# Decouple-Then-Merge: Finetune Diffusion Models as Multi-Task Learning

Qianli Ma<sup>1</sup> Xuefei Ning<sup>2</sup> Dongrui Liu<sup>3</sup> Li Niu<sup>1,4†</sup> Linfeng Zhang<sup>1†</sup>

<sup>1</sup>Shanghai Jiao Tong University <sup>2</sup>Tsinghua University <sup>3</sup>Shanghai AI Laboratory <sup>4</sup>miguo.ai

{mqlqianli,ustcnewly,zhanglinfeng}@sjtu.edu.cn

## Abstract

Diffusion models are trained by learning a sequence of models that reverse each step of noise corruption. Typically, the model parameters are fully shared across multiple timesteps to enhance training efficiency. However, since the denoising tasks differ at each timestep, the gradients computed at different timesteps may conflict, potentially degrading the overall performance of image generation. To solve this issue, this work proposes a **De**couple-then-**Me**rge (DeMe) framework, which begins with a pretrained model and finetunes separate models tailored to specific timesteps. We introduce several improved techniques during the finetuning stage to promote effective knowledge sharing while minimizing training interference across timesteps. Finally, after finetuning, these separate models can be merged into a single model in the parameter space, ensuring efficient and practical inference. Experimental results show significant generation quality improvements upon 6 benchmarks including Stable Diffusion on COCO30K, ImageNet1K, PartiPrompts, and DDPM on LSUN Church, LSUN Bedroom, and CIFAR10. Code is available at GitHub.

## 1. Introduction

Generative modeling has seen significant progress in recent years, primarily driven by the development of Diffusion Probabilistic Models (DPMs) [15, 33, 41]. These models have been applied to various tasks such as text-to-image generation [40], image-to-image translation [43], and video generation [2, 16], yielding excellent performance. Compared with other generative models such as variational autoencoders (VAEs) [21], and generative adversarial networks (GANs) [11], the most distinct characteristic of DPMs is that DPMs need to learn a *sequence* of models for denoising at multiple timesteps. Training the neural network to fit this step-wise denoising conditional distribution facilitates tractable, stable training and high-fidelity generation.

The denoising tasks at different timesteps are similar yet



Figure 1. (a) Cosine similarity between gradients at different timesteps on CIFAR10 & distribution of gradients similarity in  $t \in [0, 1000]$  and  $t \in [0, 250]$ . Non-adjacent timesteps have low similarity, indicating conflicts during their training. In contrast, adjacent timesteps have similar gradients. (b) & (c): Comparison between the traditional and our training paradigm: The previous paradigm trains one diffusion model on all timesteps, leading to conflicts in different timesteps. Our method addresses this problem by decoupling the training of diffusion models in N different timestep ranges.

different. On the one hand, the denoising tasks at different timesteps are similar in the sense that the model takes a noisy image from the same space as input and performs a denoising task. Intuitively, sharing knowledge between these tasks might facilitate more efficient training. Therefore, typical methods let the model take both the noisy image  $x_t$  and the corresponding timestep t as input, and share the model parameter across all timesteps. On the other hand, the denoising tasks at different timesteps have clear differences as the input noisy images are from different distributions, and the concrete "denoising" effect is also different. Li et al. [25] demonstrate that there is a substantial difference between the feature distributions in different timesteps. Fang et al. [8] show that the larger (noisy) timesteps tend to generate the low-frequency and the basic image content, while the smaller timesteps tend to generate the high-frequency and the image details.

We further study the conflicts of different timesteps during the training of the diffusion model. Fig. 1(a) shows the gradient similarity of different timesteps. We can observe that the diffusion models have *dissimilar gradients at dif*-

<sup>&</sup>lt;sup>†</sup>Corresponding author

*ferent timesteps, especially the non-adjacent timesteps*, indicating a conflict between the optimization direction from different timesteps, as shown in Fig. 1(b). In one word, this gradient conflict indicates that different denoising tasks might have a negative interference with each other during training, which may harm the overall performance.

Considering the similarity as well as difference of these denoising tasks, the next natural and crucial question is "how can we promote effective knowledge sharing as well as avoid negative interference between multiple denoising *tasks?*". Timestep-wise model ensemble [1, 28] solves this problem by training and inferring multiple different diffusion models at various timesteps to avoid negative interference, though introducing huge additional storage and memory overhead. For instance, Liu et al. [28] employs 6 diffusion models during inference, leading to around  $6 \times$  increase in storage and memory requirements, which renders the method impractical in application. Additionally, various loss reweighting strategies [13, 46] solve this problem by balancing different denoising tasks and mitigating negative interference. However, it may alleviate but can not truly solve the gradient conflicts in different timesteps.

