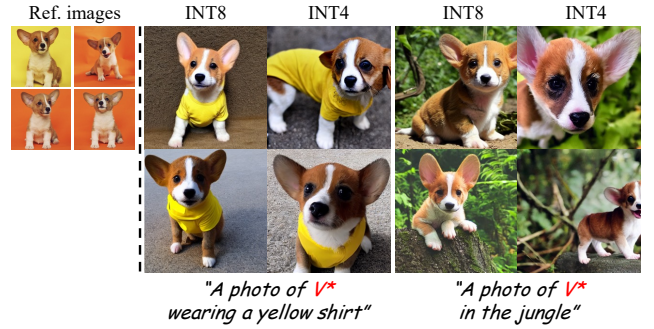


Figure S2. Qualitative results of U-Net precision at INT8 and INT4 in `<dog6>` dataset. ZOODiP works on INT4 and INT8, but performance diminishes due to degradation caused by INT4 quantization.

### S3.2. 4-bit quantized models

We applied ZOODiP to a 4-bit quantized model of the U-Net, the component with the highest VRAM usage. We observed that quantizing with INT4 precision led to a performance loss in the quantization library we utilized as seen in Fig. S1. Consequently, we also noted a degradation in the performance of the personalized results when using this quantization approach. Nevertheless, our results demonstrate that ZOODiP is fully capable of personalizing 4-bit quantized diffusion models using only **1.9 GB** of memory during the entire training process. Fig. S2 and Tab. S8 present the qualitative and quantitative results, respectively, for personalization with the 4-bit quantized model. To utilize the `qint4` data type, we employed BFloat16 (BF16) for activation in all units, including the VAE, U-Net, and text encoder. When using BF16, the default $\mu$ value of $10^{-3}$ was insufficient, so we adjusted it to $10^{-2}$ to achieve optimal performance.

### S3.3. Diversity of generated images

TI-based methods tend to exhibit less overfitting compared to DreamBooth (DB)-based methods, resulting in higher diversity in the generated images. To evaluate this tendency,
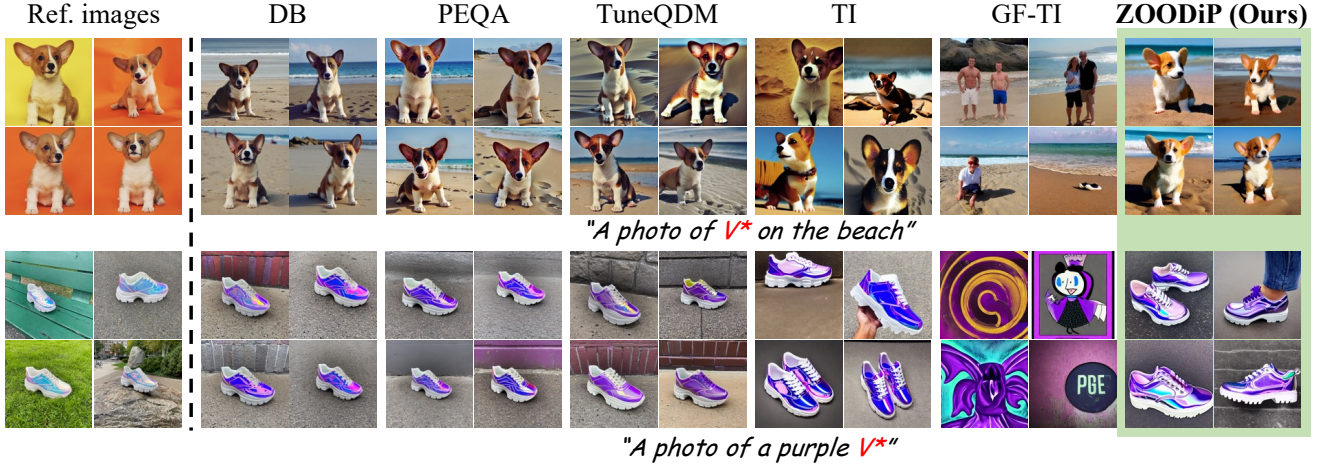
Figure S3. **Qualitative comparison of the diversity of generated images** This figure compares the diversity achieved by different personalization methods. ZOODiP demonstrates the ability to generate highly diverse images while utilizing minimal memory resources.
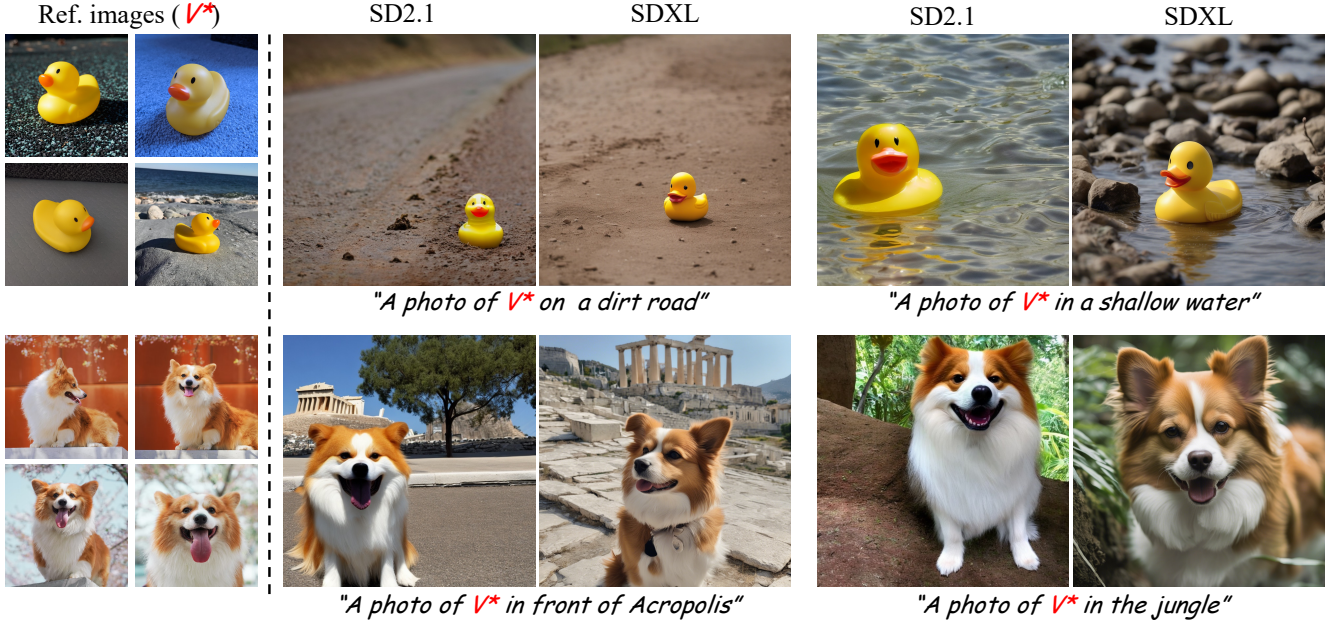


