

Faster Diffusion: Rethinking the Role of UNet Encoder in Diffusion Models

Senmao Li^{1*} Taihang Hu^{1*} Fahad Shahbaz Khan^{2,3} Linxuan Li⁴
 Shiqi Yang⁵ Yaxing Wang^{1†} Ming-Ming Cheng¹ Jian Yang¹

¹VCIP,CS, Nankai University, ²Mohamed bin Zayed University of AI, ³Linkoping University

⁴Harbin Engineering University, ⁵Universitat Autònoma de Barcelona

{senmaonk,hutaihang00}@gmail.com fahad.khan@liu.se {yaxing,csjyang}@nankai.edu.cn

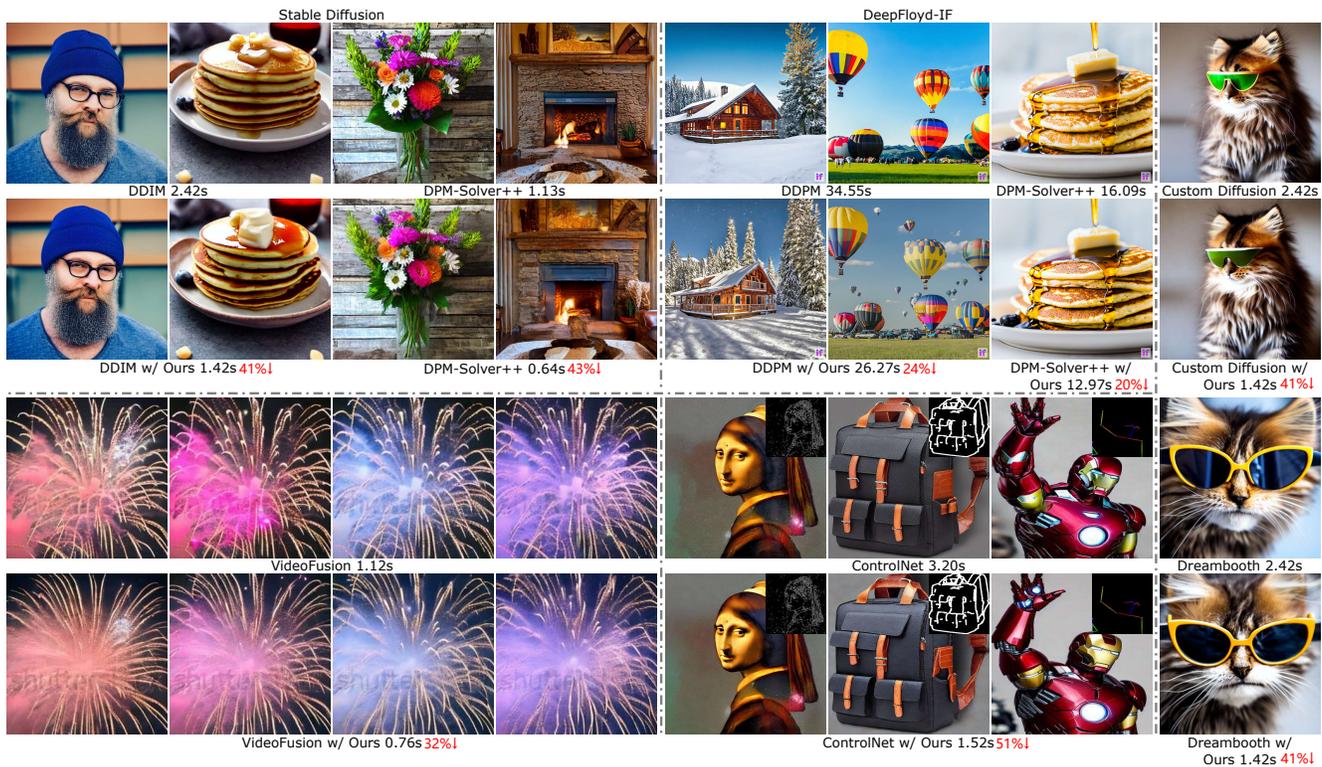


Figure 1. Our method for a diverse set of generation tasks. We could significantly increase the speed (second/image) of image generation.

Abstract

One of the key components within diffusion models is the UNet for noise prediction. While several works have explored basic properties of the UNet decoder, its encoder largely remains unexplored. In this work, we conduct the first comprehensive study of the UNet encoder. We empirically analyze the encoder features and provide insights to important questions regarding their changes at the inference process. In particular, we find that encoder features change gently, whereas the decoder features exhibit substantial variations across different time-steps. This finding inspired us to omit the encoder at certain adjacent time-

steps and reuse cyclically the encoder features in the previous time-steps for the decoder. Further based on this observation, we introduce a simple yet effective encoder propagation scheme to accelerate the diffusion sampling for a diverse set of tasks. By benefiting from our propagation scheme, we are able to perform in parallel the decoder at certain adjacent time-steps. Additionally, we introduce a prior noise injection method to improve the texture details in the generated image. Besides the standard text-to-image task, we also validate our approach on other tasks: text-to-video, personalized generation and reference-guided generation. Without utilizing any knowledge distillation technique, our approach accelerates both the Stable Diffusion (SD) and the DeepFloyd-IF models sampling by 41% and 24% respectively, while maintaining high-quality genera-

*Equal contribution.

†Corresponding author.

tion performance. Our code is available in *FasterDiffusion*.

1. Introduction

One of the popular paradigms in image generation, Diffusion Models [17, 38, 45, 53] have recently achieved significant breakthroughs in various domains, including text-to-video generation [10, 20, 33, 52, 54], personalized image generation [9, 11, 12, 22, 40] and reference-guided image generation [26, 32, 36, 57, 58].

The performance of diffusion models [17, 38, 41] somewhat relies on the UNet [39] denoising network as well as large-scale image datasets. Here, the UNet consists of an encoder \mathcal{E} , a bottleneck \mathcal{B} , and a decoder \mathcal{D} . The features obtained from the encoder \mathcal{E} are skipped to the decoder \mathcal{D} . Benefiting from the properties of the UNet architecture, several recent works [36, 43, 46, 48, 55–57] explore using the decoder features of the UNet for a diverse set of tasks. PnP [48] performs text-guided image-to-image translation by leveraging the decoder features. DIFT [46] finds an emergent correspondence phenomenon that mainly exists in the decoder features. ControlNet [57] fine-tunes an additional encoder, initialized with the encoder of the pretrained UNet, and preserves the original decoder. T2I-Adapter [36] injects the conditional information with a trainable encoder into the encoder of UNet, and preserves the initial decoder. These works indicate the crucial role of the decoder in generating high-quality images.

While the above-mentioned works explore the decoder of the UNet in diffusion models, little effort has been made to investigate the role of the UNet encoder. In this work, we pay a close attention to the characteristics of the encoder of the UNet in the pretrained text-guided diffusion model (e.g., the SD and DeepFloyd-IF). Here, an open research question is whether features from different layers have different contributions during inference time. Based on our analysis presented in Sec 3.2, we discover that UNet encoder features change gently (Fig. 2a) and have a high degree of similarity (Fig. 3 (left)), whereas the decoder features exhibit substantial variations across different time-steps (Fig. 2a and Fig. 3 (right)). Based on these findings, a natural question is how to effectively re-use cyclically the encoder features from the previous time-steps for the decoder as they barely change.

Following the aforementioned questions and with the aim to speed up the diffusion sampling, we introduce *encoder propagation* a simple yet effective propagation scheme based on encoder feature reusing during the diffusion sampling. Specifically, in certain time-steps, we do not use the encoder to get encoder features, which are the input of decoder. Instead, we directly reuse the encoder features from the previous time-steps since during these time-steps the encoder features change minimally. We show that the

proposed propagation scheme accelerates the SD sampling by 24% , and DeepFloyd-IF sampling by 18%. Furthermore, we can use the same encoder features in previous time-step as the input to the decoder of multiple later time-steps, which makes it possible to conduct multiple time-steps decoding concurrently. This parallel procession further accelerates SD sampling by 41%, and DeepFloyd-IF sampling by 24%. While the propagation scheme improves the efficiency of the diffusion sampling, we observe it to lead to some loss of texture information in the generated results (see Fig. 6 (left, middle)). To alleviate this issue, we introduce a prior noise injection strategy to preserve the texture details in the generated images. With these contributions, our proposed method achieves improved sampling efficiency along with maintaining high generation quality. We evaluate the effectiveness of our approach across a wide range of conditional diffusion-based tasks, including text-to-video generation (e.g., Text2Video-zero [20] and VideoFusion [33]), personalized image generation (e.g., Dreambooth [40]) and reference-guided image generation (e.g., ControlNet [57]). To summarize, we make the following contributions:

- 1) We conduct a thorough empirical study of the features of the UNet in the diffusion model.
- 2) Based on our analysis, we propose an encoder propagation scheme to accelerate the diffusion sampling without requiring any distillation technique.
- 3) We further introduce a parallel strategy for diffusion model sampling at adjacent time-steps. We also present a prior noise injection method to improve the image quality.
- 4) Our qualitative and quantitative experiments demonstrate the effectiveness of the proposed method.

2. Related Work

UNet. UNet is originally proposed for medical image segmentation [39]. Moreover, UNet extends to the domain of 3D medical image segmentation [13, 24, 28]. Since then, various UNet-based methods have been developed, often combining it with different architectures [19, 47].

Text-guided diffusion model. Recently, Text-to-image diffusion models [1, 8, 38, 41] have made significant advancements. Notably, Stable Diffusion and DeepFloyd-IF , stand out as two of the most successful diffusion models available within the current open-source community. These models, building upon the UNet architecture, are versatile and can be applied to a wide range of tasks, including image editing [14, 25], super-resolution [7, 50], segmentation [2, 51], and object detection [5, 37].

Diffusion model acceleration. Diffusion models use iterative denoising with UNet for image generation, which can be time-consuming. There are plenty of works diving

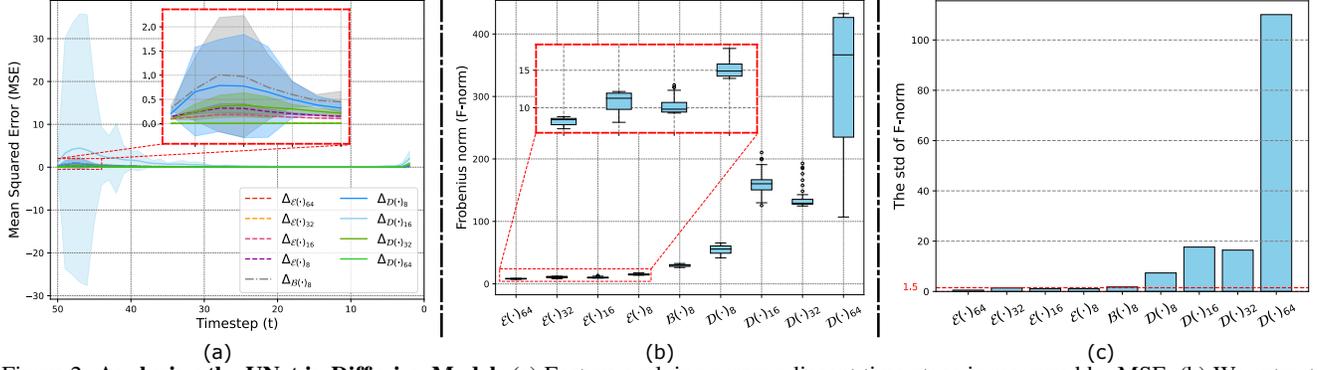


Figure 2. **Analyzing the UNet in Diffusion Model.** (a) Feature evolving across adjacent time-steps is measured by MSE. (b) We extract hierarchical features output of different layers of the UNet at each time-step, average them along the channel dimension to obtain two-dimensional hierarchical features, and then calculate their F-norms. (c) The hierarchical features of the UNet encoder show a lower standard deviation, while those of the decoder exhibit a higher standard deviation.