In this work, considering the challenges faced by timestep-wise model ensemble and loss reweighting, we propose **De**couple-then-**Me**rge (**DeMe**), a novel finetuning framework for diffusion models that achieves the best side of both worlds: mitigated training interference across different denoising tasks and inference without extra overhead. DeMe begins with a pretrained diffusion model and then finetunes its separate versions tailored to nooverlapped timestep ranges to avoid the negative interference of gradient conflicts. Several training techniques are introduced during this stage to preserve the benefits of knowledge sharing in different timesteps. Then, the post-finetuned diffusion models are merged into a single model in their parameter space, enabling effective knowledge sharing across multiple denoising tasks. Specifically, as shown in Fig. 1(c), we divide the overall timestep range [0,T) into multiple adjacent timestep ranges with no overlap as  $\{[(i-1)T/N, iT/N)\}_{i=1}^N$ , where T denotes the maximal timestep and N denotes the number of timestep ranges. Then, we finetune a pretrained diffusion model for each timestep range by only training it with the timesteps inside this range. As a result, we decouple the training of diffusion models at different timesteps. The gradients of different timesteps will not be accumulated together and their conflicts are naturally avoided. Besides, as shown in Fig. 3, we further introduce three simple but effective techniques during the finetuning stage, including Consistency Loss and Probabilistic Sampling to preserve the benefits from knowledge sharing across different timesteps, and *Channel-wise Projection* that directly enables the model to learn the channel-wise difference in different timesteps.



Figure 2. Visualization of the difference between the pre-finetuned and the post-finetuned diffusion model on the channel and spatial dimensions. We computed the difference in activation values before/after finetune along the channel and spatial dimensions of the image. (a) Visualization of channel activation, spatial activation, and their difference between the pre-finetuned and the postfinetuned model. (b) Distribution of difference for channel activation and spatial activation values. It can be observed that activation values **vary mostly in channel dimensions** during finetuning on a subset of timesteps.

After the finetuning stage, we obtain N diffusion models learned the knowledge in N different timesteps ranges, which also lead to N times costs in storage and memory. Then, we eliminate the additional costs by merging all these N models into a single model in their parameter space with the model merging technique [17]. In this way, the obtained merged model has the same computation and parameter costs as the original diffusion model while maintaining the knowledge from the N finetuned model, which indicates a notable improvement in generation quality.

Extensive experiments on 6 datasets have verified the effectiveness of DeMe for both unconditional and text-toimage generation. In summary, our contributions can be summarized as follows.

- We propose to decouple the training of diffusion models by finetuning multiple diffusion models in different timestep ranges. Three simple but effective training techniques are introduced to promote knowledge sharing between multiple denoising tasks in this stage.
- We propose to merge multiple finetuned diffusion models, each specialized for different timestep ranges, into a single diffusion model, which significantly enhances generation quality without any additional costs in computation, storage, and memory access. To the best of our knowledge, we are the first to merge diffusion models across different timesteps.
- Abundant experiments have been conducted on six datasets for both unconditional and text-to-image generation, demonstrating significant improvements in generation quality.

We note that our framework of combining task-specific training with parameter-space merging offers a novel method for multi-task learning, distinct from existing loss-balancing techniques [18, 48], and can be potentially extended to general multi-task scenarios.

## 2. Related Work

Diffusion Models. Diffusion models [6, 15, 33, 50, 52] represent a family of generative models that generate samples via a progressive denoising mechanism, starting from a random Gaussian distribution. Given that diffusion models suffer from slow generation and heavy computational costs, previous works have focused on improving diffusion models in various aspects, including model architectures [36, 41], faster sampler [29, 30, 51], prediction type and loss weighting [3, 10, 13, 46]. Besides, a few works have attempted to accelerate diffusion models generation through pruning [7], quantization [26, 49] and knowledge distillation [19, 31, 32, 46], which have achieved significant improvement on the efficiency. Motivated by the excellent generative capacity of diffusion models, diffusion models have been developed in several applications, including textto-image generation [39, 41, 44], video generation [2, 16], image restoration [45], natural language generation [24], audio synthesis [22], 3D content generation [37], ai4science such as protein structure generation [55], among others.

Training of Diffusion Models & Multi-task Learning. Multi-task Learning (MTL) is aimed at improving generalization performance by leveraging shared information across related tasks. The objective of MTL is to learn multiple related tasks jointly, allowing models to generalize better by learning representations that are useful for numerous tasks [4]. Despite its success in various applications, MTL faces significant challenges, particularly negative transfer [4, 53], which can degrade the performance of individual tasks when jointly trained. The training paradigm of diffusion models could be viewed as a multi-task learning problem: diffusion models are trained by learning a sequence of models that reverse each step of noise corruption across different noise levels. A parameter-shared denoiser is trained on different noise levels concurrently, which may cause performance degradation due to negative transfer-a phenomenon where learning multiple denoising tasks jointly hinders performance due to conflicts in timestep-specific denoising information. Previous works reweight training loss on different timesteps, improving diffusion model performance [3, 10, 15, 46] or accelerating training convergence [13]. Go et al. analyze and improve the diffusion model by exploring task clustering and applying various MTL methods to diffusion model training. Kim et al. analyze the difficulty of denoising tasks and propose a novel easy-to-hard learning scheme for progressively training diffusion models. DMP [12] integrates timestep specific learnable prompts into pretrained diffusion models, thereby enhancing their performance and enabling more effective optimization across different stages of training. Some works also reinterpret diffusion models using MTL and propose architectural improvements, such as DTR [34] and Switch-DiT [35]. Different from [10], we propose

to decouple the training of diffusion models by finetuning multiple diffusion models in different timestep ranges, and merge these models in the parameter space to mitigate gradient conflicts between timesteps.