Figure S4. Qualitative results for personalizing SD2.1 and SDXL with ZOODiP. The figure demonstrate that ZOODiP can be applied not only to SD1.5, as discussed in the main paper, but also to various other models. For SD2.1, inference were conducted with images at a resolution of $768 \times 768$, while for SDXL, image generation was performed with resolution of $1024 \times 1024$. However, for SDXL, it was observed that the model's inherent color interpretation prevents the subject's colors from being completely replicated. This indicates that the model's color rendering can vary depending on the environmental context, leading to shifts in the perceived color scheme.

we introduce the $DINO_{inter}$ score. For each subject, we generated 50 images using the same prompt with a personalized model. The DINOv2 embedding cosine similarity was then calculated for all pairs of generated images and averaged. A higher $DINO_{inter}$ score indicates that the generated images are more similar to each other, reflecting lower diversity in the generated outputs. The quantitative results for $DINO_{inter}$ are presented in the Tab. S9. Additionally, we provide qualitative diversity results in Fig. S3. Fig. S3 shows a comparison of diversity across different methods, clearly showing that TI-based methods achieve significantly higher diversity than other approaches.

| Method | DINO$_{inter}$ ↓ |
|--------|------------------|
| DB | 0.825 |
| QLoRA | 0.731 |
| PEQA | 0.806 |
| TuneQDM | 0.778 |
| TI | 0.679 |
| GF-TI | 0.150 |
| **ZOODiP (Ours)** | 0.671 |

Table S9. Comparison of DINO$_{inter}$ scores across various personalization methods. The results indicate that TI-based methods are capable of generating a diverse range of images, whereas DB-based methods exhibit relatively lower diversity. This observation was consistent across all subjects in the DreamBooth (DB) dataset, highlighting a fundamental difference in the ability of these methods to capture and reflect variations in the generated outputs.

## S3.4. Generalizability to other models

In the main paper, our experiments were conducted using Stable Diffusion v1.5 (SD1.5). However, ZOODiP is not limited to the certain model and can be extended to other models trained on different datasets and architectures, such as SD2.1 and SDXL [46] as seen in Fig. S4. For SD2.1, training was performed on a larger text token embedding dimension (1024 in `OpenCLIP-ViT-H`) using the same hyperparameters as SD1.5 except for $\mu = 10^{-2}$ and $\eta = 10^{-2}$. For SDXL, we maintained all hyperparameters except for $\mu = 10^{-2}, T_L = 700, L = 20,000$ and trained both token embeddings of `OpenCLIP-ViT-L` (768 dimension) and `OpenCLIP-ViT-G` (1280 dimension). The adjustment of $T_L$ is supported by prior research which indicates that information loss vary with dimensionality [11, 32], and the value of $T_L$ was selected empirically. Notably, full-precision (FP32) Textual Inversion on SDXL requires approximately 17 GB of memory, and DreamBooth cannot be trained on customer level GPUs, such as the RTX3090. In contrast, tuning SDXL with ZOODiP enables successful training of concepts with only **5 GB of memory**, highlighting its efficiency and scalability across diverse model configurations.