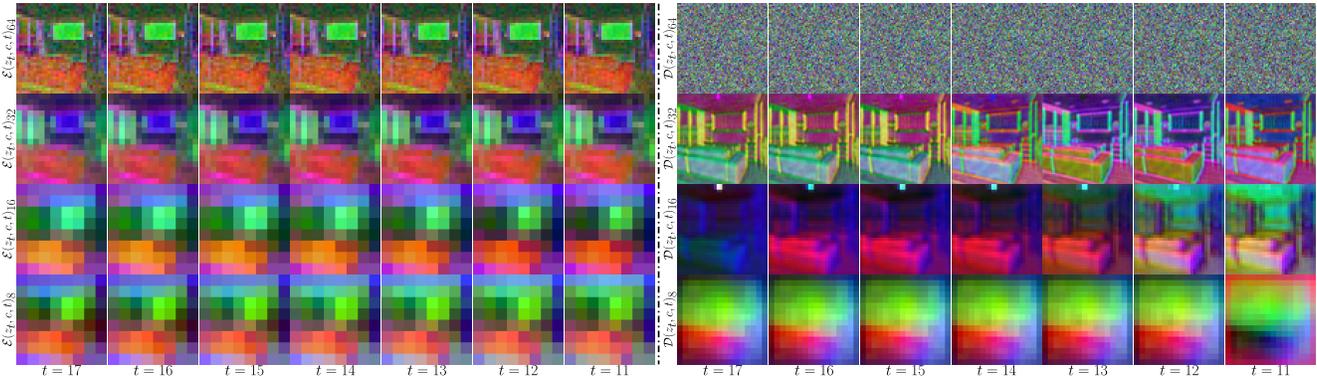


Figure 3. Visualising the hierarchical features. We applied PCA to the hierarchical features following PnP [48], and use the top three leading components as an RGB image for visualization. The encoder features changes gently and have similarity at many time-steps (left), while the decoder features exhibit substantial variations across different time-steps (right).

into this direction trying to address this issue. One approach is step distillation, progressively simplifying existing models to reduce sampling steps [34, 42]. Some recent works combine model compression and distillation for faster sampling [21, 27]. In contrast, efficient diffusion model solvers, such as DDIM [44] and DPM-Solver [30], have significantly reduced sampling steps. ToMe [3, 4] leverages token redundancy to reduce computations required for attention operation [49]. Orthogonal to these methods, we propose a novel method for efficient sampling.

3. Method

We first briefly revisit the architecture of the Stable Diffusion (SD) (Sec. 3.1), and then conduct a comprehensive analysis for the hierarchical features of the UNet (Sec. 3.2). Finally, with our finding, we introduce a novel method to accelerate the diffusion sampling without resorting to knowledge distillation technique (Sec. 3.3).

3.1. Latent Diffusion Model

In the diffusion inference stage, The denoising network ϵ_θ takes as input a text embedding c , a latent code z_t and a

time embedding, predicts noise, resulting in a latent z_{t-1} using the DDIM scheduler [44]:

$$z_{t-1} = \sqrt{\frac{\alpha_{t-1}}{\alpha_t}} z_t + \sqrt{\alpha_{t-1}} \left(\sqrt{\frac{1}{\alpha_{t-1}} - 1} - \sqrt{\frac{1}{\alpha_t} - 1} \right) \cdot \epsilon_\theta(z_t, t, c), \quad (1)$$

where α_t is a predefined scalar function at time-step t ($t = T, \dots, 1$). The denoising network uses a UNet-based architecture. It consists of an encoder \mathcal{E} , a bottleneck \mathcal{B} , and a decoder \mathcal{D} , respectively (Fig. 4b). The hierarchical features extracted from the encoder \mathcal{E} are injected into the decoder \mathcal{D} by a skip connection (Fig. 4a). For convenience of description, we divide the UNet network into specific blocks: $\mathcal{E} = \{\mathcal{E}(\cdot)_s\}$, $\mathcal{B} = \{\mathcal{B}(\cdot)_s\}$, and $\mathcal{D} = \{\mathcal{D}(\cdot)_s\}$, where $s \in \{8, 16, 32, 64\}$ (see Fig. 4b). Both $\mathcal{E}(\cdot)_s$ ¹ and $\mathcal{D}(\cdot)_s$ represent the block layers with input resolution s in both encoder and decoder, respectively.

3.2. Analyzing the UNet in Diffusion Model

In this section, we delve into the UNet which consists of the encoder \mathcal{E} , the bottleneck \mathcal{B} , and the decoder \mathcal{D} , for deeper understanding the different parts in UNet.

¹Once we replace the \cdot with specific inputs in $\mathcal{E}(\cdot)_s$, we define that it represents the feature of $\mathcal{E}(\cdot)_s$

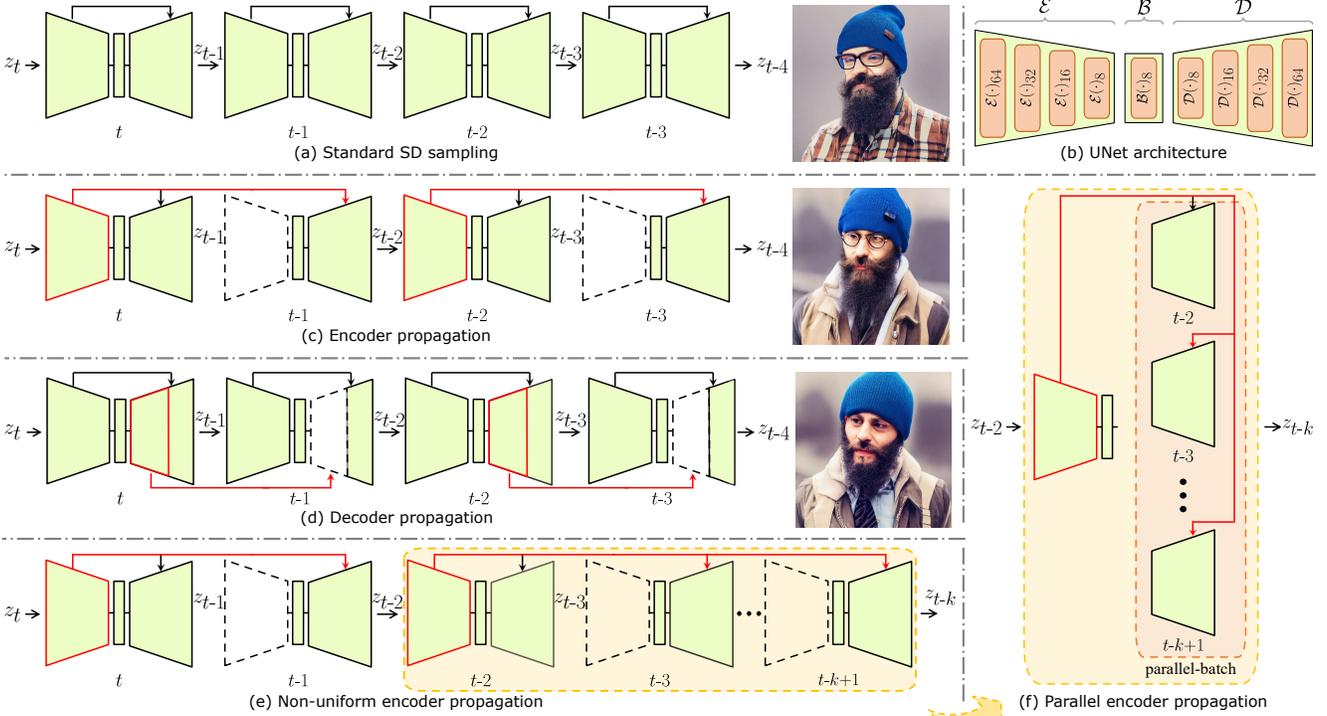


Figure 4. (a) Standard SD sampling. (b) UNet architecture. (c) Encoder propagation. We omit the encoder at certain adjacent time-steps and reuse cyclically the encoder features in the previous time-steps for the decoder. Applying encoder propagation for uniform strategy every two iterations. Note, at time-step $t-1$, predicting noise does not require z_{t-1} (i.e., Eq. 1: $z_{t-2} = \sqrt{\frac{\alpha_{t-2}}{\alpha_{t-1}}} z_{t-1} + \sqrt{\alpha_{t-2}} \left(\sqrt{\frac{1}{\alpha_{t-2}} - 1} - \sqrt{\frac{1}{\alpha_{t-1}} - 1} \right) \cdot \epsilon_{\theta}(z_{t-1}, t-1, c)$). (d) Decoder propagation. The generated images often fail to cover some specific objects in the text prompt. For example, given one prompt case “A man with a beard wearing glasses and a beanie”, this method fails to generate the **glasses** subject. (e) Applying encoder propagation for non-uniform strategy. (f) By benefiting from our propagation scheme, we are able to perform in parallel the decoder at certain adjacent time-steps.

Feature evolving across time-steps. We experimentally observe that the encoder features exhibit a subtle variation at adjacent time-steps, whereas the decoder features exhibit substantial variations across different time-steps (see Fig. 2a and Fig. 3). Specifically, given a pretrained diffusion model, we iteratively produce a latent code z_t (i.e., Eq. 1), and the corresponding hierarchical features: $\{\mathcal{E}(z_t, c, t)_s\}$, $\{\mathcal{B}(z_t, c, t)_s\}$, and $\{\mathcal{D}(z_t, c, t)_s\}$ ($s \in \{8, 16, 32, 64\}$)², as shown in Fig. 4b. We are wondering how change of the hierarchical features at adjacent time-steps. To achieve this goal, we quantify the variation of the hierarchical features as following:

$$\Delta_{\mathcal{E}(\cdot)_s} = \frac{1}{d \times s^2} \|\mathcal{E}(z_t, c, t)_s, \mathcal{E}(z_{t-1}, c, t-1)_s\|_2^2, \quad (2)$$

where d represents the number of channels in $\mathcal{E}(z_t, c, t)_s$. Similarly, we also compute $\Delta_{\mathcal{B}(\cdot)_s}$ and $\Delta_{\mathcal{D}(\cdot)_s}$.

As illustrated in Fig. 2a, for both the encoder \mathcal{E} and the decoder \mathcal{D} , these curves exhibit a similar trend: in the wake of an initial increase, the variation reaches a plateau and

²The feature resolution is half of the previous one in the encoder and two times in the decoder. Note that the feature resolutions of $\mathcal{E}(\cdot)_8$, $\mathcal{B}(\cdot)_8$ and $\mathcal{D}(\cdot)_{64}$ do not change in the SD model.

then decreases, followed by continuing growth towards the end. However, the extent of change in $\Delta_{\mathcal{E}(\cdot)_s}$ and $\Delta_{\mathcal{D}(\cdot)_s}$ is quantitatively markedly different. For example, the maximum value and variance of the $\Delta_{\mathcal{E}(\cdot)_s}$ are less than 0.4 and 0.05, respectively (Fig. 2a (zoom-in area)), while the corresponding values of the $\Delta_{\mathcal{D}(\cdot)_s}$ are about 5 and 30, respectively (Fig. 2a). Furthermore, we find that $\Delta_{\mathcal{D}(\cdot)_{64}}$, the change of the last layer of the decoder, is close to zero. This is due the output of the denoising network is similar at adjacent time-steps [35]. In conclusion, the overall feature change $\Delta_{\mathcal{E}(\cdot)_s}$ is smaller than $\Delta_{\mathcal{D}(\cdot)_s}$ throughout the inference phase.

Feature evolving across layers. We experimentally observe that the feature characteristics are significantly different between the encoder and the decoder across all time-steps. For the encoder \mathcal{E} the intensity of the change is slight, whereas it is vary drastic for the decoder \mathcal{D} . Specifically we calculate the *Frobenius norm* for hierarchical features $\mathcal{E}(z_t, c, t)_s$ across all time-steps, dubbed as $\mathcal{F}_{\mathcal{E}(\cdot)_s} = \{\mathcal{F}_{\mathcal{E}(z_T, c, T)_s}, \dots, \mathcal{F}_{\mathcal{E}(z_1, c, 1)_s}\}$. Similarly, we compute $\mathcal{F}_{\mathcal{B}(\cdot)_s}$ and $\mathcal{F}_{\mathcal{D}(\cdot)_s}$, respectively.

Fig. 2b shows the feature evolving across layers with a boxplot³. Specifically, for $\{\mathcal{F}_{\mathcal{E}(\cdot)_s}\}$ and $\{\mathcal{F}_{\mathcal{B}(\cdot)_s}\}$, the box is relatively compact, with a narrow range between their first quartile and third quartile values. For example, the maximum box height ($\mathcal{F}_{\mathcal{E}(\cdot)_{32}}$) of these features is less than 5 (see Fig. 2b (zoom-in area)). This indicates that the features from both the encoder \mathcal{E} and the bottleneck \mathcal{B} slightly change. In contrast, the box heights corresponding to $\{\mathcal{D}(\cdot)_s\}$ are relatively large. For example, for the $\mathcal{D}(\cdot)_{64}$ the box height is over 150 between the first quartile and third quartile values (see Fig. 2b). Furthermore, we also provide a standard deviation (Fig. 2c), which exhibits similar phenomena like Fig. 2b. These results show that the encoder features have relatively small discrepancy, and a high degree of similarity across all time-steps. However, the decoder features evolve drastically.

Could we omit Encoder at certain time-steps? As indicated by the previous experimental analysis, we observe that, during the denoising process, the decoder features change drastically, whereas the encoder \mathcal{E} features change mildly, and have a high degree similarities at certain adjacent time-steps. Therefore, as shown in Fig. 4c, We propose to omit the encoder at certain time-steps and reuse cyclically the encoder features in the previous time-steps for the decoder. Specifically, we delete the encoder at time-step $t - 1$ ($t - 1 < T$), and the corresponding decoder (including the skip connections) takes as input the hierarchical outputs of the encoder \mathcal{E} from the previous time-step t , instead of the ones from the current time-step $t - 1$ like the standard SD sampling (for more detail, see Sec. 3.3).

