

---

# On Architectural Compression of Text-to-Image Diffusion Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Exceptional text-to-image (T2I) generation results of Stable Diffusion models  
2 (SDMs) come with substantial computational demands. To resolve this issue, re-  
3 cent research on efficient SDMs has prioritized reducing the number of sampling  
4 steps and utilizing network quantization. Orthogonal to these directions, this study  
5 highlights the power of classical architectural compression for general-purpose T2I  
6 synthesis by introducing a block-removed knowledge-distilled SDM (BK-SDM).  
7 We eliminate several residual and attention blocks from the U-Net of SDMs, obtain-  
8 ing over a 30% reduction in the number of parameters, MACs per sampling step,  
9 and latency. We conduct distillation-based pretraining with only 0.22M LAION  
10 pairs (fewer than 0.1% of the full training pairs) on a single A100 GPU. Despite  
11 being trained with limited resources, our compact models can imitate the original  
12 SDM by benefiting from transferred knowledge and achieve competitive results  
13 against larger multi-billion parameter models on the zero-shot MS-COCO bench-  
14 mark. Moreover, we demonstrate the applicability of our lightweight pretrained  
15 models in personalized generation with DreamBooth finetuning.

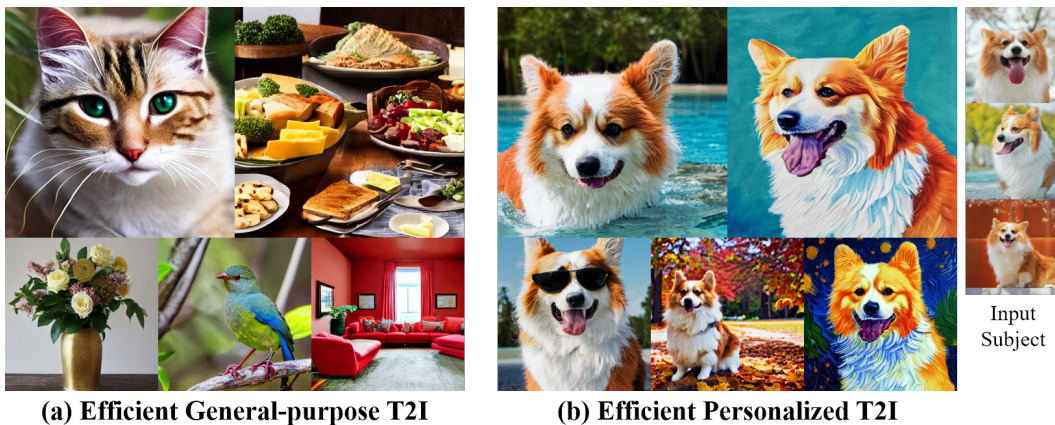


Figure 1: Our compressed stable diffusion enables efficient (a) zero-shot general-purpose text-to-image generation and (b) personalized synthesis. Selected samples from our lightest BK-SDM-Small with 36% reduced parameters and latency are shown.

# 1 Introduction

Large diffusion models [44, 51, 38, 47] have showcased groundbreaking results in text-to-image (T2I) synthesis tasks, which aim to create photorealistic images from textual descriptions. Stable Diffusion models (SDMs) [46, 47] are one of the most renowned open-source models, and their exceptional capability has begun to be leveraged as a backbone in several text-guided vision applications, e.g., text-driven image editing [2, 23] and 3D object creation [67], text-to-video generation [1, 68], and subject-driven [50, 25] and controllable [37, 71] T2I.

SDMs are T2I-specialized latent diffusion models (LDMs) [47], which employ diffusion operations [17, 59, 30] in a latent space to improve compute efficiency. Within a SDM, a U-Net [49, 6] conducts an iterative sampling procedure to gradually eliminate noise from random latents and is assisted by a text encoder [42] and an image decoder [9, 64] to produce text-aligned images. This inference process still involves excessive computational requirements (see Figure 2), which often hinder the utilization of SDMs despite their rapidly growing usage.

To alleviate this issue, numerous approaches toward efficient SDMs have been introduced.

Meng et al. [35, 34] reduce the number of denoising steps by distilling a pretrained diffusion model to guide an identically architected model with fewer sampling steps. Li et al. [28], Hou and Asghar [19], Shen et al. [57] employ post-training quantization techniques, and Chen et al. [4] enhance the implementation of SDMs for better compatibility with GPUs. However, the removal of architectural elements in diffusion models has not been investigated in spite of the established efficacy of structured pruning across discriminative models [26, 69] and generative adversarial networks (GANs) [31, 24].