## S3.5. Training time

We measured the total training time for each personalization method and observed a general inverse relationship between memory usage and training time in Tab. S10. The table shows a trade-off between time complexity and space complexity. However, Gradient-Free Textual Inversion (GF-TI) deviates from this trend, requiring significantly more training time despite its low memory usage. This behavior indicates that GF-TI struggles with learning effectively under small batch training conditions, which impacts its overall efficiency and practicality in resource-constrained devices.

| Base. | Method | Time (min) |
|-------|--------|------------|
| DB | DB | 2 |
| | QLoRA | 1.1 |
| | PEQA | 1.5 |
| | TuneQDM | 1.4 |
| TI | TI | 8.8 |
| | GF-TI | 293 |
| | ZOODiP ($n=2$) | 31 |

Table S10. Total training time with the configurations in Sec S1.5. DB-based methods consume more memory but train faster, while TI-based methods are more memory-efficient but require more iterations.
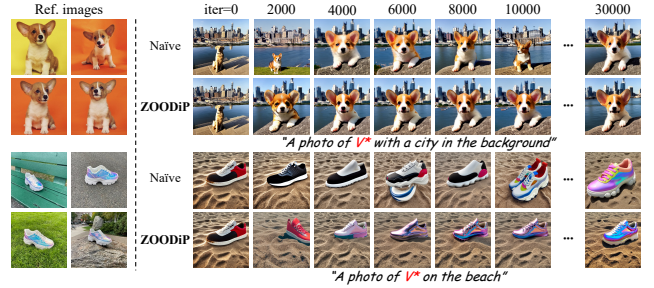


Figure S5. Qualitative comparisons on naïve ZO textual inversion without SG and PUTS to ZOODiP (Ours) over iterations. The naïve approach exhibits slower training and tends to produce images that are less aligned with the reference image. In contrast, ZOODiP achieves faster training and generates images that are closely aligned with the reference subject.

## S4. Additional Results

## S4.1. Qualitative results of SG and PUTS

Subspace Gradient (SG) and Partial Uniform Timestep Sampling (PUTS) play crucial roles in enabling efficient personalization by addressing two key challenges: SG removes noisy gradient components to focus on the most informative subspace, while PUTS suppresses sampling from ineffective timesteps to optimize training efficiency. To evaluate the impact of these components during the learning process, we generated images using the same seed with token embeddings learned at every $1,000$ iterations. As illustrated in Fig. S5, using both SG and PUTS significantly accelerates the training process compared to naïve zeroth-order optimization that does not incorporate these enhancements. The results show that models using SG and PUTS converge to high-quality personalization much faster than the baseline, underscoring the effectiveness of these components in improving the efficiency and speed of training.

## S4.2. Additional qualitative results

In this section, we present additional qualitative comparisons of personalization results that were not included in the main paper. These results are illustrated in Fig. S6, Fig. S7, Fig. S8, Fig. S9, and Fig. S10. The comparisons demonstrate that ZOODiP effectively captures the characteristics of the reference image and text prompts, achieving a level of fidelity comparable to other gradient-based personalization methods.

## S5. Limitation

In this section, we delve into the limitations of ZOODiP. While ZOODiP incorporates innovative techniques such as Subspace Gradient and Partial Uniform Timestep Sampling to mitigate the slow learning speed—a well-known challenge of zeroth-order optimization—it still requires a considerably larger number of iterations compared to gradient-based approaches. This increased iteration count stems from the fundamental nature of zeroth-order optimization, which relies on function evaluations rather than gradient backpropagation, inherently making it less sample-efficient. Although ZOODiP compensates with a faster training speed per iteration, the substantial number of iterations can introduce a significant time overhead, especially when proper hardware acceleration, such as high-performance Neural Processing Units (NPUs), is unavailable.