When omitting the encoder at certain time-step, we are able to generate similar images (Fig. 4c) like standard SD sampling (Fig. 4a, Tab. 1 (the first and second rows) and additional results on Supp.Mat.C). Alternatively if we use the similar strategy for the decoder (i.e., *decoder propagation*), we find the generated images often fail to cover some specific objects in the text prompt (Fig. 4d). For example, when provided with prompt “A man with a beard wearing glasses and a beanie”, the SD model fails to synthesize “glasses” when applying decoder propagation (Fig. 4d). This is due to the fact that the semantics are mainly contained in the features from the decoder rather than the encoder [55].

The *encoder propagation*, which uses encoder outputs from previous time-step as the input to the current decoder, could speed up the diffusion model sampling at inference time. In the following Sec. 3.3, we give a detailed illustration about encoder propagation.

³Each boxplot contains the minimum (0th percentile), the maximum (100th percentile), the median (50th percentile), the first quartile (25th percentile) and the third quartile (75th percentile) values of the feature Frobenius norm (e.g., $\{\mathcal{F}_{\mathcal{E}(z_{T,c,T})_s}, \dots, \mathcal{F}_{\mathcal{E}(z_{1,c,1})_s}\}$).

3.3. Encoder propagation

Diffusion sampling, combining iterative denoising with transformers, is time-consuming. Therefore we propose a novel and practical diffusion sampling acceleration method. During the diffusion sampling process $t = \{T, \dots, 1\}$, we refer to the time-steps where encoder propagation is deployed, as *non-key* time-steps denoted as $t^{non-key} = \{t_0^{non-key}, \dots, t_{N-1}^{non-key}\}$. The remaining time-steps are dubbed as $t^{key} = \{t_0^{key}, t_1^{key}, \dots, t_{T-1-N}^{key}\}$. In other words, we do not use the encoder at time-steps $t^{non-key}$, and leverage the hierarchical features of the encoder from the time-step t^{key} . Note we utilize the encoder \mathcal{E} at the initial time-step ($t_0^{key} = T$). Thus, the diffusion inference time-steps could be reformulated as $\{t^{key}, t^{non-key}\}$, where $t^{key} \cup t^{non-key} = \{T, \dots, 1\}$ and $t^{key} \cap t^{non-key} = \emptyset$. In the following subsections, we introduce both *uniform encoder propagation* and *non-uniform encoder propagation* strategies.

As shown in Fig. 2a, The encoder feature change is larger in the initial inference phase compared to the later phases throughout the inference process. Therefore, we select more *key* time-steps in the initial inference phase, and less *key* time-steps in later phases. We experimentally define the *key* time-steps as $t^{key} = \{50, 49, 48, 47, 45, 40, 35, 25, 15\}$ for SD model with DDIM, and $t^{key} = \{100, 99, 98, \dots, 92, 91, 90, 85, 80, \dots, 25, 20, 15, 14, 13, \dots, 2, 1\}$, $\{50, 49, \dots, 2, 1\}$ and $\{75, 73, 70, 66, 61, 55, 48, 40, 31, 21, 10\}$ for three stage of DeepFloyd-IF (see detail *key* time-steps selection on Supp.Mat.C).

The remaining time-steps are categorized as *non-key* time-steps. We define this strategy as non-uniform encoder propagation (see Fig. 4e). As shown in Fig. 4c, we also explore the time-step selection with fix stride (e.g, 2), dubbed as *uniform encoder propagation*. Tab. 3 reports the results of the ablation study, considering various combinations of *key* and *non-key* time-steps. These results indicate that the set of *key* time-steps performs better in generating images.

Parallel non-uniform encoder propagation When applying the non-uniform encoder propagation strategy, at time-step $t \in t^{non-key}$ the decoder inputs do not rely on the encoder outputs at t time-step (see Fig. 4e, right part). Instead it relies on the encoder output at the previous nearest *key* time-step. This allows us to perform *parallel non-uniform encoder propagation* at these adjacent time-steps in $t^{non-key}$. As shown in Fig. 4f, we perform decoding in parallel from $t - 2$ to $t - k + 1$ time-steps. This technique further improves the inference efficiency, since the decoder forward in multiple time-steps could be conducted concurrently. We call *these non-key time-steps* as parallel-batch *non-key* time-steps. As shown in Fig. 5 (right), we further

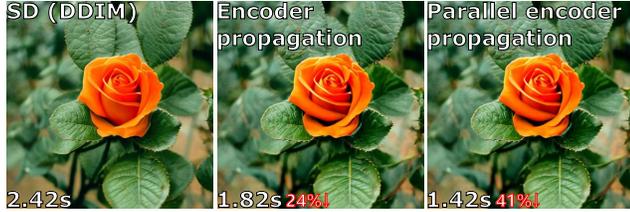


Figure 5. Comparing with SD (left), encoder propagation reduces the sampling time by 24% (middle). Furthermore, parallel encoder propagation achieves a 41% reduction in sampling time (right).

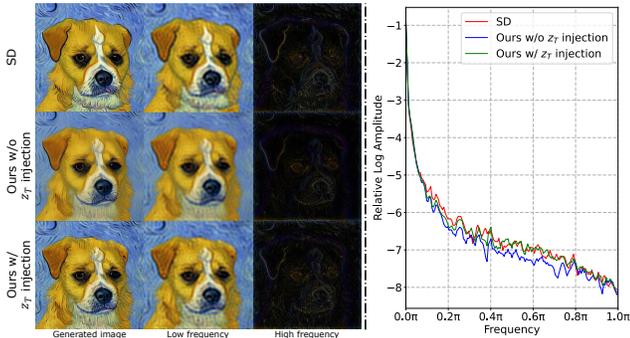


Figure 6. (left) We keep the image content through z_T injection, slightly compensating for the texture information loss caused by encoder propagation. (right) The amplitudes of the generated image through z_T injection closely resemble those from SD.

reduce the evaluation time by 41% for the SD model.

Prior noise injection Although the encoder propagation could improve the efficiency in the inference phase, we observe that it leads to a slight loss of texture information in the generated results (see Fig. 6 (left, middle)). Inspired by related works [6, 23], we propose a prior noise injection strategy. It combines the initial latent code z_T into the generative process at subsequent time-step (i.e., z_t), following $z_t = z_t + \alpha \cdot z_T$, if $t < \tau$, where $\alpha = 0.003$ is the scale parameter to control the impact of z_T . And we start to use this injection mechanism from $\tau = 25$ step. This strategic incorporation successfully improves the texture information. Importantly, it demands almost negligible extra computational resources. We observe that the loss of texture information occurs in all frequencies of the frequency domain (see Fig. 6 (right, red and blue curves)). This approach ensures a close resemblance of generated results in the frequency domain from both SD and z_T injection (see Fig. 6 (right, red and green curves)), with the generated images maintaining the desired fidelity (see Fig. 6 (left, bottom)).

4. Experiments

Datasets and evaluation metrics. We randomly select 10K prompts from the MS-COCO2017 validation dataset [29], and feed them into the diffusion model to obtain 10K generated images. For other tasks, we use the same settings as baselines (e.g., Text2Video-zero [20], VideoFu-

DM	Sampling Method	T	FID ↓	Clip-score ↑	GFLOPs/ image ↓	s/image ↓	
						Unet of DM	DM
Stable Diffusion	DDIM [44]	50	21.75	0.773	37050	2.23	2.42
	DDIM [44] w/ Ours	50	21.08	0.783	27350 _{27%↓}	1.21 _{45%↓}	1.42 _{41%↓}
	DPM-Solver [30]	20	21.36	0.780	14821	0.90	1.14
	DPM-Solver [30] w/ Ours	20	21.25	0.779	11743 _{21%↓}	0.46 _{48%↓}	0.64 _{43%↓}
Stable Diffusion	DPM-Solver++ [31]	20	20.51	0.782	14821	0.90	1.13
	DPM-Solver++ [31] w/ Ours	20	20.76	0.781	11743 _{21%↓}	0.46 _{48%↓}	0.64 _{43%↓}
DeepFloyd-IF	DDIM + ToMe [3]	50	22.32	0.782	35123	2.07	2.26
	DDIM + ToMe [3] w/ Ours	50	20.73	0.781	26053 _{26%↓}	1.15 _{44%↓}	1.33 _{41%↓}
DeepFloyd-IF	DDPM [17]	225	23.89	0.783	734825	33.91	34.55
	DDPM [17] w/ Ours	225	23.73	0.782	626523 _{15%↓}	25.61 _{25%↓}	26.27 _{24%↓}
DeepFloyd-IF	DPM-Solver++ [31]	100	20.79	0.784	370525	15.19	16.09
	DPM-Solver++ [31] w/ Ours	100	20.85	0.785	313381 _{15%↓}	12.02 _{21%↓}	12.97 _{20%↓}

Table 1. Quantitative evaluation in both SD model and DeepFloyd-IF diffusion model on MS-COCO 2017 10K subset.

sion [33], Dreambooth [40] and ControlNet [57]). We use the Fréchet Inception Distance (FID) [16] metric to assess the visual quality of the generated images, and the Clip-score [15] to measure the consistency between image content and text prompt. Furthermore, we report the average values for both the computational workload (GFLOPs/image) and sampling time (s/image) to represent the resource demands for a single image. See more detailed implementation information on Supp.Mat.A. As shown in Code. 1, in the standard SD sampling code, adding 3 lines of code with green comments can achieve encoder propagation.

Code 1. Encoder propagation for SD (DDIM)

```

from diffusers import StableDiffusionPipeline
import torch
from utils import register_parallel_pipeline,
register_faster_forward # 1. import package

model_id = "runwayml/stable-diffusion-v1-5"
pipe =
    StableDiffusionPipeline.from_pretrained(model_id,
torch_dtype=torch.float16)
pipe = pipe.to("cuda")

register_parallel_pipeline(pipe) # 2. enable parallel
register_faster_forward(pipe.unet) # 3. encoder
propagation

prompt = "a photo of an astronaut riding a horse on
mars"
image = pipe(prompt).images[0]
image.save("astronaut_rides_horse.png")

```

4.1. Text-to-image Generation

We first evaluate the proposed encoder propagation method for the standard text-to-image generation task on both the latent space (i.e., SD) and pixel space (i.e., DeepFloyd-IF) diffusion models. As shown on Tab. 1, we significantly accelerate the diffusion sampling with negligible performance degradation. Specifically, our proposed method

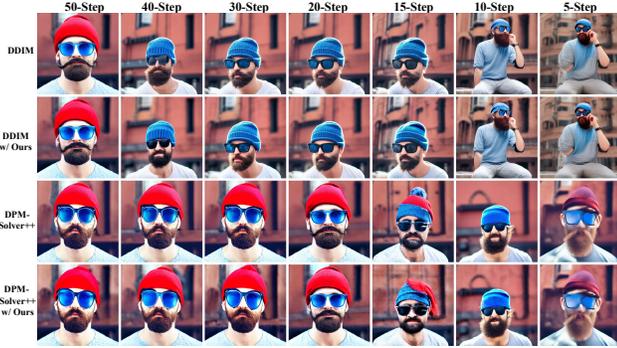


Figure 7. Generated images at different time-steps, with the prompt “A man with a beard wearing sunglasses and a beanie”.

Method	T	FID ↓	Clip-score ↑	GFLOPs/ image ↓	Unet of SD	s/image ↓ SD
Text2Video-zero [20]	50	-	0.732	39670	12.59/8	13.65/8
Text2Video-zero [20] w/ Ours	50	-	0.731	30690 _{22%↓}	9.46/8 _{25%↓}	10.54/8 _{23%↓}
VideoFusion [33]	50	-	0.700	224700	16.71/16	17.93/16
VideoFusion [33] w/ Ours	50	-	0.700	148680 _{33%↓}	11.10/16 _{34%↓}	12.27/16 _{32%↓}
ControlNet [57] (edges)	50	13.78	0.769	49500	3.09	3.20
ControlNet [57] (edges) w/ Ours	50	14.65	0.767	31400 _{37%↓}	1.43 _{54%↓}	1.52 _{51%↓}
ControlNet [57] (scribble)	50	16.17	0.775	56850	3.85	3.95
ControlNet [57] (scribble) w/ Ours	50	16.42	0.775	35990 _{37%↓}	1.83 _{53%↓}	1.93 _{51%↓}
Dreambooth [40]	50	-	0.640	37050	2.23	2.42
Dreambooth [40] w/ Ours	50	-	0.660	27350 _{27%↓}	1.21 _{45%↓}	1.42 _{41%↓}
Custom Diffusion [22]	50	-	0.640	37050	2.21	2.42
Custom Diffusion [22] w/ Ours	50	-	0.650	27350 _{27%↓}	1.21 _{45%↓}	1.42 _{41%↓}

Table 2. Quantitative evaluation on text-to-video, personalized generation and reference-guided generation tasks.

decreases the computational burden (GFLOPs) by a large margin (27%), and greatly reduces sampling time to 41% when compared to standard DDIM sampling in SD. Similarly, in DeepFloyd-IF, the reduction in both computational burden and time (s) reaches 15% and 24%, respectively. Furthermore, our method could be combined with the latest sampling techniques like DPM-Solver [30], DPM-Solver++ [31], and ToMe [3]. Our method enhances sampling efficiency while preserving good model performance, with the negligible variations of both FID and Clipscore values (Tab. 1 (the third to eighth rows)). And our method achieves good performance across different sampling steps (Fig. 7). Importantly, these results show that our method is orthogonal and compatible with these acceleration techniques. As shown in Fig. 1, we visualize the generated images with different sampling techniques. Our method still generates high-quality results (see Supp.Mat.C. for additional results).