This study unlocks the immense potential of classical architectural compression in attaining smaller and faster diffusion models. We eliminate multiple residual and attention blocks from the U-Net of a SDM and pretrain it with feature-level knowledge distillation (KD) [48, 13] for general-purpose T2I synthesis. Despite being trained with only 0.22M LAION pairs (less than 0.1% of the entire training pairs) [55] on a single A100 GPU, our compact models can mimic the original SDM by leveraging transferred knowledge. On the popular zero-shot MS-COCO benchmark [29], our work achieves a FID [15] score of 15.76 with 0.76B parameters and 16.98 with 0.66B parameters, which are on par with multi-billion parameter models [43, 7, 8]. Furthermore, we present the practical application of our lightweight pretrained models in customized T2I with DreamBooth finetuning [50].

Our contributions are summarized as follows:

- To the best of our knowledge, this is the first study to architecturally compress large-scale diffusion models. Our work is orthogonal to prior directions for efficient diffusion, e.g., enabling less sampling steps and employing quantization, and can be readily integrated with them.
- We compress SDMs by removing architectural blocks from the U-Net and achieve more than 30% reduction in model size and inference speed. We also introduce an interesting finding on the minor role of innermost blocks.
- We demonstrate the advantage of distillation-based pretraining, which allows us to attain competitive zero-shot T2I results even with very limited training resources.
- We highlight the capability of our light pretrained backbones in customized generation. Our models can lower the finetuning cost by 30% while retaining 97% scores of the original SDM.

## 2 Related work

**Large T2I diffusion models.** By gradually removing noise from corrupted data, diffusion-based generative models [18, 59, 6] enable high-fidelity synthesis with broad mode coverage. Integrating these merits with the advancement of pretrained language models [42, 41, 5] has significantly improved the quality of T2I synthesis. In GLIDE [38] and Imagen [51], a text-conditional diffusion

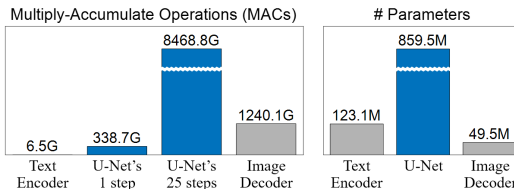


Figure 2: Computation of the major components in Stable Diffusion v1. The denoising U-Net is the main processing bottleneck. THOP [75] is used to measure MACs in generating a 512×512 image.

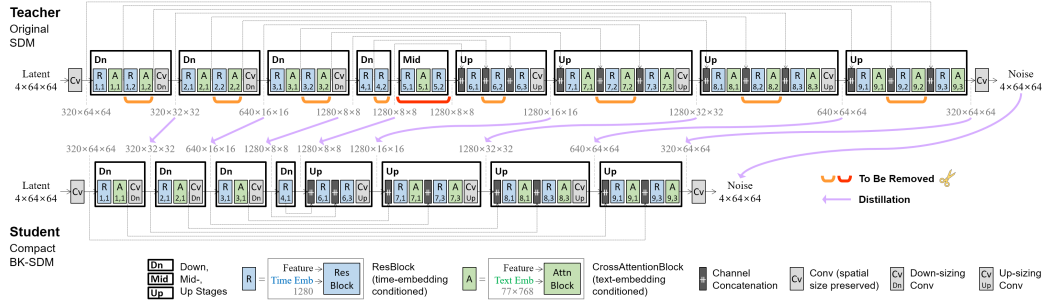


Figure 3: U-Net architectures of SDMs and KD-based pretraining process. The compact U-Net student is built by eliminating several residual and attention blocks from the original U-Net teacher. Through the feature and output distillation from the teacher, the student can be trained effectively yet rapidly. See Appendix for the details of block components.

68 model generates a  $64 \times 64$  image, which is upsampled via super-resolution modules. In DALL-E-2 [44],  
 69 a text-conditional prior network produces an image embedding, which is transformed into a  $64 \times 64$   
 70 image via a diffusion decoder and further upsampled into higher resolutions. SDMs [46, 47] perform  
 71 the diffusion modeling in a  $64 \times 64$  latent space constructed through a pixel-space autoencoder. We use  
 72 SDM as our baseline because of its open-access and gaining popularity over numerous downstream  
 73 tasks [2, 67, 1, 50].