Furthermore, as ZOODiP is built upon the Textual Inversion framework, its performance is influenced by the strengths and weaknesses of Textual Inversion. This dependency implies that ZOODiP may face challenges in cases where Textual Inversion struggles, such as subjects with complex or highly variable visual characteristics, or when adapting to certain models that inherently perform poorly in personalization tasks. For example, if the base model lacks sufficient representational capacity or if the dataset used for training does not adequately capture the nuances of the subject, the effectiveness of ZOODiP can diminish.
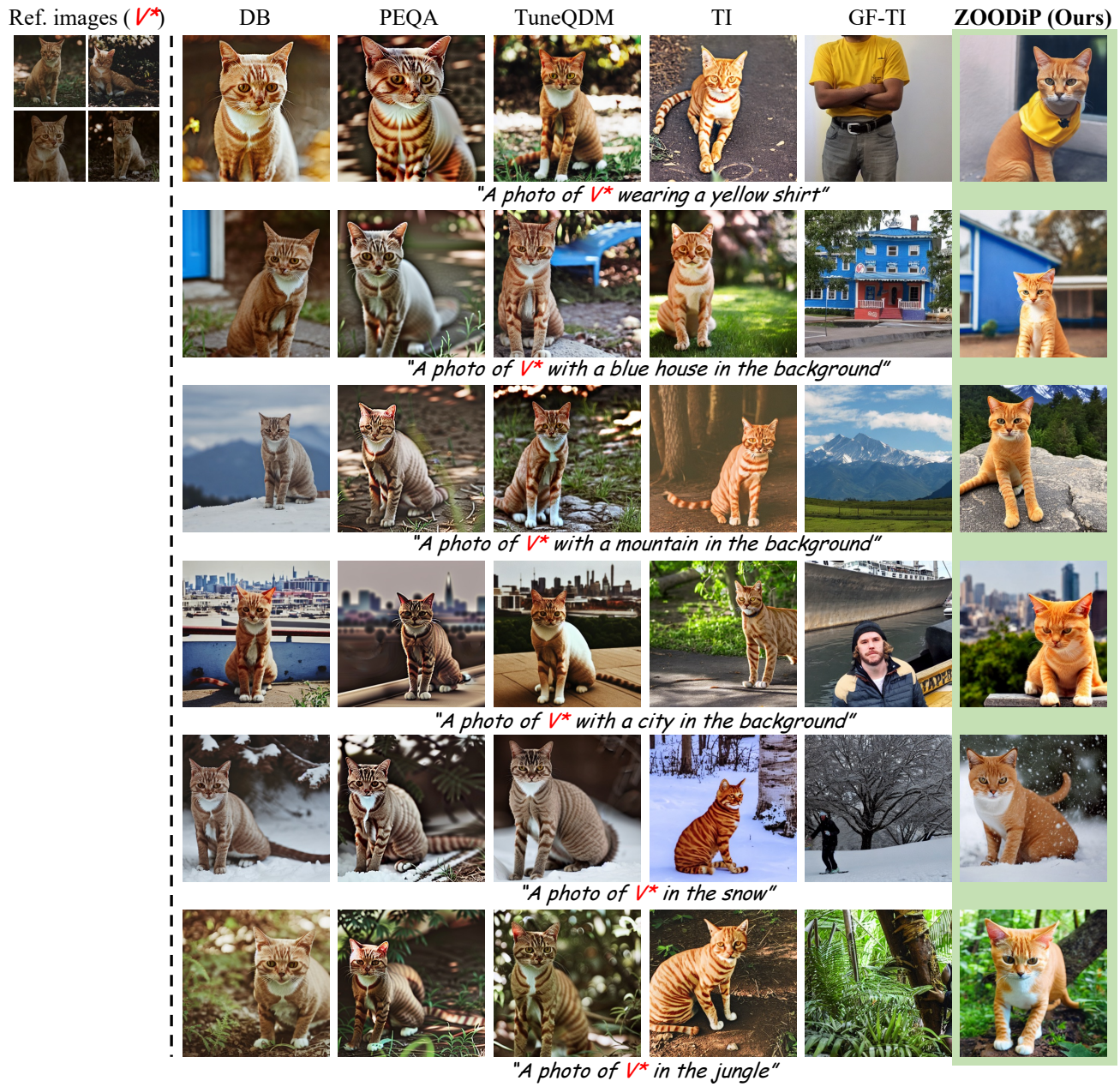
Figure S6. Qualitative comparison of image and text alignment on the `<cat>` subset of DB dataset.
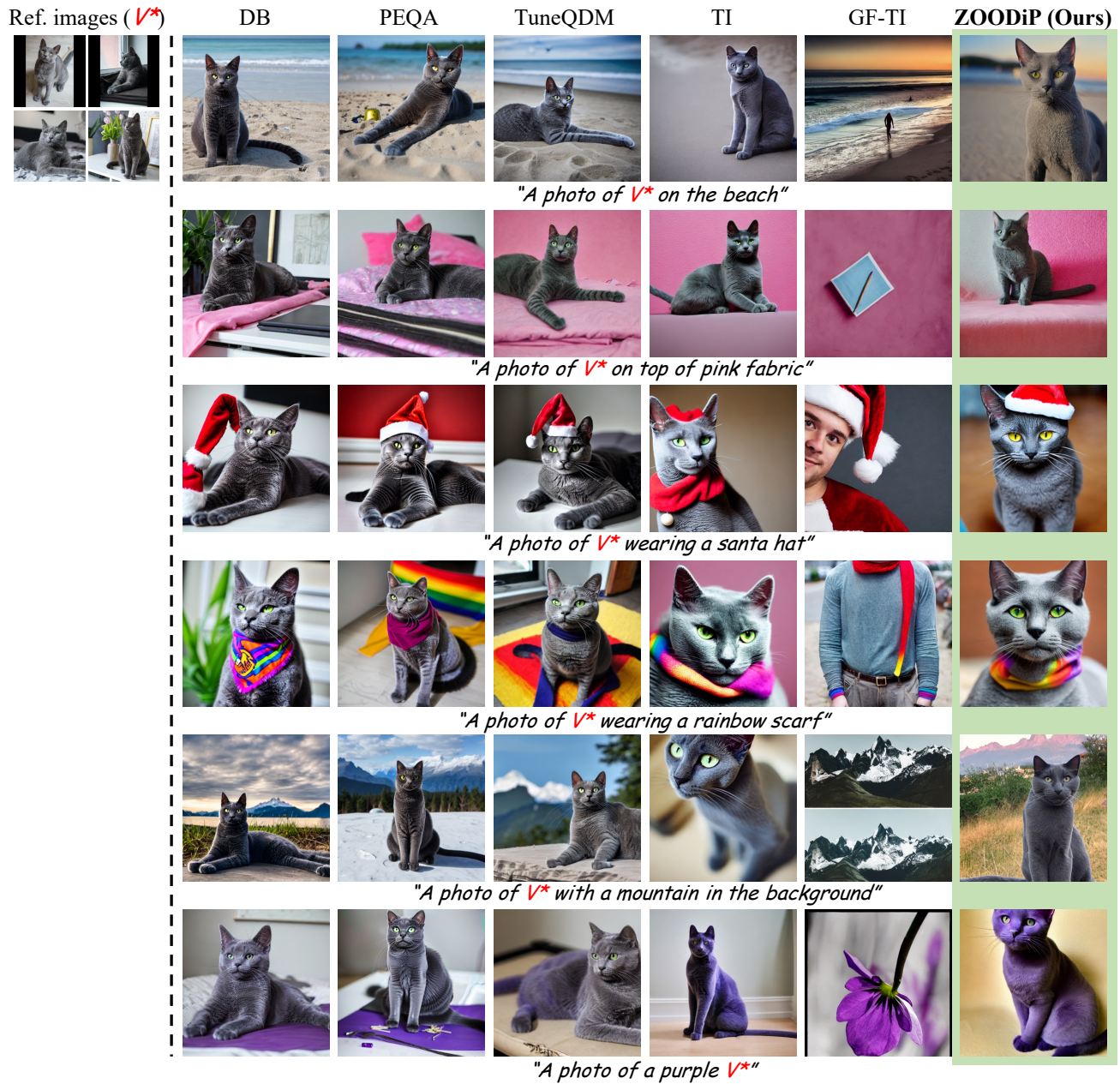
Figure S7. Qualitative comparison of image and text alignment on the `<cat2>` subset of DB dataset.