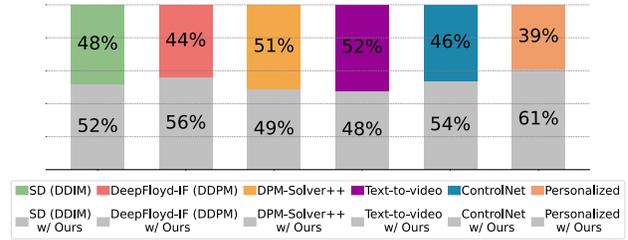


Figure 8. User study results.

4.2. Other tasks with text-guided diffusion model

Besides the standard text-to-image task, we also validate our proposed approach on other tasks: *text-to-video generation*, *personalized generation* and *reference-guided image generation*.

text-to-video generation, (i.e., Text2Video-zero [20] and VideoFusion [33]), *personalized generation* (i.e., Dreambooth [40]) and *reference-guided image generation* (i.e., ControlNet [57]).

Text-to-video To evaluate our method, we combine our method with both Text2Video-zero [20] and VideoFusion [33]. As reported in Tab. 2 (the first and fourth rows), when combined with our method, the two methods have the reduction of approximately 22% to 33% in both computational burden and generation time (s). These results indicate that we are able to enhance the efficiency of generative processes in the text-to-video task while preserving video fidelity at the same time (Fig. 1 (left, bottom)). As an example, when generating a video using the prompt “Fireworks bloom in the night sky”, the VideoFusion model takes 17.92 seconds with 16 frames for the task (1.12s/frame), when combining with our method it only takes 12.27s (0.76s/frame) to generate a high-quality video (Fig. 1 (left, bottom)).

Personalized image generation Dreambooth [40] and Custom Diffusion [22] are two approaches for customizing tasks by fine-tuning text-to-image diffusion models. As reported in Tab. 2 (the ninth to twelfth rows), our method, working in conjunction with the two customization approaches, accelerates the image generation and reduces computational demands. Visually, it maintains the ability to generate images with specific contextual relationships based on reference images. (Fig. 1 (right))

Reference-guided image generation ControlNet [57] incorporates a trainable encoder, successfully generates a text-guided image, and preserves similar content with conditional information. Our approach can be applied concurrently to two encoders of ControlNet. In this paper, we validate the proposed method with two conditional controls: *edge* and *scribble*. Tab. 2 (the fifth to eighth row) reports

Propagation strategy	FID ↓ Clipscore ↑		GFLOPs /image ↓	s/image ↓		FTC ↓
				Unet of SD	SD	
SD	21.75	0.773	37050	2.23	2.42	68.1
Uniform	II	$t^{key} = \{50, 48, 46, 44, 42, 40, 38, 36, 34, 32, 30, 28, 26, 24, 22, 20, 18, 16, 14, 12, 10, 8, 6, 4, 2\}$				
		22.38	0.788	31011 _{16%↓}	1.62 _{27%↓}	1.81 _{25%↓}
	III	$t^{key} = \{50, 44, 38, 32, 26, 20, 14, 8, 2\}$				
		21.54	0.779	27350 _{27%↓}	1.26 _{43%↓}	1.46 _{40%↓}
IIII	$t^{key} = \{50, 38, 26, 14, 2\}$					
	24.61	0.766	26370 _{29%↓}	1.12 _{50%↓}	1.36 _{44%↓}	43.1
Non-uniform	II	$t^{key} = \{50, 40, 39, 38, 30, 25, 20, 15, 5\}$				
		22.94	0.776	27350 _{27%↓}	1.26 _{43%↓}	1.42 _{41%↓}
	III	$t^{key} = \{50, 30, 25, 20, 15, 14, 5, 4, 3\}$				
		35.25	0.742	27350 _{27%↓}	1.25 _{43%↓}	1.42 _{41%↓}
	IIII	$t^{key} = \{50, 41, 37, 35, 22, 21, 18, 14, 5\}$				
		22.14	0.778	27350 _{27%↓}	1.22 _{45%↓}	1.42 _{41%↓}
IV (Ours)	21.08	0.783	27350 _{27%↓}	1.21 _{45%↓}	1.42 _{41%↓}	38.2

Table 3. Quantitative evaluation in various propagation strategies on MS-COCO 2017 10K subset. FTC=FID×Time/Clipscore.

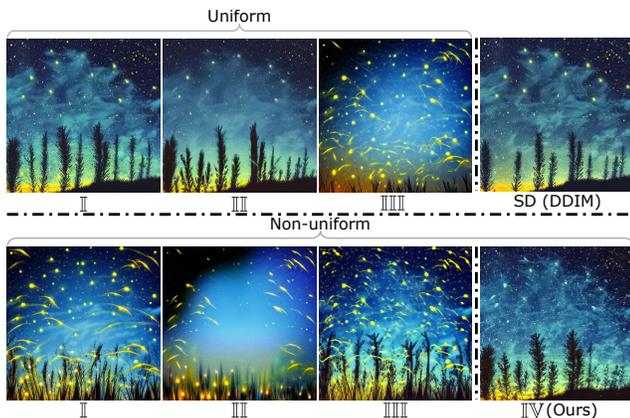


Figure 9. Generating image with uniform and non-uniform encoder propagation. The result of uniform strategy III yields smooth and loses textual compared with SD. Both uniform strategy IIIII and non-uniform strategy II, III and IIIII generate images with unnatural saturation levels.

quantitative results. We observe that it leads to a significant decrease in both generation time and computational burden. Furthermore, Fig. 1 (middle, bottom) qualitatively shows that our method successfully preserves the given structure information and achieves similar results as ControlNet.

User study We conducted a user study, as depicted in Fig. 8, and ask subjects to select results. We apply pairwise comparisons (forced choice) with 18 users (35 pairs images or videos/user). The results demonstrate that our method performs equally well as the baseline methods.

Sampling Method	FID ↓	Clipscore ↑
SD (DDIM [17])	21.75	0.773
SD (DDIM [17]) + Ours w/o z_T injection	21.71	0.779
SD (DDIM [17]) + Ours w/ z_T injection	21.08	0.783

Table 4. Quantitative evaluation for prior noise injection.

4.3. Ablation study

We ablate the results with different selections of both uniform and non-uniform encoder propagation. Tab. 3 reports that the performance of the non-uniform setting largely outperforms the uniform one in terms of both FID and Clipscore (see Tab. 3 (the third and eighth rows)). Furthermore, we explore different configurations within the non-uniform strategy. The strategy, using the set of *key* time-steps we established, yields better results in the generation process (Tab. 3 (the eighth row)). We further present qualitative results stemming from the above choices. As shown in Fig. 9, given the same number of *key* time-steps, the appearance of nine-step non-uniform strategy II, III and IIIII settings do not align with the prompt “Fireflies dot the night sky”. Although the generated image in the two-step setting exhibits a pleasing visual quality, its sampling efficiency is lower than our chosen setting (see Tab. 3 (the second and eighth rows)).

Effectiveness of prior noise injection. We evaluate effectiveness of injecting initial z_T . As reported in Tab. 4, the differences in FID and Clipscores without z_T (the second row), when compared to DDIM and Ours (the first and third rows), are approximately 0.01%, which can be considered negligible. While it is not the case for the visual expression of the generated image, it is observed that the output contains complete semantic information with smoothing texture (refer to Fig. 6 (left, the second row)). Injecting the z_T aids in maintaining fidelity in the generated results during encoding propagation (see Fig. 6 (left, the third row) and Fig. 6 (right, red and green curves)).

5. Conclusion

In this work, We explore the characteristics of the encoder and decoder in UNet of the text-to-image diffusion model and find that the encoder is negligible at many time-steps, while the decoder plays a significant role across all time-steps. Building upon this finding, we propose encoder propagation for efficient diffusion sampling, reducing time on both Stable Diffusion and DeepFloyd-IF on diverse set of generation tasks. We conduct extensive experiments and validate that our approach can achieve improved sampling efficiency while maintaining image quality. **Limitations:** Although our approach achieves efficient diffusion sampling, it faces challenges in generating quality when using a limited number of sampling steps (e.g., 5).

References

- [1] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. [2](#)
- [2] Dmitry Baranchuk, Ivan Rubachev, Andrey Voynov, Valentin Khruikov, and Artem Babenko. Label-efficient semantic segmentation with diffusion models. *arXiv preprint arXiv:2112.03126*, 2021. [2](#)
- [3] Daniel Bolya and Judy Hoffman. Token merging for fast stable diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4598–4602, 2023. [3](#), [6](#), [7](#), [1](#)
- [4] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your vit but faster. *arXiv preprint arXiv:2210.09461*, 2022. [3](#)
- [5] Shoufa Chen, Peize Sun, Yibing Song, and Ping Luo. Diffusiondet: Diffusion model for object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19830–19843, 2023. [2](#)
- [6] 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. [6](#)
- [7] Hyungjin Chung, Byeongsu Sim, and Jong Chul Ye. Come-closer-diffuse-faster: Accelerating conditional diffusion models for inverse problems through stochastic contraction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12413–12422, 2022. [2](#)
- [8] Xiaoliang Dai, Ji Hou, Chih-Yao Ma, Sam Tsai, Jialiing Wang, Rui Wang, Peizhao Zhang, Simon Vandenhende, Xiaofang Wang, Abhimanyu Dubey, et al. Emu: Enhancing image generation models using photogenic needles in a haystack. *arXiv preprint arXiv:2309.15807*, 2023. [2](#)
- [9] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022. [2](#)
- [10] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features for consistent video editing. *arXiv preprint arxiv:2307.10373*, 2023. [2](#)
- [11] Jing Gu, Yilin Wang, Nanxuan Zhao, Tsu-Jui Fu, Wei Xiong, Qing Liu, Zhifei Zhang, He Zhang, Jianming Zhang, Hyun-Joon Jung, and Xin Eric Wang. Photoswap: Personalized subject swapping in images, 2023. [2](#)
- [12] Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris Metaxas, and Feng Yang. Svdiff: Compact parameter space for diffusion fine-tuning. *arXiv preprint arXiv:2303.11305*, 2023. [2](#)
- [13] Ali Hatamizadeh, Yucheng Tang, Vishwesh Nath, Dong Yang, Andriy Myronenko, Bennett Landman, Holger Roth, and Daguang Xu. Unetr: Transformers for 3d medical image segmentation. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 574–584, 2022. [2](#)
- [14] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. [2](#), [4](#), [11](#)
- [15] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021. [6](#)
- [16] 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. [6](#)
- [17] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. [2](#), [6](#), [8](#), [3](#), [7](#)
- [18] Ziqi Huang, Tianxing Wu, Yuming Jiang, Kelvin C.K. Chan, and Ziwei Liu. ReVersion: Diffusion-based relation inversion from images. *arXiv preprint arXiv:2303.13495*, 2023. [4](#), [11](#)
- [19] Debesh Jha, Pia H Smedsrud, Michael A Riegler, Dag Johansen, Thomas De Lange, Pål Halvorsen, and Håvard D Johansen. Resunet++: An advanced architecture for medical image segmentation. In *2019 IEEE international symposium on multimedia (ISM)*, pages 225–2255. IEEE, 2019. [2](#)
- [20] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. *arXiv preprint arXiv:2303.13439*, 2023. [2](#), [6](#), [7](#), [1](#), [3](#), [4](#), [8](#)
- [21] Bo-Kyeong Kim, Hyoung-Kyu Song, Thibault Castells, and Shinkook Choi. On architectural compression of text-to-image diffusion models. *arXiv preprint arXiv:2305.15798*, 2023. [3](#)
- [22] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1931–1941, 2023. [2](#), [7](#), [1](#), [3](#), [4](#), [10](#)
- [23] Mingi Kwon, Jaeseok Jeong, and Youngjung Uh. Diffusion models already have a semantic latent space. *arXiv preprint arXiv:2210.10960*, 2022. [6](#)
- [24] Van-Linh Le and Olivier Saut. Rrc-unet 3d for lung tumor segmentation from ct scans of non-small cell lung cancer patients. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2316–2325, 2023. [2](#)
- [25] Senmao Li, Joost van de Weijer, Taihang Hu, Fahad Shahbaz Khan, Qibin Hou, Yaxing Wang, and Jian Yang. Stylediffusion: Prompt-embedding inversion for text-based editing. *arXiv preprint arXiv:2303.15649*, 2023. [2](#)
- [26] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. In *Pro-*

- ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22511–22521, 2023. 2
- [27] Yanyu Li, Huan Wang, Qing Jin, Ju Hu, Pavlo Chemerys, Yun Fu, Yanzhi Wang, Sergey Tulyakov, and Jian Ren. Snap-fusion: Text-to-image diffusion model on mobile devices within two seconds. *arXiv preprint arXiv:2306.00980*, 2023. 3
- [28] Junjie Liang, Cihui Yang, Jingting Zhong, and Xiaoli Ye. Btswin-unet: 3d u-shaped symmetrical swin transformer-based network for brain tumor segmentation with self-supervised pre-training. *Neural processing letters*, 55(4): 3695–3713, 2023. 2
- [29] 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. 6, 1
- [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, 6, 7, 1
- [31] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022. 6, 7, 1, 3
- [32] Shilin Lu, Yanzhu Liu, and Adams Wai-Kin Kong. Tf-icon: Diffusion-based training-free cross-domain image composition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2294–2305, 2023. 2
- [33] Zhengxiong Luo, Dayou Chen, Yingya Zhang, Yan Huang, Liang Wang, Yujun Shen, Deli Zhao, Jingren Zhou, and Tieniu Tan. Videofusion: Decomposed diffusion models for high-quality video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10209–10218, 2023. 2, 6, 7, 1, 3, 4, 9
- [34] 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
- [35] Daiki Miyake, Akihiro Iohara, Yu Saito, and Toshiyuki Tanaka. Negative-prompt inversion: Fast image inversion for editing with text-guided diffusion models. *arXiv preprint arXiv:2305.16807*, 2023. 4
- [36] Chong Mou, Xintao Wang, Liangbin Xie, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023. 2
- [37] Walter HL Pinaya, Mark S Graham, Robert Gray, Pedro F Da Costa, Petru-Daniel Tudosiu, Paul Wright, Yee H Mah, Andrew D MacKinnon, James T Teo, Rolf Jager, et al. Fast unsupervised brain anomaly detection and segmentation with diffusion models. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 705–714. Springer, 2022. 2
- [38] 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. 2
- [39] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: 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. 2
- [40] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023. 2, 6, 7, 1, 3, 4, 10
- [41] 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. 2
- [42] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022. 3
- [43] Yujun Shi, Chuhui Xue, Jiachun Pan, Wenqing Zhang, Vincent YF Tan, and Song Bai. Dragdiffusion: Harnessing diffusion models for interactive point-based image editing. *arXiv preprint arXiv:2306.14435*, 2023. 2
- [44] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 3, 6, 4
- [45] 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. 2
- [46] Luming Tang, Menglin Jia, Qianqian Wang, Cheng Perng Phoo, and Bharath Hariharan. Emergent correspondence from image diffusion. *arXiv preprint arXiv:2306.03881*, 2023. 2
- [47] Dmitrii Torbunov, Yi Huang, Haiwang Yu, Jin Huang, Shinjae Yoo, Meifeng Lin, Brett Viren, and Yihui Ren. Uvcgan: Unet vision transformer cycle-consistent gan for unpaired image-to-image translation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 702–712, 2023. 2
- [48] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1921–1930, 2023. 2, 3
- [49] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3
- [50] Jianyi Wang, Zongsheng Yue, Shangchen Zhou, Kelvin CK Chan, and Chen Change Loy. Exploiting diffusion prior

- for real-world image super-resolution. *arXiv preprint arXiv:2305.07015*, 2023. [2](#)
- [51] Julia Wolleb, Robin Sandkühler, Florentin Bieder, Philippe Valmaggia, and Philippe C Cattin. Diffusion models for implicit image segmentation ensembles. In *International Conference on Medical Imaging with Deep Learning*, pages 1336–1348. PMLR, 2022. [2](#)
- [52] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7623–7633, 2023. [2](#)
- [53] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Yingxia Shao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *arXiv preprint arXiv:2209.00796*, 2022. [2](#)
- [54] Shuai Yang, Yifan Zhou, Ziwei Liu, and Chen Change Loy. Rerender a video: Zero-shot text-guided video-to-video translation. *arXiv preprint arXiv:2306.07954*, 2023. [2](#)
- [55] Guanqi Zhan, Chuanxia Zheng, Weidi Xie, and Andrew Zisserman. What does stable diffusion know about the 3d scene? *arXiv preprint arXiv:2310.06836*, 2023. [2](#), [5](#)
- [56] Junyi Zhang, Charles Herrmann, Junhwa Hur, Luisa Polania Cabrera, Varun Jampani, Deqing Sun, and Ming-Hsuan Yang. A tale of two features: Stable diffusion complements dino for zero-shot semantic correspondence. *arXiv preprint arXiv:2305.15347*, 2023.
- [57] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. [2](#), [6](#), [7](#), [1](#), [3](#), [4](#), [10](#)
- [58] Yuxin Zhang, Weiming Dong, Fan Tang, Nisha Huang, Haibin Huang, Chongyang Ma, Tong-Yee Lee, Oliver Deussen, and Changsheng Xu. Prospect: Expanded conditioning for the personalization of attribute-aware image generation. *arXiv preprint arXiv:2305.16225*, 2023. [2](#)

Supplementary Material

In the supplemental material, we provide a detailed description of the experimental implementation (Sec. A). Subsequently, an analysis of the parameter quantities for the encoder and decoder is conducted (Sec. B). Following that, we present additional experiments and results, including more detailed ablation studies and comparative experiments (Sec. C).

A. Implementation Details

A.1. Configure

We use the Stable Diffusion v1.5 pre-trained model⁴ and DeepFloyd-IF⁵. All of our experiments are conducted using an A40 GPU (48GB of VRAM).

We randomly select 100 captions from the MS-COCO 2017 validation dataset [29] as prompts for generating images. The analytical results presented in Sec. 3 are based on the statistical outcomes derived from these 100 generated images.

A.2. Details about the layers in the UNet

The UNet in Stable Diffusion (SD) consists of an encoder \mathcal{E} , a bottleneck \mathcal{B} , and a decoder \mathcal{D} , respectively. We divide the UNet into specific blocks: $\mathcal{E} = \{\mathcal{E}(\cdot)_s\}$, $\mathcal{B} = \{\mathcal{B}(\cdot)_s\}$, and $\mathcal{D} = \{\mathcal{D}(\cdot)_s\}$, where $s \in \{8, 16, 32, 64\}$. $\mathcal{E}(\cdot)_s$ and $\mathcal{D}(\cdot)_s$ represent the block layers with input resolution s in the encoder and decoder, respectively. Tab. 5 presents detailed information about the block architecture.

A.3. Time and memory consumption ratios

We report the run time and GPU memory consumption ratios for text-to-image task. As shown in Tab. 6, we significantly accelerate the diffusion sampling with encoder propagation while maintaining a comparable memory demand to the baselines (Tab. 6 (the last two columns)). Specifically, our proposed method reduces the spending time (s/image) by 24% and require a little additional memory compared to standard DDIM sampling in SD (DDIM vs. Ours: 2.62GB vs. 2.64GB). The increased GPU memory requirement is for caching the features of the encoder from the previous time-step. Though applying parallel encoder propagation results in an increase in memory requirements by 51%, it leads to a more remarkable acceleration of 41% (DDIM vs. Ours: 2.62GB vs. 3.95GB). In conclusion, applying encoder propagation reduces the sampling time, accompanied by a negligible increase in memory requirements. Parallel encoder propagation on text-to-image tasks yields a sampling speed improvement of 20% to 43%, requiring an additional acceptable amount of memory.

⁴<https://huggingface.co/runwayml/stable-diffusion-v1-5>

⁵<https://github.com/deep-floyd/IF>

Besides the standard text-to-image task, we also validate our proposed approach on other tasks: *text-to-video generation* (i.e., Text2Video-zero [20] and VideoFusion [33]), *personalized generation* (i.e., Dreambooth [40] and Custom Diffusion [22]) and *reference-guided image generation* (i.e., ControlNet [57]). We present the time and memory consumption ratios for these tasks in Tab. 7.

As reported in Tab. 7 (top), when combined with our method, there is a reduction in sampling time by 23% and 32% for Text2Video-zero [20] and VideoFusion [33], respectively, while the memory requirements increased slightly by 3% (0.2GB) and 0.9% (0.11GB). The time spent by reference-guided image generation (i.e., ControlNet [57]) is reduced by more than 20% with a negligible increase in memory (1%). When integrated with our parallel encoder propagation, the sampling time in this task can be reduced by more than half (51%) (Tab. 7 (middle)). Dreambooth [40] and Custom Diffusion [22] are two approaches for customizing tasks by fine-tuning text-to-image diffusion models. As reported in Tab. 7 (bottom), our method, working in conjunction with the two customization approaches, accelerates the image generation with an acceptable increase in memory.

Note that our method, either utilizing encoder propagation or parallel encoder propagation, improves sampling speed without compromising image quality (Sec. 4). We conducted GPU memory consumption ratios using the official method provided by PyTorch, `torch.cuda.max_memory_allocated`⁶, which records the peak allocated memory since the start of the program.

A.4. GFLOPs

We use the `fvcore`⁷ library to calculate the GFLOPs required for a single forward pass of the diffusion model. Multiplying this by the number of sampling steps gives us the total computational burden for sampling one image. Additionally, using `fvcore`, we can determine the computational load required for each layer of the diffusion model. Based on the key time-steps set in our experiments, we subtract the computation we save from the original total computational load, which then represents the GFLOPs required by our method. Similarly, `fvcore` also supports parameter count statistics.

A.5. Baseline Implementations

For the comparisons of text-to-image generation in Sec. 4, we use the official implementation of DPM-Solver [30], DPM-Solver++ [31]⁸, and ToMe [3]⁹. For the other tasks

⁶https://pytorch.org/docs/stable/generated/torch.cuda.max_memory_allocated.html

⁷<https://github.com/facebookresearch/fvcore>

⁸<https://github.com/LuChengTHU/dpm-solver>

⁹<https://github.com/dbolya/tomesd>

UNet	Layer number	Type of layer	Layer name	Input resolution of layer	Output resolution of layer	
\mathcal{E}	$\mathcal{E}(\cdot)_{64}$	0	resnets	down_blocks.0.resnets.0	(320, 64 , 64)	(320, 64, 64)
		1	attention	down_blocks.0.attentions.0	(320, 64, 64)	(320, 64, 64)
		2	resnet	down_blocks.0.resnets.1	(320, 64, 64)	(320, 64, 64)
		3	attention	down_blocks.0.attentions.1	(320, 64, 64)	(320, 64, 64)
		4	downsamplers	down_blocks.0.downsamplers.0	(320, 64, 64)	(320, 32, 32)
	$\mathcal{E}(\cdot)_{32}$	5	resnet	down_blocks.1.resnets.0	(320, 32 , 32)	(640, 32, 32)
		6	attention	down_blocks.1.attentions.0	(640, 32, 32)	(640, 32, 32)
		7	resnet	down_blocks.1.resnets.1	(640, 32, 32)	(640, 32, 32)
		8	attention	down_blocks.1.attentions.1	(640, 32, 32)	(640, 32, 32)
		9	downsamplers	down_blocks.1.downsamplers.0	(640, 32, 32)	(640, 16, 16)
	$\mathcal{E}(\cdot)_{16}$	10	resnet	down_blocks.2.resnets.0	(640, 16 , 16)	(1280, 16, 16)
		11	attention	down_blocks.2.attentions.0	(1280, 16, 16)	(1280, 16, 16)
		12	resnet	down_blocks.2.resnets.1	(1280, 16, 16)	(1280, 16, 16)
		13	attention	down_blocks.2.attentions.1	(1280, 16, 16)	(1280, 16, 16)
		14	downsamplers	down_blocks.2.downsamplers.0	(1280, 16, 16)	(1280, 8, 8)
	$\mathcal{E}(\cdot)_8$	15	resnet	down_blocks.3.resnets.0	(1280, 8 , 8)	(1280, 8, 8)
16		resnet	down_blocks.3.resnets.1	(1280, 8, 8)	(1280, 8, 8)	
\mathcal{B}	$\mathcal{B}(\cdot)_8$	17	resnet	mid_blocks.resnets.0	(1280, 8 , 8)	(1280, 8, 8)
		18	attention	mid_blocks.attentions.0	(1280, 8, 8)	(1280, 8, 8)
		19	resnet	mid_blocks.resnets.1	(1280, 8, 8)	(1280, 8, 8)
\mathcal{D}	$\mathcal{D}(\cdot)_8$	20	resnet	up_blocks.0.resnets.0	(1280, 8 , 8)	(1280, 8, 8)
		21	resnet	up_blocks.0.resnets.1	(1280, 8, 8)	(1280, 8, 8)
		22	resnet	up_blocks.0.resnets.2	(1280, 8, 8)	(1280, 8, 8)
		23	upsamplers	up_blocks.0.upsamplers.0	(1280, 8, 8)	(1280, 16, 16)
	$\mathcal{D}(\cdot)_{16}$	24	resnet	up_blocks.1.resnets.0	(1280+1280, 16 , 16)	(1280, 16, 16)
		25	attention	up_blocks.1.attentions.0	(1280, 16, 16)	(1280, 16, 16)
		26	resnet	up_blocks.1.resnets.1	(1280+1280, 16, 16)	(1280, 16, 16)
		27	attention	up_blocks.1.attentions.1	(1280, 16, 16)	(1280, 16, 16)
		28	resnet	up_blocks.1.resnets.2	(1280+640, 16, 16)	(1280, 16, 16)
		29	attention	up_blocks.1.attentions.2	(1280, 16, 16)	(1280, 16, 16)
		30	upsamplers	up_blocks.1.upsamplers.0	(1280, 16, 16)	(1280, 32, 32)
	$\mathcal{D}(\cdot)_{32}$	31	resnet	up_blocks.2.resnets.0	(1280+640, 32 , 32)	(640, 32, 32)
		32	attention	up_blocks.2.attentions.0	(640, 32, 32)	(640, 32, 32)
		33	resnet	up_blocks.2.resnets.1	(640+640, 32, 32)	(640, 32, 32)
34		attention	up_blocks.2.attentions.1	(640, 32, 32)	(640, 32, 32)	
35		resnet	up_blocks.2.resnets.2	(640+320, 32, 32)	(640, 32, 32)	
36		attention	up_blocks.2.attentions.2	(640, 32, 32)	(640, 32, 32)	
37		upsamplers	up_blocks.2.upsamplers.0	(640, 32, 32)	(640, 64, 64)	
$\mathcal{D}(\cdot)_{64}$	38	resnet	up_blocks.3.resnets.0	(640+320, 64 , 64)	(320, 64, 64)	
	39	attention	up_blocks.3.attentions.0	(320, 64, 64)	(320, 64, 64)	
	40	resnet	up_blocks.3.resnets.1	(320+320, 64, 64)	(320, 64, 64)	
	41	attention	up_blocks.3.attentions.1	(320, 64, 64)	(320, 64, 64)	
	42	resnet	up_blocks.3.resnets.2	(320+320, 64, 64)	(320, 64, 64)	
	43	attention	up_blocks.3.attentions.2	(320, 64, 64)	(320, 64, 64)	