## 3. Methodology

## 3.1. Preliminary

The fundamental concept of diffusion models is to generate images by progressively applying denoising steps, starting from random Gaussian noise  $x_T$ , and gradually transforming it into a structured image  $x_0$ . Diffusion models consist of two phases: the forward process and the reverse process. In the forward process, a data point  $\mathbf{x}_0 \sim q(\mathbf{x})$  is randomly sampled from the real data distribution, then gradually corrupted by adding noise step-by-step  $q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$ , where t is the current timestep and  $\beta_t$  is a pre-defined variance schedule that schedules the noise. In the reverse process, diffusion models transform a random Gaussian noise  $x_T \sim \mathcal{N}(0, \mathbf{I})$  into the target distribution by modeling conditional probability  $q(x_{t-1} | x_t)$ , which denoises the latent  $x_t$  to get  $x_{t-1}$ . Formally, the conditional probability  $p_{\theta}(x_{t-1} | x_t)$  can be modeled as:

$$\mathcal{N}\left(x_{t-1}; \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(x_t, t)\right), \beta_t \mathbf{I}\right), \quad (1)$$

where  $\alpha_t = 1 - \beta_t$ ,  $\bar{\alpha}_t = \prod_{i=1}^T \alpha_i$ .  $\epsilon_{\theta}$  denotes a noise predictor, which is usually an U-Net [42] autoencoder in diffusion models, with current timestep t and previous latent  $x_t$  as input. It is usually trained with the objective function:

$$\mathcal{L}_{\theta} = \mathbb{E}_{t \sim U[0,T], x_0 \sim q(x), \epsilon \sim \mathcal{N}(0,1)} \left[ \left\| \epsilon - \epsilon_{\theta} \left( x_t, t \right) \right\|^2 \right], \quad (2)$$

where T denotes the number of timesteps and U denotes a uniform distribution. After training, a clean image  $x_0$  can be obtained via an iterative denoising process from the random Gaussian noise  $x_T \sim \mathcal{N}(0, \mathbf{I})$  with the modeled distribution  $x_{t-1} \sim p_\theta(x_{t-1} \mid x_t)$  in Equation 1.

## 3.2. Decouple the Training of Diffusion Model

In this section, we demonstrate how to decouple the training of diffusion model. As illustrated in Fig. 1(c), we first divide the timesteps of [0, T) into N multiple continuous and non-overlapped timesteps ranges, which can be formulated as  $\{[(i-1)T/N, iT/N)\}_{i=1}^N$ . Subsequently, based on a diffusion model pretrained by Equation 1, we finetune a group of N diffusion models  $\{\epsilon_{\theta_i}\}_{i=1}^N$  on each of the N timestep ranges. The training objective of  $\epsilon_{\theta_i}$  which can be formulated as

$$\mathbb{E}_{t \sim U\left[ ^{(i-1)T/N, iT/N} \right], x_0 \sim q(x), \epsilon \sim \mathcal{N}(0, 1)} \left[ \left\| \epsilon - \epsilon_{\theta_i} \left( x_t, t \right) \right\|^2 \right].$$
(3)



Figure 3. Pipeline of our framework. The following training techniques are incorporated into the finetuning process: **Consistency loss** preserves the original knowledge of diffusion models learned at all timesteps by minimizing the difference between pre-finetuned and post-finetuned diffusion models. **Probabilistic sampling** strategy samples from both the corresponding and other timesteps with different probabilities, helping the diffusion model overcome forgetting knowledge from other timesteps. **Channel-wise projection** enables the diffusion model to directly capture the feature difference in channel dimension. Model merging scheme merges the parameters of all the finetuned models into one unified model to promote the knowledge sharing across different timestep ranges.

However, although Equation 3 can decouple the training of the diffusion model in different timesteps and avoid the negative interference between multiple denoising tasks, it also eliminates the positive benefits of learning from different timesteps, which may make the finetuned diffusion model overfit a specific timestep range and lose its knowledge in the other timesteps. Besides, it is also challenging for the diffusion model to capture the difference in different timesteps during finetuning. To address these problems, we further introduce the following techniques shown in Fig. 3. Channel-wise Projection. Fig. 2 shows the difference between the pre-finetuned and the post-finetuned diffusion models, demonstrating that there is a significant difference in the channel dimension instead of the spatial dimension, which further implies that the knowledge learned during finetuning in a timestep range is primarily captured by channel-wise mapping instead of spatial mapping. Based on this observation, we further apply a channel-wise projection layer to facilitate the training process by directly formulating the channel-wise mapping. Let  $\mathbf{F}_t \in \mathbb{R}^{C \times H \times W}$ denote the intermediate feature map of the noise predictor  $\epsilon_{\theta}(x_t, t)$  at the timestep t, where C, H, W denote the number of channels, height, and width of the feature map  $\mathbf{F}_t$ , respectively. The channel-wise projection is designed as  $\mathbb{P}(\mathbf{F}_t) = \mathbf{W} \cdot \mathbf{F}_t$ , where  $\mathbf{W} \in \mathbb{R}^{C \times C}$  is a learnable projection matrix that enables the diffusion model to directly capture the feature difference in the channel dimension. Please note that we initialize W as an identity matrix to stabilize the training process. It is worth noting that the parameter of channel-wise projection layer is small, accounting for 1.06% of the diffusion model.