Figure S8. Qualitative comparison of image and text alignment on the `<dog6>` subset of DB dataset.

| Ref. images ( $V^*$ ) | DB | PEQA | TuneQDM | TI | GF-TI | **ZOODiP (Ours)** |
|---|---|---|---|---|---|---|

*"A photo of $V^*$ in the snow"*

*"A photo of $V^*$ on top of pink fabric"*

*"A photo of $V^*$ floating on top of water"*

*"A photo of $V^*$ on top of a dirt road"*

*"A photo of a purple $V^*$"*

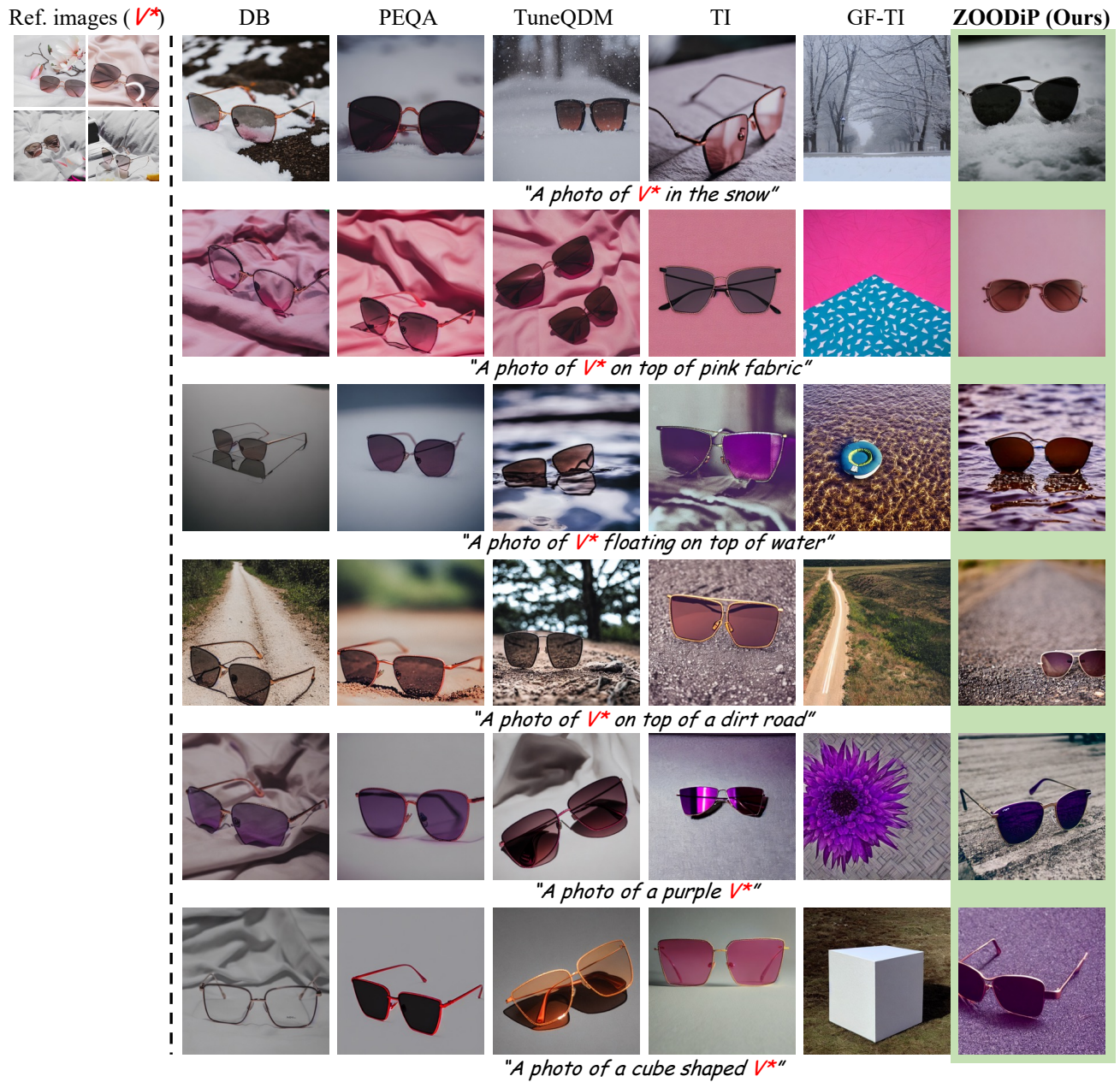*"A photo of a cube shaped $V^*$"*

Figure S9. Qualitative comparison of image and text alignment on the `<pink_sunglasses>` subset of DB dataset.

Figure S10. Qualitative comparison of image and text alignment on the `<shiny_sneaker>` subset of DB dataset.

# References

[1] Yuval Alaluf, Elad Richardson, Gal Metzer, and Daniel Cohen-Or. A neural space-time representation for text-to-image personalization. *ACM TOG*, 42(6):1–10, 2023. 1

[2] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Qinsheng Zhang, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, et al. ediff-i: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv:2211.01324*, 2022. 2, 5

[3] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv:1308.3432*, 2013. 2, 4

[4] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv:2311.15127*, 2023. 1

[5] Raghu Bollapragada, Richard Byrd, and Jorge Nocedal. Adaptive sampling strategies for stochastic optimization. *SIAM Journal on Optimization*, 28(4):3312–3343, 2018. 3

[6] HanQin Cai, Daniel McKenzie, Wotao Yin, and Zhenliang Zhang. Zeroth-order regularized optimization (zoro): Approximately sparse gradients and adaptive sampling. *SIAM Journal on Optimization*, 32(2):687–714, 2022. 3, 4

[7] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *ACM TOG*, 42(4):1–10, 2023. 2, 5

[8] Aochuan Chen, Yimeng Zhang, Jinghan Jia, James Diffenderfer, Jiancheng Liu, Konstantinos Parasyris, Yihua Zhang, Zheng Zhang, Bhavya Kailkhura, and Sijia Liu. Deepzero: Scaling up zeroth-order optimization for deep model training. *arXiv:2310.02025*, 2023. 3, 4