74 **Efficient diffusion models.** Several studies have addressed the slow sampling process of diffusion  
 75 models. Diffusion-tailored distillation approaches [35, 34, 52] progressively transfer knowledge from  
 76 a pretrained diffusion model to a fewer-step model with the same architecture. Fast high-order solvers  
 77 [32, 33, 73] for diffusion ordinary differential equations boost the sampling speed. Orthogonal  
 78 to these directions for less sampling steps, our network compression approach reduces per-step  
 79 computation and can be easily integrated with them. Leveraging quantization techniques [28, 19, 57]  
 80 and implementation optimizations [4] has been applied for SDMs and also can be combined with our  
 81 models for further efficiency gains.

82 **Distillation-based compression.** KD enhances the performance of small-size models by exploiting  
 83 output-level [16, 39] and feature-level [48, 13, 70] information of large source models. Although  
 84 this classical distillation has been actively used toward efficient GANs [27, 45, 31, 22, 72], its power  
 85 has not been explored for structurally compressed diffusion models. Distillation-based pretraining  
 86 enables small yet capable general-purpose language models [54, 61, 21] and vision transformers  
 87 [63, 11]. Beyond such models, we show that its success can be extended to diffusion models with  
 88 iterative sampling steps. Concurrently with our study, a recently released small SDM without paper  
 89 evidence [40] similarly utilizes KD pretraining for a block-eliminated architecture, but it relies on  
 90 significantly more training resources along with multi-stage distillation. In contrast, our lightest  
 91 model achieves further reduced computation, and we show that competitive results can be obtained  
 92 even with much less data and single-stage distillation.

### 93 3 BK-SDM: block-removed knowledge-distilled SDM

94 We compress the U-Net [49] of a SDM [46, 47], which is the most compute-heavy component (see  
 95 Figure 2). Conditioned on the text and time-step embeddings, the U-Net performs multiple denoising  
 96 steps on latent representations. At each denoising step, the U-Net produces the noise residual to  
 97 compute the latent for the next step (see the top part of Figure 3). We reduce this per-step computation  
 98 by exploiting block-level elimination and feature distillation.

#### 99 3.1 Compressed U-Net architecture

100 The proposed models are referred to as:

- 101 ○ BK-SDM-Base (0.76B parameters) obtained with Section 3.1.1 (fewer blocks in outer stages).
- 102 ○ BK-SDM-Small (0.66B) with Section 3.1.1 (fewer blocks) and Section 3.1.2 (mid-stage removal).

Table 1: Minor impact of eliminating the mid-stage from the U-Net of SDM on zero-shot MS-COCO performance. Any retraining is not performed for the mid-stage removed model. For evaluation details, see Section 5.1.1.

Model	Performance		# Params	
	FID ↓	IS ↑	U-Net	Whole
SDM-v1.4 [46]	13.05	36.76	859.5M	1032.1M
Mid-Stage Removal	15.60	32.33	762.5M (-11.3%)	935.1M (-9.4%)

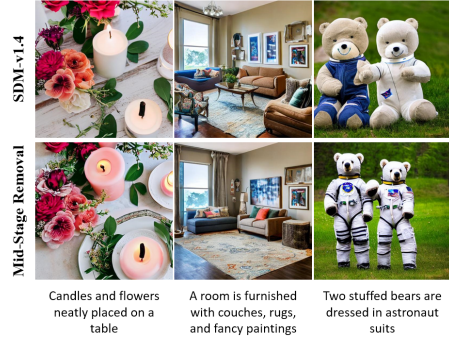


Figure 4: Visual results of the mid-stage removed U-Net without retraining.

### 103 3.1.1 Fewer blocks in the down and up stages

104 Our design philosophy is closely aligned with that of DistilBERT [54] which halves the number of  
 105 layers for improved computational efficiency and initializes the compact model with the original  
 106 weights by benefiting from the shared dimensionality. In the original U-Net, each stage with a  
 107 common spatial size consists of multiple blocks, and most stages contain pairs of residual (R) [12]  
 108 and cross-attention (A) [65, 20] blocks. We hypothesize the existence of some unnecessary pairs and  
 109 use the following removal strategies, as shown in Figure 3.