**Consistency Loss.** A consistency loss is introduced into the training process to minimize the difference between

the pre-finetuned and post-finetuned diffusion model, which can be formulated as

$$\mathbb{E}_{t \sim U\left[\left(i-1\right)T/N, iT/N\right]} \left| \left\| \epsilon_{\theta}\left(x_{t}, t\right) - \epsilon_{\theta_{i}}\left(x_{t}, t\right) \right\|^{2} \right|, \quad (4)$$

where  $\epsilon_{\theta}(x_t, t)$  denotes the output of the original diffusion model.  $\epsilon_{\theta_i}(x_t, t)$  denotes the output of  $i_{th}$  post-finetuned diffusion model. Minimizing the consistency loss preserves the initial knowledge of the diffusion model, and ensures that the finetuned diffusion models do not differ significantly from the pre-finetuned diffusion model. Besides, the consistency loss also enhances the stability of the training process for finetuning diffusion models in the timestep range. Combining Equation 3 and Equation 4, we can derive the overall loss:

**Probabilistic Sampling.** To further preserve the initial knowledge learned at all the timesteps, we design a *Probabilistic Sampling* strategy which enables the finetuned model to mainly learn from its corresponding timestep range, but still possible to preserve the knowledge in the other timestep ranges. Concretely, during the finetuning of  $i_{th}$  diffusion model, we sample t from the timestep range [(i-1)T/N, iT/N) with a probability of 1-p, while sampling from the overall range [0, T) with a probability p. The overall sampling strategy can be expressed as follows:

$$t \sim \begin{cases} \left[ {^{(i-1)T}/N}, {^{iT}/N} \right], i \in [1, N] & \text{with probability } 1 - p, \\ \left[ 0, T \right] & \text{with probability } p. \end{cases}$$
(6)



Figure 4. Loss landscape of the pretrained diffusion model in different timestep ranges on CIFAR10. We use dimension reduction methods to handle high-dimensional neural networks. Contour line density reflects the frequency of loss variations (*i.e.*, gradients), with blue representing low loss and red representing high loss. The pretrained model resides at the critical point (with zero gradients) with sparse contour lines for the overall timesteps  $t \in [0, 1000)$ , but when the training process is decoupled, it tends to be located in regions with densely packed contour lines, suggesting that there still exists gradients that enable pretrained model to escape from the critical point.

### **3.3. Merging Models in Different Timestep Ranges**

After finetuning N diffusion models in their corresponding timesteps, it is a natural step to ensemble these finetuned diffusion models in inference stage. The sampling process under a timestep-wise model ensemble scheme is achieved by inferring each post-finetuned diffusion model in its corresponding timestep range, which can be formulated as

$$p_{\theta}(x_{t-1} \mid x_t) = \mathcal{N}\left(x_{t-1}; \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_{\theta_i}(x_t, t)\right), \beta_t \mathbf{I}\right),$$
(7)

where  $i = \lfloor t \times N/T \rfloor$ . For instance, the  $i_{th}$  finetuned diffusion model is only utilized in timestep  $t \in \lfloor (i-1)T/N, T/N \rfloor$ . This inference scheme does not introduce additional computation costs during the inference period but does incur additional storage costs. Since model merging methods [17, 54] can integrate diverse knowledge from each model, we propose to merge multiple finetuned diffusion models into *a single diffusion model*, which avoids additional computation or storage costs during inference while significantly improving generation quality.

**Model Merging Scheme.** Fig. 3 shows the overview of the model merge scheme. Inspired by model merging methods [17, 54] that aim to merge the parameters of models finetuned in different datasets and tasks, we propose to merge multiple post-finetuned diffusion models. Specifically, we first compute the *task vectors* of different post-finetuned diffusion models, which indicates the difference in their parameters compared with the pre-finetuned version. The task vector  $\tau_i$  of the  $i_{th}$  finetuned diffusion model can be denoted as  $\tau_i = \theta_i - \theta$ , where  $\theta$  and  $\theta_i$  denote the parameters of the pre-finetuned and the  $i_{th}$  post-finetuned diffusion model. Following previous work [17], the model merging can be achieved by adding all the task vectors to

the pre-finetuned model, which can be formulated as

$$\theta_{\text{merged}} = \theta + \sum_{i=1}^{N} w_i \tau_i, \quad \text{where} \quad \tau_i = \theta_i - \theta, \quad (8)$$

where  $w_i$  means merging weights of task vectors. We use grid search algorithm to obtain the optimal combination of  $w_i$ . In this scheme, we finally obtain  $\theta_{\text{merged}}$  which can be applied across all timesteps in [0, T), following the same inference process as in traditional diffusion models. As a result, the model merge scheme also leads to significant enhancement in generation quality without introducing any additional costs in computation or storage during inference.