[9] Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart-\sigma: Weak-to-strong training of diffusion transformer for 4k text-to-image generation. *arXiv:2403.04692*, 2024. 1

[10] Jooyoung Choi, Jungbeom Lee, Chaehun Shin, Sungwon Kim, Hyunwoo Kim, and Sungroh Yoon. Perception prioritized training of diffusion models. In *CVPR*, pages 11472–11481, 2022. 2, 5

[11] CrossLabs. Diffusion with offset noise. Technical report, CrossLabs, 2023. 14

[12] Giannis Daras and Alexandros G Dimakis. Multiresolution textual inversion. *NeurIPSW*, 2022. 2, 5

[13] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *NeurIPS*, 36, 2024. 3, 5, 6, 9

[14] John C Duchi, Michael I Jordan, Martin J Wainwright, and Andre Wibisono. Optimal rates for zero-order convex optimization: The power of two function evaluations. *IEEE Transactions on Information Theory*, 61(5):2788–2806, 2015. 3

[15] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *ICML*, 2024. 1

[16] Zhengcong Fei, Mingyuan Fan, and Junshi Huang. Gradient-free textual inversion. In *ACM MM*, pages 1364–1373, 2023. 1, 3, 5, 6, 9, 11

[17] 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:2208.01618*, 2022. 1, 2, 3, 4, 5, 6, 9

[18] Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM journal on optimization*, 23(4):2341–2368, 2013. 3

[19] Hyojun Go, Yunsung Lee, Seunghyun Lee, Shinhyeok Oh, Hyeongdon Moon, and Seungtaek Choi. Addressing negative transfer in diffusion models. *NeurIPS*, 36, 2024. 2, 5

[20] Yuchao Gu, Xintao Wang, Jay Zhangjie Wu, Yujun Shi, Yunpeng Chen, Zihan Fan, Wuyou Xiao, Rui Zhao, Shuning Chang, Weijia Wu, et al. Mix-of-show: Decentralized low-rank adaptation for multi-concept customization of diffusion models. *NeurIPS*, 36, 2024. 1

[21] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 33:6840–6851, 2020. 1, 4

[22] 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:2210.02303*, 2022. 1

[23] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *NeurIPS*, 35:8633–8646, 2022. 1

[24] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv:2106.09685*, 2021. 2

[25] Gwanghyun Kim, Alonso Martinez, Yu-Chuan Su, Brendan Jou, José Lezama, Agrim Gupta, Lijun Yu, Lu Jiang, Aren Jansen, Jacob Walker, et al. A versatile diffusion transformer with mixture of noise levels for audiovisual generation. *arXiv:2405.13762*, 2024. 1

[26] Jeonghoon Kim, Jung Hyun Lee, Sungdong Kim, Joonsuk Park, Kang Min Yoo, Se Jung Kwon, and Dongsoo Lee. Memory-efficient fine-tuning of compressed large language models via sub-4-bit integer quantization. *NeurIPS*, 36, 2024. 1, 3, 5, 6, 9

[27] Raghuraman Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A whitepaper. *arXiv:1806.08342*, 2018. 4

[28] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. In *CVPR*, pages 1931–1941, 2023. 1, 2

[29] Kyungmin Lee, Sangkyung Kwak, Kihyuk Sohn, and Jinwoo Shin. Direct consistency optimization for compositional text-

to-image personalization. *arXiv preprint arXiv:2402.12004*, 2024. 9

[30] Yunsung Lee, JinYoung Kim, Hyojun Go, Myeongho Jeong, Shinhyeok Oh, and Seungtaek Choi. Multi-architecture multi-expert diffusion models. In *AAAI*, pages 13427–13436, 2024. 2, 5

[31] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *CVPR*, pages 300–309, 2023. 1

[32] Shanchuan Lin, Bingchen Liu, Jiashi Li, and Xiao Yang. Common diffusion noise schedules and sample steps are flawed. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 5404–5411, 2024. 14

[33] Zhenqing Ling, Daoyuan Chen, Liuyi Yao, Yaliang Li, and Ying Shen. On the convergence of zeroth-order federated tuning for large language models. In *ACM SIGKDD*, pages 1827–1838, 2024. 8

[34] Enshu Liu, Xuefei Ning, Zinan Lin, Huazhong Yang, and Yu Wang. Oms-dpm: Optimizing the model schedule for diffusion probabilistic models. In *ICML*, pages 21915–21936. PMLR, 2023. 2, 5

[35] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *ICCV*, pages 9298–9309, 2023. 1

[36] Sijia Liu, Pin-Yu Chen, Bhavya Kailkhura, Gaoyuan Zhang, Alfred O Hero III, and Pramod K Varshney. A primer on zeroth-order optimization in signal processing and machine learning: Principals, recent advances, and applications. *IEEE Signal Processing Magazine*, 37(5):43–54, 2020. 3