110 For the down stages, we maintain the first R-A pairs while eliminating the second pairs, because the  
 111 first pairs process the changed spatial information and would be more important than the second pairs.  
 112 This design choice does not harm the dimensionality of the original U-Net, enabling the use of the  
 113 corresponding pretrained weights for initialization [54].

114 For the up stages, while adhering to the aforementioned scheme, we retain the third R-A pairs. This  
 115 allows us to utilize the output feature maps at the end of each down stage and the corresponding skip  
 116 connections between the down and up stages. The same process is applied to the innermost down and  
 117 up stages that contain only R blocks.

### 118 3.1.2 Removal of the entire mid-stage

119 Surprisingly, removing the entire mid-stage from the original U-Net (marked with red in Figure 3)  
 120 does not noticeably degrade the generation quality for many text prompts while effectively reducing  
 121 the number of parameters (see Table 1 and Figure 4). This observation is consistent with the minor  
 122 role of inner layers in the U-Net generator of GANs [24].

123 Integrating the mid-stage removal with fewer blocks in Section 3.1.1 further decreases computational  
 124 burdens (Table 3) at the cost of a slight decline in performance (Table 2). Therefore, we offer  
 125 this mid-stage elimination as an option, depending on the priority between compute efficiency and  
 126 generation quality.

### 127 3.2 Distillation-based pretraining

128 For general-purpose T2I generation, we train the compact U-Net to mimic the behavior of the original  
 129 U-Net. Following Rombach et al. [47], we use the pretrained-and-frozen encoders to obtain the inputs  
 130 of the U-Net.

131 Given the latent representation  $z$  of an image and its paired text embedding  $y$ , the task loss for the  
 132 reverse denoising process [18, 47] is computed as:

$$\mathcal{L}_{\text{Task}} = \mathbb{E}_{z, \epsilon, y, t} \left[ \|\epsilon - \epsilon_S(z_t, y, t)\|_2^2 \right], \quad (1)$$

133 where  $\epsilon \sim N(0, I)$  and  $t \sim \text{Uniform}(1, T)$  denote the noise and time step sampled from the diffusion  
 134 process, respectively, and  $\epsilon_S(\circ)$  indicates the output of our compact U-Net student. For brevity, we  
 135 omit the subscripts of  $\mathbb{E}_{z, \epsilon, y, t}[\circ]$  in the following notations.



136 The compact student is also trained to imitate the outputs of the original U-Net teacher,  $\epsilon_T(\circ)$ , with  
 137 the following output-level KD objective [16]:

$$\mathcal{L}_{\text{OutKD}} = \mathbb{E} \left[ \|\epsilon_T(z_t, y, t) - \epsilon_S(z_t, y, t)\|_2^2 \right]. \quad (2)$$

138 A key to our approach is the utilization of feature-level KD [48, 13] that provides abundant guidance  
 139 for the student’s training:

$$\mathcal{L}_{\text{FeatKD}} = \mathbb{E} \left[ \sum_l \|f_T^l(z_t, y, t) - f_S^l(z_t, y, t)\|_2^2 \right], \quad (3)$$

140 where  $f_T^l(\circ)$  and  $f_S^l(\circ)$  represent the feature maps of the  $l$ -th layer in a predefined set of distilled layers  
 141 from the teacher and the student, respectively. While learnable regressors (e.g.,  $1 \times 1$  convolutions  
 142 to match the number of channels) have been commonly used in existing studies [58, 45, 48], our  
 143 approach circumvents this requirement. By applying distillation at the end of each stage in both  
 144 models, we ensure that the dimensionality of the feature maps already matches, thus eliminating the  
 145 need for additional regressors.

146 The final objective is formalized as below, and we simply set the loss weights  $\lambda_{\text{OutKD}}$  and  $\lambda_{\text{FeatKD}}$   
 147 as 1. Without any hyperparameter tuning, our approach is effective in empirical validation.

$$\mathcal{L} = \mathcal{L}_{\text{Task}} + \lambda_{\text{OutKD}} \mathcal{L}_{\text{OutKD}} + \lambda_{\text{FeatKD}} \mathcal{L}_{\text{FeatKD}}. \quad (4)$$

### 148 3.3 Application: faster and smaller personalized SDMs

149 To emphasize the benefit of our lightweight pretrained SDMs, we use a popular finetuning scenario  
 150 for personalized generation. DreamBooth [50] enables T2I diffusion models to create contents about  
 151 a particular subject using just a few input images. Our compact models not only accelerate inference  
 152 speed but also reduce finetuning cost. Moreover, they produce high-quality images based on the  
 153 inherited capability of the original SDM.