### 4. Experiments

#### 4.1. Experiment Setting

Datasets and Metrics. For unconditional image generation datasets CIFAR10 [23], LSUN-Church, and LSUN-Bedroom [56], we generated 50K images for evaluation. For text-to-image generation, following the previous work [19], we finetune each model on a subset of LAION-Aesthetics V2 (L-Aes) 6.5+ [47] and test model's capacity of zero-shot text-to-image generation on MS-COCO validation set [27], ImageNet1K [5] and PartiPrompts [57]. Fréchet Inception Distance [14] is used to evaluate the quality of generated images. CLIP score computed by CLIP-ViT-g/14 [38] is used to evaluate the text-image alignment. **Baselines.** We choose some loss reweighting methods as baselines for comparison: SNR+1, truncated SNR [46], *Min-SNR-\gamma* [13], P2 weighting [3]. We also select ANT [10] to apply MTL methods to different timestep intervals for comparison. We prove that the aforementioned diffusion loss weights can be unified under the same prediction target with different weight forms(Proof in supplementary material). We also demonstrate that our decouple-thenmerge framework can be formally transformed into the loss reweighting framework (proof in supplementary material). To ensure a fair comparison, the baseline models are trained with an equal number of iterations with our training frame-

Table 1. Quantitative results (FID, lower is better) on CIFAR10, LSUN-Church, and LSUN-Bedroom with DDPM. Numbers in the brackets indicate the FID difference compared with DDPM.

Method	CIFAR10	LSUN-Church	LSUN-Bedroom	#Iterations
Before-finetuning [15]	4.42	10.69	6.46	-
SNR+1 [46]	5.41	10.80	6.41	80K
Trun-SNR [46]	4.49	10.81	6.42	80K
Min-SNR- $\gamma$ [13]	5.77	10.82	6.41	80K
P2 Weighting [3]	5.63	10.77	6.53	80K
ANT-NashMTL [10]	4.24	10.45	6.43	80K
ANT-UW [10]	4.21	10.43	6.48	80K
DeMe (Before Merge)	$3.79_{(-0.63)}$	9.57 (-1.12)	5.87 (-0.59)	20K×4
DeMe (After Merge)	$3.51_{(-0.91)}$	7.27 <sub>(-3.42)</sub>	$5.84_{(-0.62)}$	20K×4

work. Additionally, we also ensemble finetuned diffusion models and compare them with the merging scheme for a more detailed comparison. Please refer to the supplementary material for details on the implementation.

## 4.2. Quantitative Study

Results on Unconditional Generation. Table 1 presents quantitative results on unconditional generation, demonstrating great improvement in generation quality across various unconditional image generation benchmarks. The model merging scheme achieves performance comparable to, or even better than, the ensemble scheme with a unified diffusion model, highlighting the superiority of the merging approach. Concretely, 0.63, 1.12, and 0.59 FID reduction can be observed on CIFAR10, LSUN-Church, and LSUN-Bedroom with the model ensemble scheme, respectively. The model merging scheme leads to 0.91, 3.42, and 0.62 FID reductions on CIFAR10, LSUN-Church, and LSUN-Bedroom, respectively. In contrast, previous loss weighting methods obtain very few FID reductions under the same finetuning setting and even harm the generation quality during fine-tuning.

Results on Text-to-Image Generation. Table 2 shows that DeMe outperforms the baselines in both image quality and text-image alignment, as demonstrated by the experiment results on text-to-image generation benchmarks for Stable Diffusion [41]. Specifically, on MS COCO, our ensemble method achieves a 0.64 FID reduction and a 0.03 CLIP score reduction, while merging method yields a 0.36 FID reduction along with a 0.23 CLIP score increase. On ImageNet1k, ensemble method results in a 1.26 FID reduction and a 0.17 CLIP score reduction, whereas merging method produces a 0.39 FID reduction and a 0.02 CLIP score increase. Additionally, on PartiPrompts, both the ensemble and merging schemes show improvements in CLIP score, with increases of 0.24 and 0.20, respectively. These results validate the effectiveness of DeMe, showing significant improvements in both image quality and text-image alignment.

rumpt I: "A graceful white horse galloping through a dig of wildflowers, its mane flowing in the wind as the next behind it."
 Prompt II: "A trapical beach wild it dig of wildflowers, its mane flowing in the wind as the next behind it."
 Profession of test-image alignment: as the san sets behind it."
 Before Finetungs: used test-image alignment, lifelike horse end to be hold it.
 Before Finetungs: used test-image alignment is the san sets behind it.





reze under a clear blue sky."

Finetuning After Finetuning Be

Before Finetuning After Finetuning

Figure 5. Qualitative comparison between DeMe and the original Stable Diffusion on various prompts. More images based on various text prompts could be found in supplementary material.

## 4.3. Qualitative Study

Fig. 5 depicts some fancy generated images given detailed prompts, which illustrates that our method effectively generates images that align with the provided text descriptions, resulting in generated images that are both more detailed and photorealistic. Prompts highlighted in bold indicate where Stable Diffusion fails to align the image with the text, whereas our method generates images with better text-image alignment. For example, in the middle image pair of Fig. 5, Stable Diffusion fails to generate asmall wooden cabin in the image, while our method successfully captures the subject and preserves the detailed information described in the prompt. The finetuned Stable Diffusion model demonstrates an improved ability to generate visually coherent and contextually accurate images that closely match the nuances of the prompts, as highlighted in the comparison between before- and after-finetuning results, showcasing its enhanced capacity for text-to-image synthesis. More figures based on various text prompts could be found in supplementary material. Besides, we also provide images generated on LSUN in supplementary material.