[37] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv:2209.03003*, 2022. 1

[38] Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan, Yue Huang, Hanchi Sun, Jianfeng Gao, et al. Sora: A review on background, technology, limitations, and opportunities of large vision models. *arXiv:2402.17177*, 2024. 1

[39] Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev Arora. Finetuning language models with just forward passes. *NeurIPS*, 36:53038–53075, 2023. 1, 2, 3, 5, 8, 12

[40] Yurii Nesterov and Vladimir Spokoiny. Random gradientfree minimization of convex functions. *Foundations of Computational Mathematics*, 17(2):527–566, 2017. 3, 4

[41] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv:2304.07193*, 2023. 5, 9

[42] NaHyeon Park, Kunhee Kim, and Hyunjung Shim. Textboost: Towards one-shot personalization of text-to-image models via fine-tuning text encoder. *arXiv:2409.08248*, 2024. 1, 2, 5, 8

[43] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *NeurIPS*, 32, 2019. 9

[44] Or Patashnik, Daniel Garibi, Idan Azuri, Hadar AverbuchElor, and Daniel Cohen-Or. Localizing object-level shape variations with text-to-image diffusion models. In *CVPR*, pages 23051–23061, 2023. 2, 5

[45] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, pages 4195–4205, 2023. 1

[46] 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:2307.01952*, 2023. 1, 14

[47] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv:2209.14988*, 2022. 1

[48] Qualcomm. Unlocking on-device generative ai with an npu and heterogeneous computing. Technical report, Qualcomm, 2024. 1

[49] 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 *ICML*, pages 8748–8763. PMLR, 2021. 5, 9

[50] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv:2204.06125*, 1(2):3, 2022. 1

[51] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022. 1, 3, 5, 9

[52] 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 *CVPR*, pages 22500–22510, 2023. 1, 2, 5, 6, 9, 10

[53] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Wei Wei, Tingbo Hou, Yael Pritch, Neal Wadhwa, Michael Rubinstein, and Kfir Aberman. Hyperdreambooth: Hypernetworks for fast personalization of text-to-image models. In *CVPR*, pages 6527–6536, 2024. 1

[54] Hyogon Ryu, Seohyun Lim, and Hyunjung Shim. Memory-efficient fine-tuning for quantized diffusion model. *arXiv:2401.04339*, 2024. 1, 3, 5, 6, 9

[55] Hoigi Seo, Hayeon Kim, Gwanghyun Kim, and Se Young Chun. Ditto-nerf: Diffusion-based iterative text to omnidirectional 3d model. *arXiv:2304.02827*, 2023. 1

[56] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, pages 2256–2265. PMLR, 2015. 1

[57] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv:2010.02502*, 2020.

[58] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *NeurIPS*, 32, 2019. 1

[59] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv:2011.13456*, 2020. 1

[60] James C Spall. Multivariate stochastic approximation using a simultaneous perturbation gradient approximation. *IEEE transactions on automatic control*, 37(3):332–341, 1992. 3, 4

[61] Yoad Tewel, Rinon Gal, Gal Chechik, and Yuval Atzmon. Key-locked rank one editing for text-to-image personalization. In *ACM SIGGRAPH*, pages 1–11, 2023. 1

[62] Andrey Voynov, Qinghao Chu, Daniel Cohen-Or, and Kfir Aberman. p+: Extended textual conditioning in text-to-image generation. *arXiv:2303.09522*, 2023. 1, 2

[63] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *CVPR*, pages 12619–12629, 2023. 1

[64] Daniel Watson, William Chan, Jonathan Ho, and Mohammad Norouzi. Learning fast samplers for diffusion models by differentiating through sample quality. In *ICLR*, 2022. 1

[65] Chendong Xiang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. A closer look at parameter-efficient tuning in diffusion models. *arXiv:2303.18181*, 2023. 1

[66] Penghang Yin, Jiancheng Lyu, Shuai Zhang, Stanley Osher, Yingyong Qi, and Jack Xin. Understanding straight-through estimator in training activation quantized neural nets. *arXiv:1903.05662*, 2019. 2

[67] Pengyun Yue, Long Yang, Cong Fang, and Zhouchen Lin. Zeroth-order optimization with weak dimension dependency. In *COLT*, pages 4429–4472. PMLR, 2023. 8

[68] Yuxin Zhang, Weiming Dong, Fan Tang, Nisha Huang, Haibin Huang, Chongyang Ma, Tong-Yee Lee, Oliver Deussen, and Changsheng Xu. Prospect: Prompt spectrum for attribute-aware personalization of diffusion models. *ACM TOG*, 42(6):1–14, 2023. 1, 2, 5