## 154 4 Experimental setup

### 155 4.1 Datasets and evaluation metrics

156 **Pretraining.** We train our compact SDM with only 0.22M image-text pairs from LAION-Aesthetics  
 157 V2 6.5+ [55, 56], which are significantly fewer than the original training data used for SDM-v1.4  
 158 [46] (i.e., 600M pairs of LAION-Aesthetics V2 5+ [55] for the resumed training).

159 **Zero-shot T2I evaluation.** Following the popular protocol [43, 47, 51] to assess general-purpose T2I  
 160 with pretrained models, we use 30K prompts from the MS-COCO validation split [29] and compare  
 161 the generated images to the whole validation set. We compute Fréchet Inception Distance (FID) [15]  
 162 and Inception Score (IS) [53] to assess visual quality. Moreover, we measure CLIP score [42, 14]  
 163 with CLIP-ViT-g/14 model to assess text-image correspondence.

164 **Finetuning for personalized generation.** We use the DreamBooth dataset [50] that covers 30  
 165 subjects, each of which is associated with 25 prompts and 4~6 images. Through individual finetuning  
 166 for each subject, 30 personalized models are obtained. For evaluation, we follow the protocol of Ruiz  
 167 et al. [50] based on four synthesized images per subject and per prompt. We consider CLIP-I and  
 168 DINO scores to measure how well subject details are maintained in generated images (i.e., subject  
 169 fidelity) and CLIP-T scores to measure text-image alignment (i.e., text fidelity). We use ViT-S/16  
 170 embeddings [3] for DINO scores and CLIP-ViT-g/14 embeddings for CLIP-I and CLIP-T.

### 171 4.2 Implementation

172 We use the released version v1.4 of SDM [46] as our compression target. We remark that our approach  
 173 is also applicable to other versions in v1.1–v1.5 with the same architecture and to SDM-v2 with a  
 174 similarly designed architecture.

Table 2: Zero-shot results on 30K prompts from MS-COCO validation set [29] at 256×256 resolution. Despite being trained with a smaller dataset and having fewer parameters, our compressed models achieve results on par with prior approaches for general-purpose T2I. For our models, the results with the minimum FID and the final 50K-th iteration are reported (see Section 5.1.3 for detailed analysis).

Model	Type	FID ↓	IS ↑	# Params	Data Size
SDM-v1.4 [47]	DF	13.05	36.76	1.04B	600M
Small Stable Diffusion [40]	DF	12.76	32.33	0.76B	229M
BK-SDM-Base (Ours) @ Min FID	DF	13.57	29.22	0.76B	0.22M
BK-SDM-Base (Ours) @ Final Iter	DF	15.76	33.79	0.76B	0.22M
BK-SDM-Small (Ours) @ Min FID	DF	15.93	29.61	0.66B	0.22M
BK-SDM-Small (Ours) @ Final Iter	DF	16.98	31.68	0.66B	0.22M
DALL-E <sup>†*</sup> [43]	AR	27.5	17.9	12B	250M
CogView <sup>‡*</sup> [7]	AR	27.1	18.2	4B	30M
CogView2 <sup>†*</sup> [8]	AR	24.0	22.4	6B	30M
Make-A-Scene <sup>‡</sup> [10]	AR	11.84	-	4B	35M
LAFITE <sup>‡‡</sup> [74]	GAN	26.94	26.02	0.23B	3M
GALIP (CC3M) <sup>†</sup> [62]	GAN	16.12	-	0.32B	3M
GALIP (CC12M) <sup>†</sup> [62]	GAN	12.54	-	0.32B	12M
GLIDE <sup>‡</sup> [38]	DF	12.24	-	5B	250M
LDM-KL-8-G <sup>‡‡</sup> [47]	DF	12.63	30.29	1.45B	400M
DALL-E-2 <sup>†</sup> [44]	DF	10.39	-	5.2B	250M

<sup>†</sup> and <sup>‡</sup>: FID from [62] and [47], respectively. \* and <sup>‡</sup>: IS from [8] and [47], respectively. DF and AR: diffusion and autoregressive models. ↓ and ↑: lower and higher values are better.