### 4.4. Ablation Study

Our framework applies three training techniques to finetune diffusion model in different timesteps. As shown in Table 3, we conducted ablation studies on training techniques individually. All experiments are conducted on CIFAR10, with a 100-step DDIM sampler [51]. Several key observations can be made: (i) The traditional training paradigm results in the poorest performance. With N set to 1 and none of the specialized training techniques applied-following the traditional diffusion training paradigm—the model yields a poor results, with a FID of 4.40. Gradient conflicts lead to negative interference across different denoising tasks, adversely affecting overall training. (ii) Channel-wise pro-

	MS-COCO		ImageNet		PartiPrompts	
Method	FID↓	CLIP Score↑	FID↓	CLIP Score↑	CLIP Score↑	#Iterations
Before-finetuning [41]	13.42	29.88	27.62	27.07	29.78	-
SNR+1 [46]	13.92	29.96	27.56	27.03	29.86	80K
Trun-SNR [46]	13.93	29.95	27.60	27.05	29.85	80K
Min-SNR- $\gamma$ [13]	13.92	29.93	27.59	27.02	29.87	80K
P2 Weighting [3]	13.23	29.93	26.92	26.44	29.50	80K
ANT-NashMTL [10]	13.39	29.81	27.41	26.99	29.90	80K
ANT-UW [10]	13.17	29.94	26.91	26.78	29.98	80K
DeMe (Before Merge)	$12.78_{(-0.64)}$	29.85 (-0.03)	26.36 (-1.26)	26.90 (-0.17)	30.02 (+0.24)	20K×4
DeMe (After Merge)	$13.06_{(-0.36)}$	$30.11_{(+0.23)}$	$27.23_{(-0.39)}$	$27.09_{(+0.02)}$	<b>29.98</b> (+0.20)	20K×4

Table 2. Quantitative studies on MS COCO, PartiPrompts and ImageNet with Stable Diffusion. Numbers in the brackets indicate the FID or CLIP Score difference compared with Stable Diffusion.

Table 3. Ablation study on CIFAR10. N denotes the number of finetuned models.

Ν	Probabilistic Sampling	Consistency Loss	Channel-wise Projection	FID↓
1	× ×	× ×	× •	4.40 4.45
8	2 2 2	× •	× × •	4.32 4.27 3.87

jection struggles to capture feature differences in the channel dimension without alleviating gradient conflicts. With N set to 1 and Channel-wise projection applied, model yields a worse results, with a FID of 4.45. In contrast, with N set to 8 and Channel-wise projection applied additionally, model yields the best results, with a FID of 3.87. We posit that channel-wise projection struggles to capture feature changes due to the significant differences across the timesteps. (iii) Dividing overall timesteps into N nonoverlapping ranges effectively alleviates gradient conflicts, resulting in a significant reduction in FID. For instance, with N set to 8, introducing Probabilistic Sampling achieves a 0.08 FID reduction, while applying Consistency Loss yields an additional 0.05 FID reduction. When all techniques are applied during finetuning, a total FID reduction of 0.53 is achieved. Our experimental results demonstrates that dividing overall timestep into non-overlapping ranges serves as a necessary condition. Building on this foundation, our training techniques significantly improve model performance. Sensitive studies on influence of N and p have been conducted in supplementary material, demonstrating that our method is robust to variations in the choices of N and p.

## 5. Discussion

**DeMe Enables Pretrained Model Escaping from the Critical Point.** We explore how DeMe can guide pre-

trained models to escape from critical points, leading to further optimization. We refer to the approach of [9, 54], visualizing the relationship between model parameters and training loss by plotting the loss landscape. Fig. 4 presents some visualization results on the training loss landscape that support our claims. Two significant findings can be drawn from Fig. 4: (i) The pretrained diffusion model has converged when  $t \in [0, 1000)$ , residing at the critical point with sparse contour lines (i.e., no gradient). However, it is evident that the pretrained model is not at an optimal point, as there are nearby points with lower training loss, suggesting a potential direction for further optimization. (ii) For different timestep ranges, the pretrained model tends to be situated in regions with densely packed contour lines (i.e., larger gradient), suggesting that there exists an optimization direction. For instance, when  $t \in [0, 250)$ , the pretrained model stays at a point with frequent loss variations, indicating a potential direction for lower training loss. The decoupled training framework facilities the diffusion model to optimize more efficiently. Based on the above observation, DeMe decouples the training process, enabling the pretrained model to move away from the critical point, resulting in further improvement.

Loss Landscape Visualization for Task Vectors. To provide some intuitions, we visualize a two-dimensional training loss representation when applying two task vectors to merge finetuned models across various datasets, shown in Fig. 6(a). We utilize pretrained model  $\theta$ , two finetuned model  $\theta_i$  (i = 1, 2) to obtain two task vectors  $\tau_i$  (i = 1, 2), which span a plane in parameter space. We evaluate the diffusion training loss on this plane, and there are three key observations obtained from Fig. 6(a): (i) For both CI-FAR10 and LSUN-Church, the training loss contours are basin-shaped and none of the model parameters are optimal, which means there exists a direction towards a better model parameters. (ii) The weighted sum of task vectors  $\tau$ 



Figure 6. (a): Loss landscape for applying task vectors. The optimal model parameters are neither the pretrained one nor the finetuned one, but lie within the plane spanned by the task vectors computed in Sec. 3.3. We utilize the pretrained and two finetuned model parameters to obtain the two task vectors, respectively. Following [9, 54], we compute an orthonormal basis from the plane spanned by the task vectors. Axis denotes the movement direction in the parameter space. (b): Box plot of task vector distribution over different layers on LSUN-Church. Task vectors exhibit notable value in  $t \in [500, 1000)$  but only slight value in  $t \in [0, 500)$ .