Table 5. Detailed information about the layers of the encoder \mathcal{E} , bottleneck \mathcal{B} and decoder \mathcal{D} in the UNet.

with text-guided diffusion model, we use the official imple-

mentation of Text2Video-zero [20]¹⁰, VideoFusion [33]¹¹,

¹⁰<https://github.com/Picsart-AI-Research/Text2Video-Zero>

DM	Sampling Method	T	s/image	memory (GB)
Stable Diffusion	DDIM [44]		2.42	2.62
	DDIM [44] w/ Ours	†	1.82 _{24%↓}	2.64
		‡	1.42 _{41%↓}	3.95
	DPM-Solver [30]		1.14	2.62
	DPM-Solver [30] w/ Ours	†	0.92 _{19%↓}	2.64
		‡	0.64 _{43%↓}	2.69
	DPM-Solver++ [31]		1.13	2.61
	DPM-Solver++ [31] w/ Ours	†	0.91 _{19%↓}	2.65
‡		0.64 _{43%↓}	2.68	
DDIM + ToMe [3]		2.26	2.62	
DDIM + ToMe [3] w/ Ours	†	1.72 _{24%↓}	2.64	
	‡	1.33 _{41%↓}	3.95	
DeepFloyd-IF	DDPM [17]		34.55	40.5
	DDPM [17] w/ Ours	†	29.45 _{15%↓}	41.1
		‡	26.27 _{24%↓}	41.1
	DPM-Solver++ [31]		16.09	40.5
DPM-Solver++ [31] w/ Ours	†	14.13 _{12%↓}	40.8	
	‡	12.97 _{20%↓}	40.8	

Table 6. Time and GPU memory consumption ratios in both SD model and DeepFloyd-IF diffusion model. †: Encoder propagation, ‡: Parallel encoder propagation.

Sampling Method	T	s/image	memory (GB)
Text2Video-zero [20]		13.65	6.59
Text2Video-zero [20] w/ Ours	†	10.54 _{23%↓}	6.79
	‡	—	—
VideoFusion [33]		17.93	11.87
VideoFusion [33] w/ Ours	†	12.27 _{32%↓}	11.98
	‡	—	—
ControlNet [57] (edges)		3.20	3.81
ControlNet [57] (edges) w/ Ours	†	2.01 _{37%↓}	3.85
	‡	1.52 _{51%↓}	5.09
ControlNet [57] (scribble)		3.95	3.53
ControlNet [57] (scribble) w/ Ours	†	3.18 _{20%↓}	3.57
	‡	1.93 _{51%↓}	4.45
Dreambooth [40]		2.42	2.61
Dreambooth [40] w/ Ours	†	1.81 _{24%↓}	2.65
	‡	1.42 _{41%↓}	3.93
Custom Diffusion [22]		2.42	2.61
Custom Diffusion [22] w/ Ours	†	1.82 _{24%↓}	2.64
	‡	1.42 _{41%↓}	3.94

Table 7. Time and GPU memory consumption ratios in text-to-video, personalized generation and reference-guided generated tasks. †: Encoder propagation, ‡: Parallel encoder propagation.

ControlNet [57]¹², Dreambooth [40]¹³, and Custom Diffusion [22]¹⁴. We maintain the original implementations of

¹¹https://huggingface.co/docs/diffusers/api/pipelines/text_to_video

¹²<https://github.com/lllyasviel/ControlNet>

¹³<https://github.com/google/dreambooth>

¹⁴<https://github.com/adobe-research/custom-diffusion>

	Parameter (billion)	FLOPs (million)
$\mathcal{E} + \mathcal{B}$	0.25 + 0.097	224.2+6.04
\mathcal{D}	0.52 _{1.47×}	504.4 _{2.2×}

Table 8. Model complexity comparison regarding the encoder \mathcal{E} , the bottleneck \mathcal{B} and the decoder \mathcal{D} in terms of parameter count and FLOPs.

these baselines and directly integrate code into their existing implementations to implement our method.

B. Parameter Count and FLOPs of SD

We take into account model complexity in terms of parameter count and FLOPs (see Tab. 8). It’s noteworthy that the decoder \mathcal{D} exhibits a significantly greater parameter count, totaling 0.51 billion. This figure is approximately 1.47 times the number of parameter combinations for the encoder \mathcal{E} (250 million) and the bottleneck \mathcal{B} (97 million). This substantial parameter discrepancy suggests that the decoder \mathcal{D} carries a more substantial load in terms of model complexity.

Furthermore, when we consider the computational load, during a single forward inference pass of the SD model, the decoder \mathcal{D} incurs a considerable 504.4 million FLOPs. This value is notably higher, approximately 2.2 times, than the cumulative computational load of the encoder \mathcal{E} , which amounts to 224.2 million FLOPs, and the bottleneck \mathcal{B} , which require 6.04 million FLOPs. This observation strongly implies that the decoder \mathcal{D} plays a relatively more important role in processing and transforming the data within the UNet architecture, emphasizing its critical part to the overall functionality of the model.

C. Ablation experiments and additional results

C.1. The definition of *key* time-steps in various tasks

We utilize 50, 20 and 20 time-steps for DDIM scheduler [44], DPM-solver [30] and DPM-solver++ [31], respectively. We follow the official implementation of DeepFloyd-IF, setting the time-steps for the three sampling stages to 100, 50, and 75, respectively.

Text-to-image generation. We experimentally define the *key* time-steps as $t^{key} = \{50, 49, 48, 47, 45, 40, 35, 25, 15\}$ for SD model with DDIM [44], and $t^{key} = \{100, 99, 98, \dots, 92, 91, 90, 85, 80, \dots, 25, 20, 15, 14, 13, \dots, 2, 1\}$, $\{50, 49, \dots, 2, 1\}$ and $\{75, 73, 70, 66, 61, 55, 48, 40, 31, 21, 10\}$ for three stages of DeepFloyd-IF with DDPM [17]. For SD with both DPM-Solver [30] and DPM-Solver++ [31], we experimentally set the *key* time-steps to $t^{key} = \{20, 19, 18, 17, 15, 10, 5\}$.

Other tasks with text-guided diffusion model. In addi-

tion to standard text-to-image tasks, we further validate our approach on other tasks with text-guided diffusion model (Sec. 4.2). These tasks are all based on the SD implementation of DDIM [44]. Through experiments, we set the *key* time-steps to $t^{key} = \{50, 49, 48, 47, 45, 40, 35, 25, 15\}$ for Text2Video-zero [20], VideoFusion [33], and ControlNet [57]. For personalized tasks (i.e., Dreambooth [40] and Custom Diffusion [22]), we set the *key* time steps to $t^{key} = \{50, 49, 48, 47, 45, 40, 35, 25, 15, 10\}$.

C.2. The effectiveness of encoder propagation

In Sec 3.2, we have demonstrated that encoder propagation (Fig. 10c) can preserve semantic consistency with standard SD sampling (Fig. 10a). However, images generated through decoder propagation (Fig. 10d) often fail to match certain specific objects mentioned in the text prompt.

We extended this strategy to include encoder and decoder propagation (Fig. 10e) as well as decoder and encoder dropping (Fig. 10f). Similarly, encoder and decoder propagation often fail to cover specific objects mentioned in the text prompt, leading to a degradation in the quality of the generated results (Fig. 10e). On the other hand, decoder and encoder dropping is unable to completely denoise, resulting in the generation of images with noise (Fig. 10f).

C.3. Additional results

In Figs. 11, 12 and 13, we show additional text-to-image generation results. Further results regarding other tasks with text-guided diffusion model are illustrated in Figs. 14, 15 and 16.

C.4. Additional tasks

ReVersion [18] is a relation inversion method that relies on the Stable Diffusion (SD). Our method retains the capacity to generate images with specific relations based on exemplar images, as illustrated in Fig. 17 (top).

P2P [14] is an image editing method guided solely by text. When combined with our approach, it enhances sampling efficiency while preserving the editing effect (Fig. 17 (bottom)).

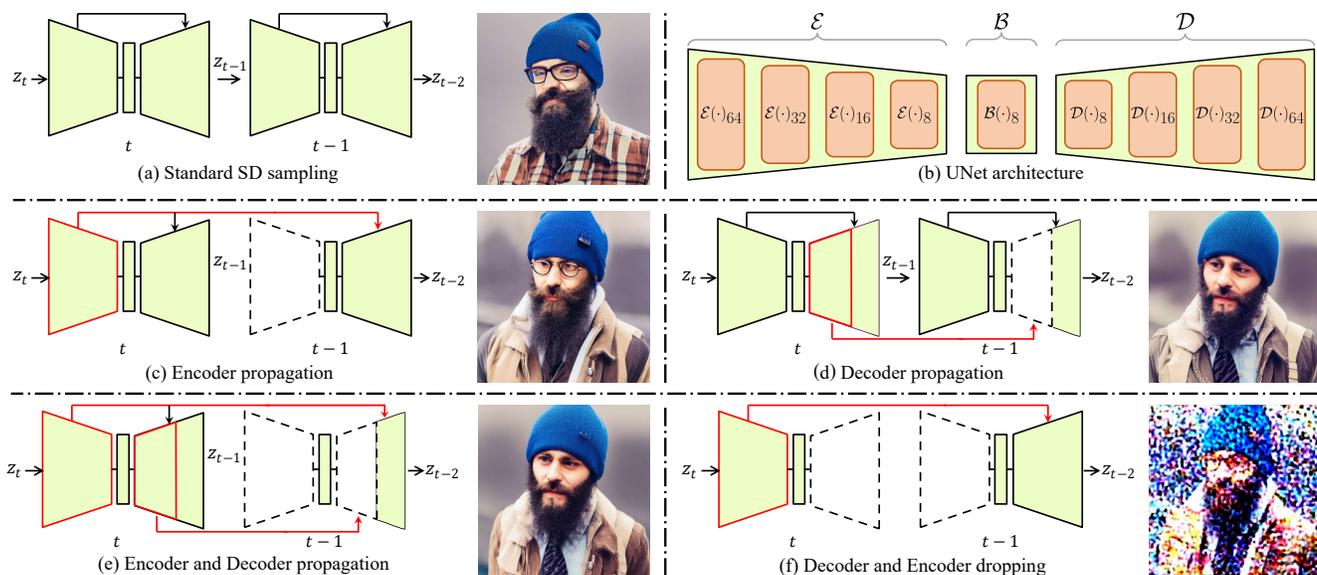


Figure 10. (a) Standard SD sampling. (b) UNet architecture. (c) Encoder propagation. (d) Decoder propagation. (e) Encoder and Decoder propagation. (f) Decoder and Encoder dropping.