175 We adjust the codes in Diffusers library [66] for pretraining our models and those in PEFT library  
 176 [60] for DreamBooth-finetuning, both of which adopt the training process of DDPM [18] in latent  
 177 spaces. We use a single NVIDIA A100 80G GPU for 50K-iteration pretraining with a constant  
 178 learning rate of 5e-5. For DreamBooth, we use a single NVIDIA GeForce RTX 3090 GPU to finetune  
 179 each personalized model for 800 iterations with a constant learning rate of 1e-6.

180 Following the default inference setup, we use PNDM scheduler [30] for zero-shot T2I generation  
 181 and DPM-Solver [32, 33] for DreamBooth results. For compute efficiency, we always opt for 25  
 182 denoising steps of the U-Net at the inference phase. The classifier-free guidance scale [17, 51] is set  
 183 to the default value of 7.5, except the analysis in Figure 7.

## 184 5 Results

### 185 5.1 General-purpose T2I generation

#### 186 5.1.1 Main results

187 Table 2 shows the zero-shot T2I results on 30K samples from the MS-COCO 256×256 validation  
 188 set. Despite being trained with only 0.22M samples and having fewer than 1B parameters, our  
 189 compressed models demonstrate competitive performance on par with previous large pretrained  
 190 models. Despite the absence of a paper support, we include the model [40] that is identical in  
 191 structure to BK-SDM-Base for comparison. This model benefits from far more training resources,  
 192 i.e., two-stage KD relied on two teachers (SDM-v1.4 and v1.5) and a much larger volume of data  
 193 with significantly longer iterations.

194 Figure 5 depicts synthesized images of different models with some MS-COCO captions. Our  
 195 compressed models inherit the superior ability of SDM and produce more photorealistic images  
 196 compared to the AR-based [8] and GAN-based [74, 62] baselines. Noticeably, the same latent code  
 197 results in a shared visual style between the original and our compact SDMs (4th–6th columns in  
 198 Figure 5), similar to the observation in transfer learning for GANs [36].

199 Table 3 summarizes how the computational reduction for each sampling step of the U-Net impacts the  
 200 overall compute of the entire SDM. The per-step reduction effectively decreases MACs and inference  
 201 time by more than 30% as well as the number of parameters.



Figure 5: Visual comparison on zero-shot MS-COCO benchmark. The results of previous studies [8, 74, 62] were obtained with their official codes and released models. We do not apply any CLIP-based reranking for SDM and our models.

Table 3: The impact of per-step compute reduction of the U-Net on the entire SDM. The number of sampling steps is indicated with the parentheses, e.g., U-Net (1) for one step. The full computation (denoted by “Whole”) covers the text encoder, U-Net, and image decoder. All corresponding values are obtained on the generation of a single 512×512 image with 25 denoising steps. The latency was measured on Xeon Silver 4210R CPU 2.40GHz and NVIDIA GeForce RTX 3090 GPU.

Model	# Params		MACs			CPU Latency			GPU Latency		
	U-Net	Whole	U-Net (1)	U-Net (25)	Whole	U-Net (1)	U-Net (25)	Whole	U-Net (1)	U-Net (25)	Whole
SDM-v1.4 [46]	860M	1033M	339G	8469G	9716G	5.63s	146.28s	153.02s	0.049s	1.28s	1.41s
BK-SDM-Base (Ours)	580M	752M	224G	5594G	6841G	3.84s	99.95s	106.62s	0.032s	0.83s	0.96s
	(-32.6%)	(-27.1%)	(-33.9%)	(-33.9%)	(-29.5%)	(-31.8%)	(-31.7%)	(-30.3%)	(-34.6%)	(-35.2%)	(-31.9%)
BK-SDM-Small (Ours)	483M	655M	218G	5444G	6690G	3.45s	89.78s	96.52s	0.030s	0.77s	0.90s
	(-43.9%)	(-36.5%)	(-35.7%)	(-35.7%)	(-31.1%)	(-38.7%)	(-38.6%)	(-36.9%)	(-38.7%)	(-39.8%)	(-36.1%)

## 202 5.1.2 Ablation study

203 Table 4 presents the ablation study with the zero-shot MS-COCO benchmark dataset. The common  
 204 default settings for the models N1–N7 involve the usage of fewer blocks in the down and up stages  
 205 (Section 3.1.1) and the denoising task loss (Eq. 1). All the models are drawn at the 50K-th training  
 206 iteration. We made the following observations.