(*i.e.*, the interpolation of finetuned model parameters  $\theta_i$ ) can yield parameters with a lower training loss. For instance, on CIFAR10, the weighted sum of the task vectors can produce optimal model parameters, outperforming the two individual finetuned model parameters. (iii) The loss variation is relatively smooth, opening up the possibility to employ advanced search methods, such as evolutionary search, which could serve as a potential avenue for further improvement. In brief summary, the above observations indicate that by applying a weighted sum to the task vectors, a more optimal set of model parameters can be achieved, leading to a lower training loss.

Task Vector Analysis. DeMe employs the model merge technique to merge the multiple finetuned diffusion models by calculating the *linear combination of task vectors* introduced in Equation 8. Here we visualize the task vectors in Fig. 6(b), which shows significant differences between the task vectors in different timestep ranges. Specifically, the magnitude of task vectors has a larger value for  $t \in [500, 1000)$  and a smaller value for  $t \in [0, 500)$ , indicating that there are more significant differences in parameters for diffusion models finetuned for  $t \in [500, 1000)$ . We suggest this because the original SNR loss term [15] has lower values in larger t. As a result, the original diffusion model bias to the gradients in smaller t when larger t and smaller thave conflicts in gradients, leading to poor optimization for larger t. In contrast, DeMe decouples the training of diffusion models across larger t and smaller t, allowing different timestep ranges to be optimized separately. Hence, the diffusion model finetuned on larger t exhibits a more significant difference compared with the original model, which leads to better generalization quality.

## 6. Conclusion

Motivated by the observation that different timesteps in the diffusion model training have low similarity in their gradients, this paper proposes DeMe, which decouples the training of diffusion models in different timesteps and merge the finetuned diffusion models in parameter-space, thereby mitigating the negative impacts of gradient conflicts. Besides, three simple but effective training techniques have been introduced to facilitate the finetuning process, which preserve the benefits of knowledge sharing in different timesteps. Our experimental results on six datasets with both unconditional and text-to-image generation demonstrate that our approach leads to substantial improvements in generation quality without incurring additional computation or storage costs during sampling. The effectiveness of DeMe may promote more research work on the optimization of diffusion models. Additionally, the feasibility of combining taskspecific training with parameter-space merging presented in this work may stimulate more research into diffusion model merging, and can be potentially extended to general multitask learning scenarios.

## Acknowledgement

The work is supported by the National Natural Science Foundation of China (Grant No.62471287)

## References

- [1] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Qinsheng Zhang, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, et al. ediff-i: Text-to-image diffusion models with an ensemble of expert denoisers. arXiv preprint arXiv:2211.01324, 2022. 2
- [2] 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 preprint arXiv:2311.15127, 2023. 1, 3
- [3] Jooyoung Choi, Jungbeom Lee, Chaehun Shin, Sungwon Kim, Hyunwoo Kim, and Sungroh Yoon. Perception prioritized training of diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11472–11481, 2022. 3, 5, 6, 7
- [4] Michael Crawshaw. Multi-task learning with deep neural networks: A survey. *arXiv preprint arXiv:2009.09796*, 2020.
   3
- [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pages 248–255. Ieee, 2009. 5
- [6] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. Advances in neural information processing systems, 34:8780–8794, 2021. 3
- [7] Gongfan Fang, Xinyin Ma, and Xinchao Wang. Structural pruning for diffusion models. In Advances in Neural Information Processing Systems, 2023. 3
- [8] Gongfan Fang, Xinyin Ma, and Xinchao Wang. Structural pruning for diffusion models. arXiv preprint arXiv:2305.10924, 2023. 1
- [9] Timur Garipov, Pavel Izmailov, Dmitrii Podoprikhin, Dmitry P Vetrov, and Andrew G Wilson. Loss surfaces, mode connectivity, and fast ensembling of dnns. *Advances in neural information processing systems*, 31, 2018. 7, 8
- [10] Hyojun Go, Yunsung Lee, Seunghyun Lee, Shinhyeok Oh, Hyeongdon Moon, and Seungtaek Choi. Addressing negative transfer in diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024. 3, 5, 6, 7
- [11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 1
- [12] Seokil Ham, Sangmin Woo, Jin-Young Kim, Hyojun Go, Byeongjun Park, and Changick Kim. Diffusion model patching via mixture-of-prompts. *arXiv preprint arXiv:2405.17825*, 2024. 3
- [13] Tiankai Hang, Shuyang Gu, Chen Li, Jianmin Bao, Dong Chen, Han Hu, Xin Geng, and Baining Guo. Efficient diffusion training via min-snr weighting strategy. In *Proceedings*

of the IEEE/CVF International Conference on Computer Vision, pages 7441–7451, 2023. 2, 3, 5, 6, 7