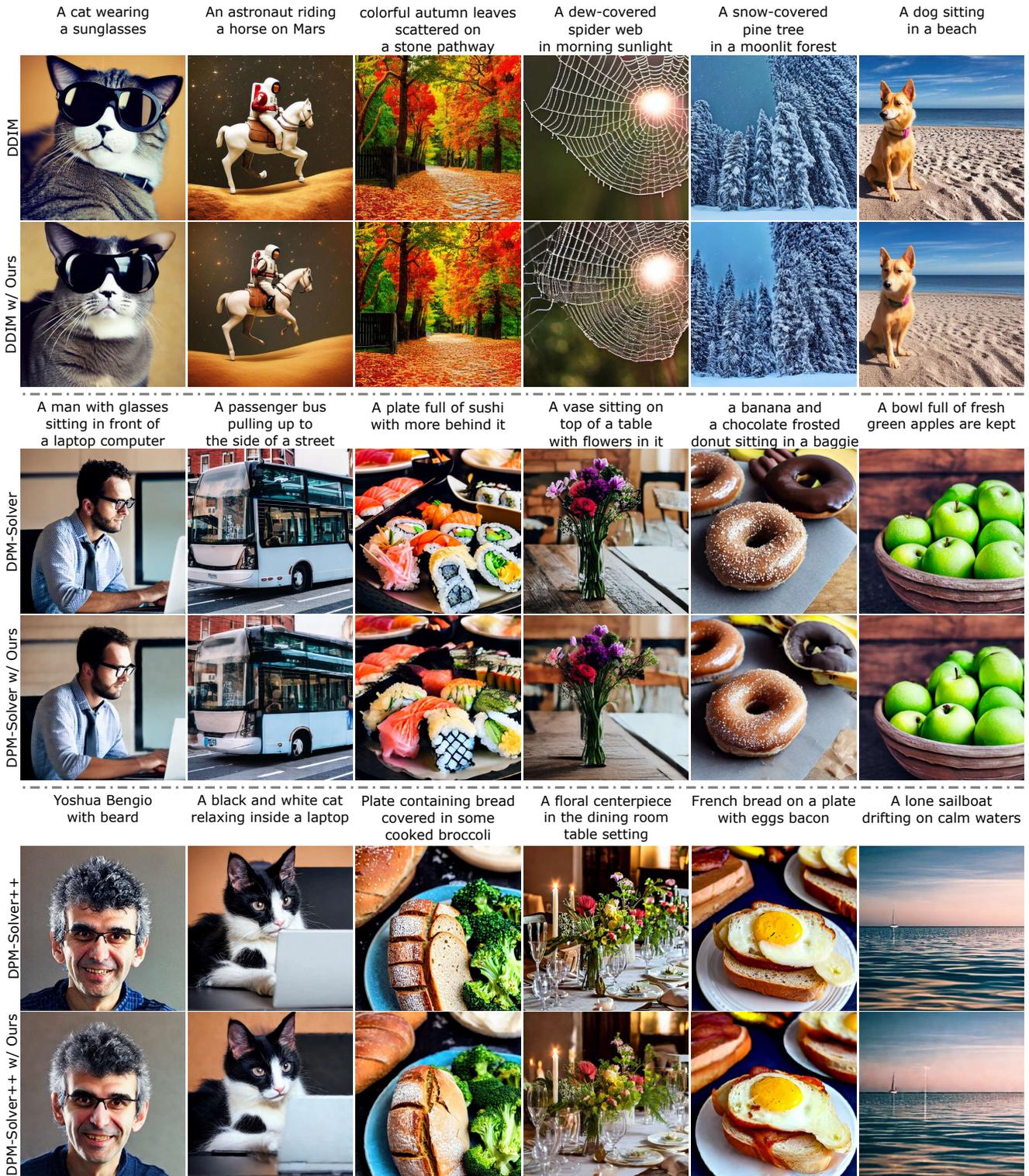


Figure 11. Additional results of text-to-image generation combining SD with DDIM [44], DPM-Solver [30], DPM-Solver++ [31], and these methods in conjunction with our proposed approach.

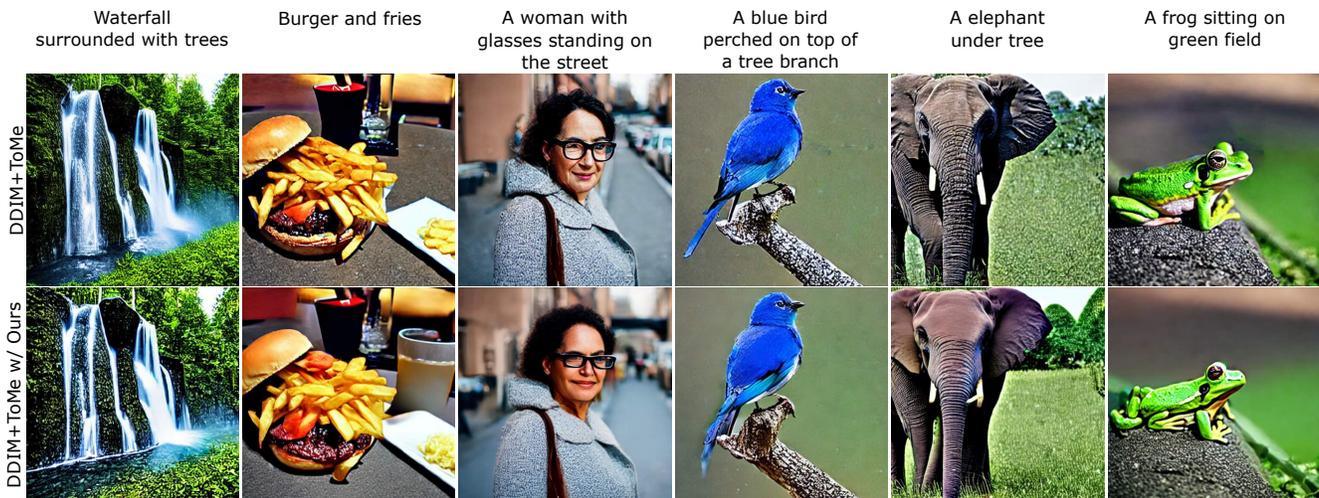


Figure 12. Additional results of text-to-image generation combining SD with DDIM+ToMe [3], and this method in conjunction with our proposed approach.

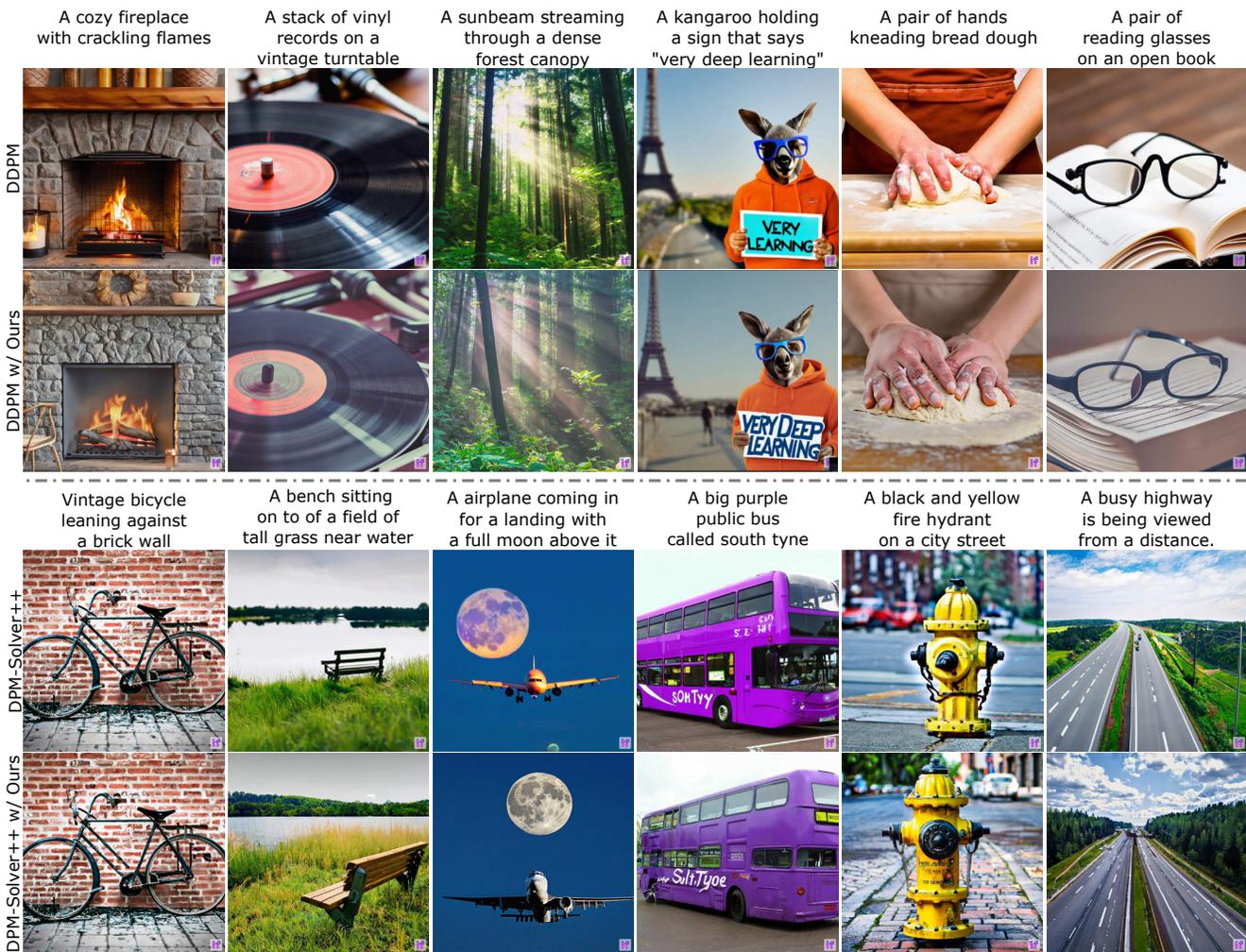


Figure 13. Additional results of text-to-image generation combining DeepFloyd-IF with DDPM [17] and DPM-Solver++ [31], and these methods in conjunction with our proposed approach.

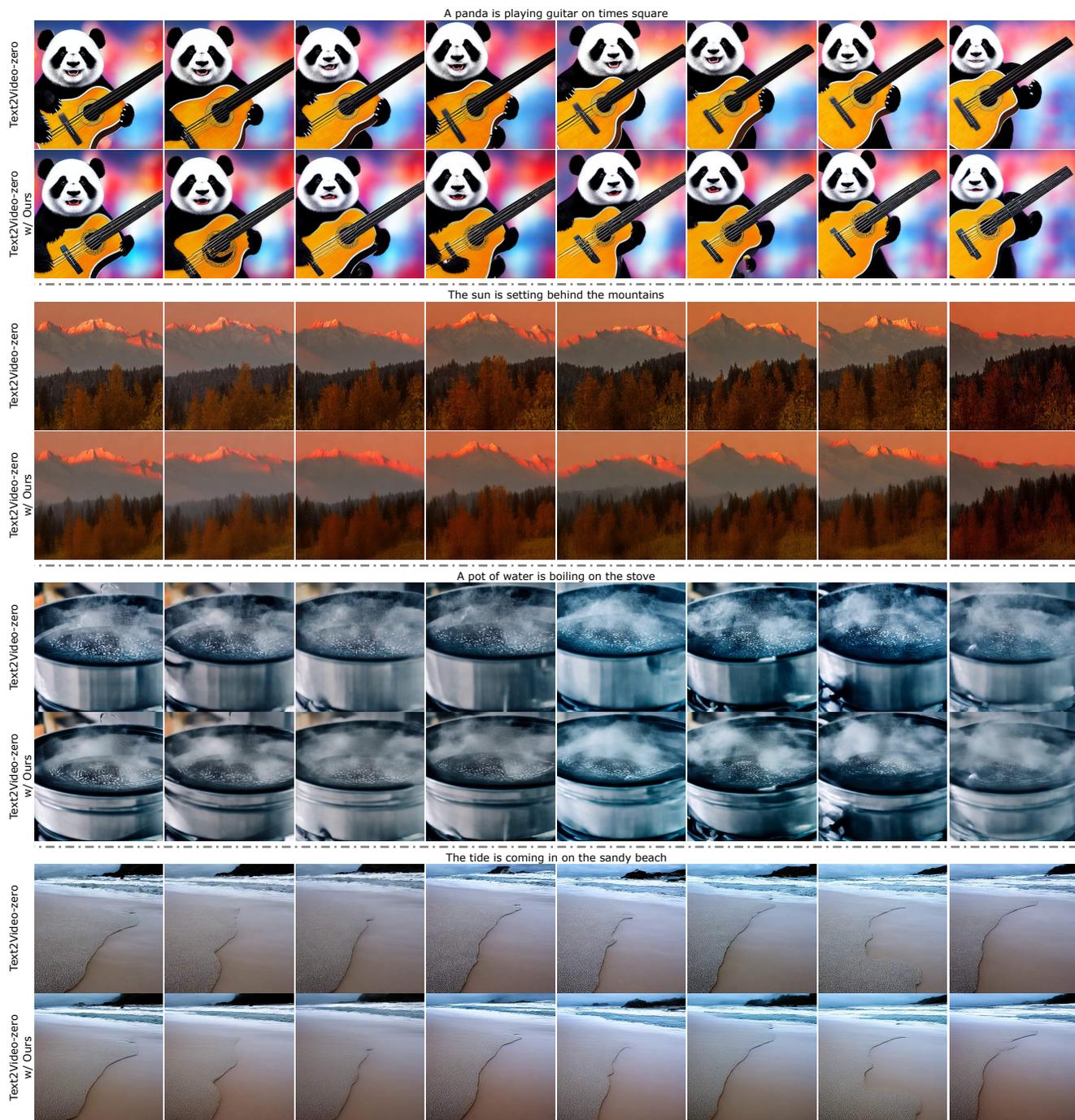


Figure 14. Additional results of Text2Video-zero [20] both independently and when combined with our proposed method.

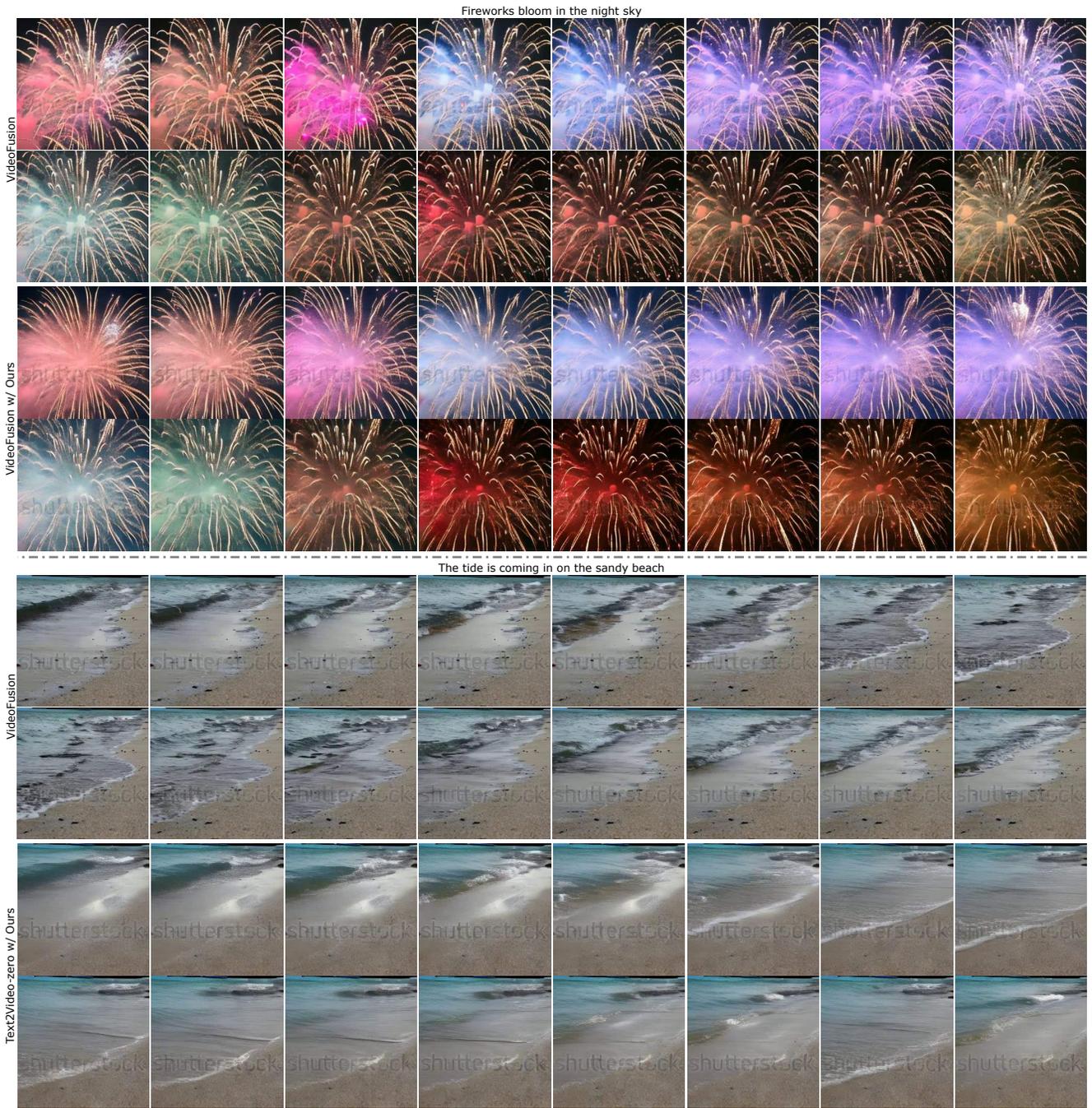


Figure 15. Additional results of VideoFusion [33] both independently and when combined with our proposed method.

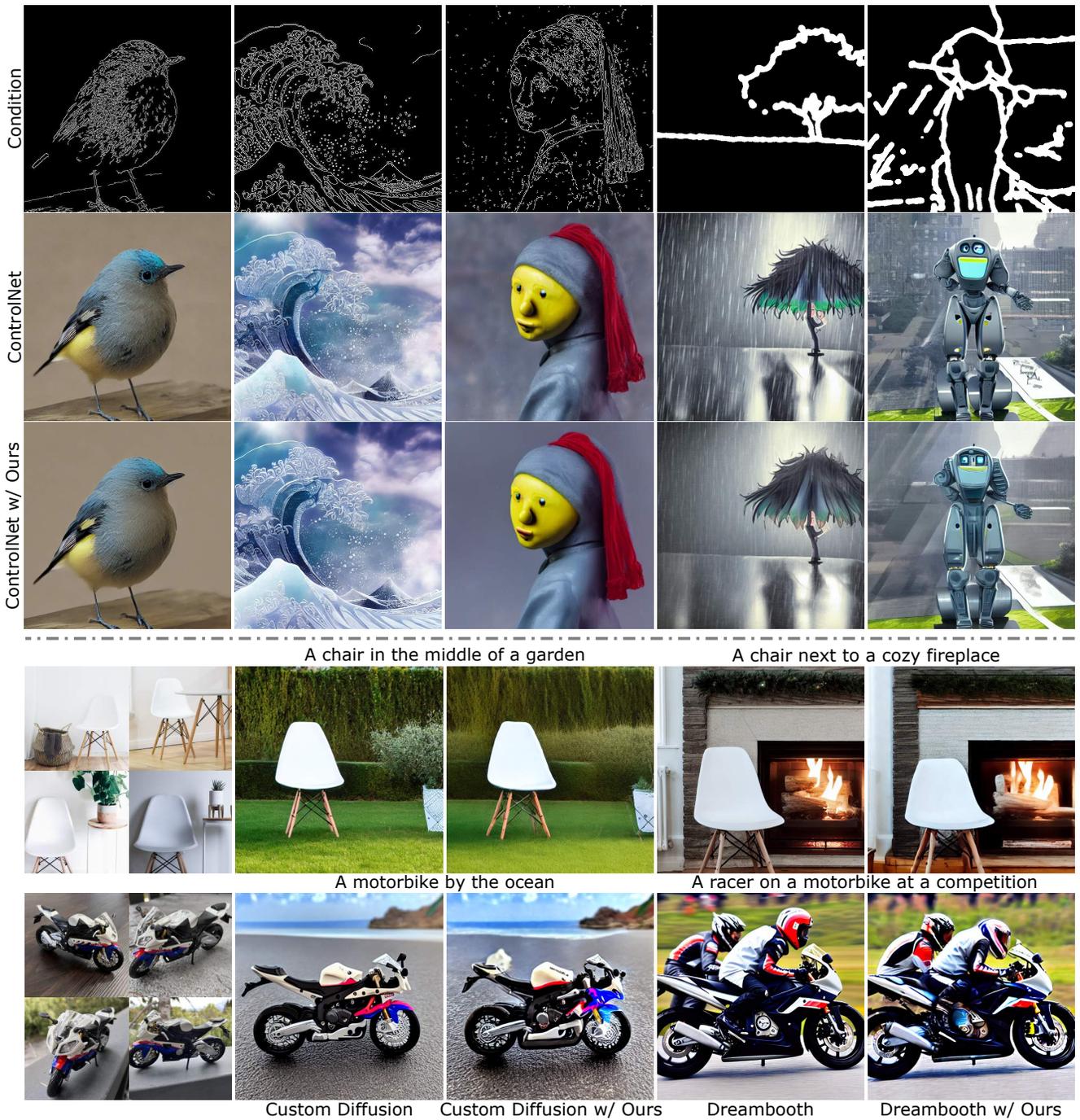


Figure 16. Additional results of ControlNet [57] (top) and personalized tasks [22, 40] (bottom) obtained both independently and in conjunction with our proposed method..

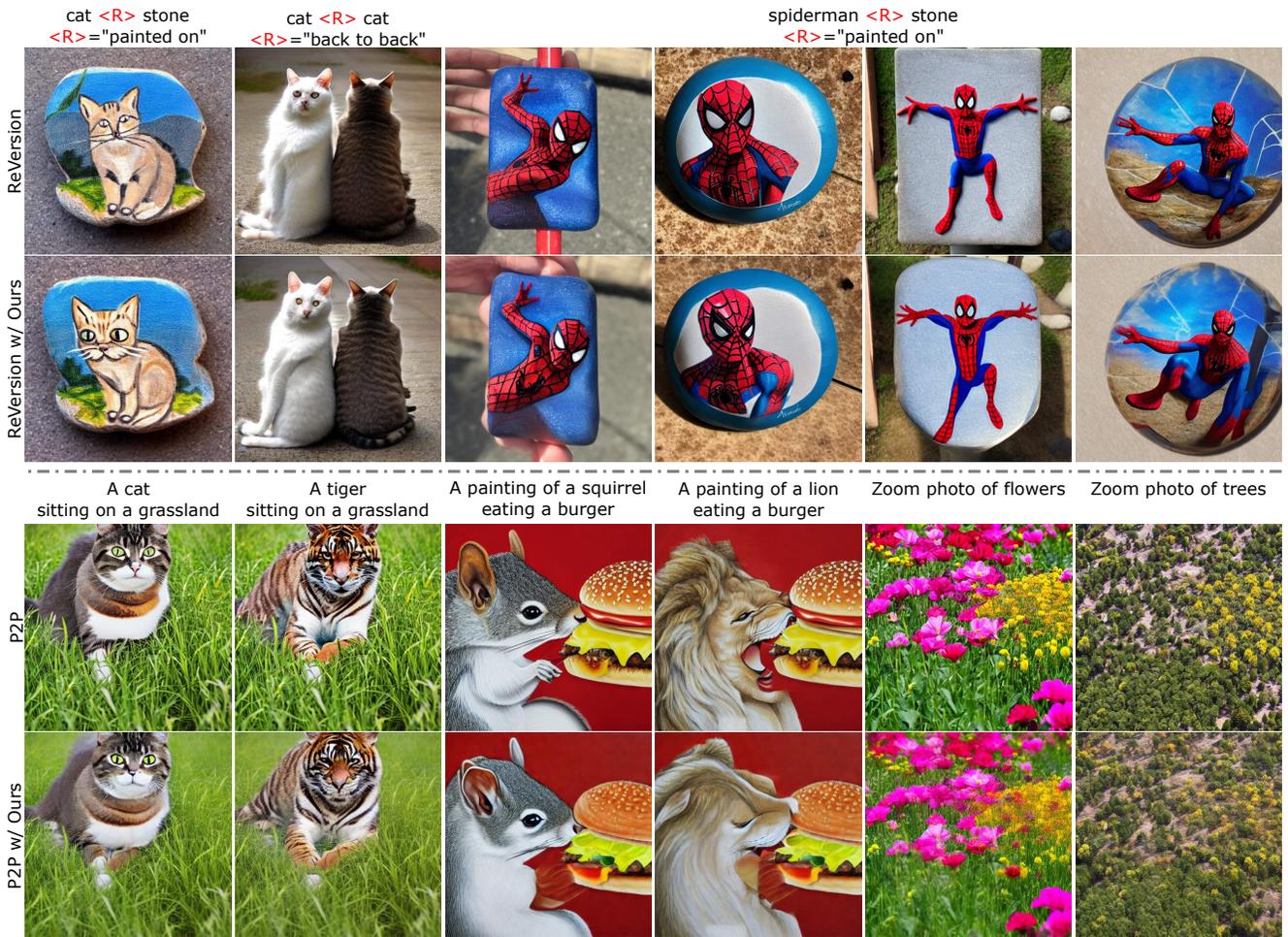


Figure 17. Edited results of ReVersion [18] (top) and P2P [14] (bottom) obtained both independently and in conjunction with our proposed method.