- [14] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems, 30, 2017. 5
- [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1, 3, 6, 8
- [16] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022. 1, 3
- [17] Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. arXiv preprint arXiv:2212.04089, 2022. 2, 5
- [18] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7482–7491, 2018. 2
- [19] Bo-Kyeong Kim, Hyoung-Kyu Song, Thibault Castells, and Shinkook Choi. Bk-sdm: A lightweight, fast, and cheap version of stable diffusion. *arXiv preprint arXiv:2305.15798*, 2023. 3, 5
- [20] Jin-Young Kim, Hyojun Go, Soonwoo Kwon, and Hyun-Gyoon Kim. Denoising task difficulty-based curriculum for training diffusion models. arXiv preprint arXiv:2403.10348, 2024. 3
- [21] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114, 2013. 1
- [22] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. arXiv preprint arXiv:2009.09761, 2020. 3
- [23] Alex Krizhevsky et al. Learning multiple layers of features from tiny images. 2009. 5
- [24] Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-Im improves controllable text generation. Advances in Neural Information Processing Systems, 35:4328–4343, 2022. 3
- [25] Xiuyu Li, Long Lian, Yijiang Liu, Huanrui Yang, Zhen Dong, Daniel Kang, Shanghang Zhang, and Kurt Keutzer. Q-diffusion: Quantizing diffusion models. arXiv preprint arXiv:2302.04304, 2023. 1
- [26] Xiuyu Li, Yijiang Liu, Long Lian, Huanrui Yang, Zhen Dong, Daniel Kang, Shanghang Zhang, and Kurt Keutzer. Q-diffusion: Quantizing diffusion models. In *Proceedings* of the IEEE/CVF International Conference on Computer Vision, pages 17535–17545, 2023. 3
- [27] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13, pages 740–755. Springer, 2014. 5

- [28] Enshu Liu, Xuefei Ning, Zinan Lin, Huazhong Yang, and Yu Wang. Oms-dpm: Optimizing the model schedule for diffusion probabilistic models. In *International Conference* on Machine Learning, pages 21915–21936. PMLR, 2023. 2
- [29] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. arXiv preprint arXiv:2202.09778, 2022. 3
- [30] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022. 3
- [31] Eric Luhman and Troy Luhman. Knowledge distillation in iterative generative models for improved sampling speed. *arXiv preprint arXiv:2101.02388*, 2021. 3
- [32] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 14297–14306, 2023. 3
- [33] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pages 8162–8171. PMLR, 2021. 1, 3
- [34] Byeongjun Park, Sangmin Woo, Hyojun Go, Jin-Young Kim, and Changick Kim. Denoising task routing for diffusion models. arXiv preprint arXiv:2310.07138, 2023. 3
- [35] Byeongjun Park, Hyojun Go, Jin-Young Kim, Sangmin Woo, Seokil Ham, and Changick Kim. Switch diffusion transformer: Synergizing denoising tasks with sparse mixture-ofexperts. In *European Conference on Computer Vision*, pages 461–477. Springer, 2024. 3
- [36] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023. 3
- [37] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. arXiv preprint arXiv:2209.14988, 2022. 3
- [38] 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*, pages 8748–8763. PMLR, 2021. 5
- [39] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1 (2):3, 2022. 3
- [40] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 10684–10695, 2022. 1
- [41] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of*

the IEEE/CVF conference on computer vision and pattern recognition, pages 10684–10695, 2022. 1, 3, 6, 7

- [42] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. Unet: Convolutional networks for biomedical image segmentation. In Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18, pages 234–241. Springer, 2015. 3
- [43] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In ACM SIGGRAPH 2022 conference proceedings, pages 1–10, 2022. 1
- [44] 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. Advances in neural information processing systems, 35:36479–36494, 2022. 3
- [45] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image superresolution via iterative refinement. *IEEE transactions on pattern analysis and machine intelligence*, 45(4):4713–4726, 2022. 3
- [46] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022. 2, 3, 5, 6, 7
- [47] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. Advances in Neural Information Processing Systems, 35:25278–25294, 2022. 5
- [48] Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. Advances in neural information processing systems, 31, 2018. 2
- [49] Yuzhang Shang, Zhihang Yuan, Bin Xie, Bingzhe Wu, and Yan Yan. Post-training quantization on diffusion models. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 1972–1981, 2023. 3
- [50] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015. 3
- [51] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. arXiv preprint arXiv:2010.02502, 2020. 3, 6
- [52] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. arXiv preprint arXiv:2011.13456, 2020. 3
- [53] Zirui Wang, Yulia Tsvetkov, Orhan Firat, and Yuan Cao. Gradient vaccine: Investigating and improving multi-task optimization in massively multilingual models. arXiv preprint arXiv:2010.05874, 2020. 3
- [54] Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos,

Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *International conference on machine learning*, pages 23965–23998. PMLR, 2022. 5, 7, 8

- [55] Kevin E Wu, Kevin K Yang, Rianne van den Berg, Sarah Alamdari, James Y Zou, Alex X Lu, and Ava P Amini. Protein structure generation via folding diffusion. *Nature Communications*, 15(1):1059, 2024. 3
- [56] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. arXiv preprint arXiv:1506.03365, 2015. 5
- [57] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv* preprint arXiv:2206.10789, 2(3):5, 2022. 5