207 **N1 vs. N2.** Importing the pretrained weights for initialization clearly improves the performance of  
 208 block-removed SDMs. Transferring knowledge from well-trained models, a popularized practice in  
 209 machine learning, is also beneficial for T2I generation with SDMs.

210 **N2 vs. N3 vs. N4.** Exploiting output-level KD (Eq. 2) effectively boosts the generation quality  
 211 compared to using only the denoising task loss. Leveraging feature-level KD (Eq. 3) further improves  
 212 the performance by offering sufficient guidance over multiple stages in the student.

213 **N4 vs. N5.** An increased batch size leads to a better IS and CLIP score but with a minor drop in FID.  
 214 We opt for a batch size of 256 based on the premise that more samples per batch would enhance the  
 215 model’s understanding ability.

216 **N6 and N7.** Despite slight performance drop, the models N6 and N7 with the mid-stage removal  
 217 have fewer parameters (0.66B) than N4 and N5 (0.76B), offering improved compute efficiency.

Table 4: Ablation study on zero-shot MS-COCO 256×256 30K. The common settings include fewer blocks in the down and up stages and the denoising task loss. N5 and N7 correspond to BK-SDM-Base and BK-SDM-Small, respectively

Model						Performance		
No.	Initialize Weights	Output KD	Feature KD	Batch Size	Remove Mid	FID ↓	IS ↑	CLIP score ↑
N1	Random	✗	✗	64	✗	43.80	13.61	0.1622
N2	Pretrained	✗	✗	64	✗	20.45	22.68	0.2444
N3	Pretrained	✓	✗	64	✗	16.48	27.30	0.2620
N4	Pretrained	✓	✓	64	✗	14.61	31.44	0.2826
N5	Pretrained	✓	✓	256	✗	15.76	33.79	0.2878
N6	Pretrained	✓	✓	64	✓	16.87	29.51	0.2644
N7	Pretrained	✓	✓	256	✓	16.98	31.68	0.2677
Original SDM-v1.4 [46, 47]						13.05	36.76	0.2958

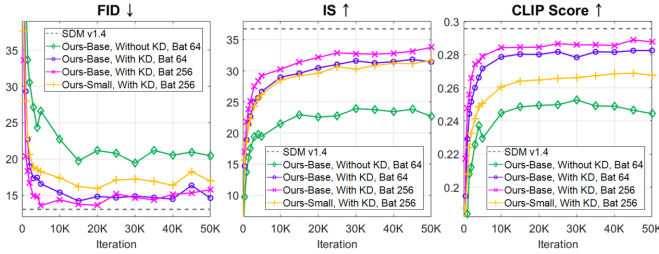


Figure 6: Results on zero-shot MS-COCO 256×256 30K over training progress. For our models, the architecture size, usage of KD, and batch size are denoted.

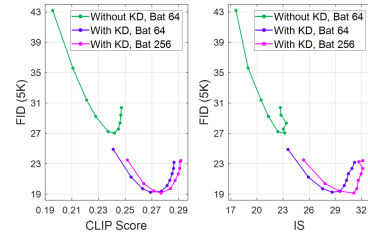


Figure 7: Effect of different classifier-free guidance scales on MS-COCO 512×512 5K.

218 **5.1.3 Impact of distillation on pretraining phase**

219 We further analyze the merits of transferred knowledge via distillation, with the models from the  
 220 pretrained weight initialization. Figure 6 shows zero-shot T2I performance over training iterations.  
 221 Compared to the absence of KD (indicated with green), distillation (purple and pink) accelerates the  
 222 training process and leads to improved generation scores, demonstrating the benefits of providing  
 223 sufficient hints for training guidance. Notably, our small-size model trained with KD (yellow)  
 224 outperforms the bigger base-size model without KD (green). Additionally, while the best FID score  
 225 is observed early on for our models, IS and CLIP score exhibit ongoing improvement, implying that  
 226 judging models solely with FID may be suboptimal.

227 Figure 7 shows the trade-off curves from different classifier-free guidance scales [17, 51]  
 228 {2.0, 2.5, 3.0, 3.5, 4.5, 5.5, 6.5, 7.5, 8.5, 9.5}. For the analysis, we use 5K samples from the MS-  
 229 COCO validation set and our base-size models from the 50K-th iteration. Higher guidance scales  
 230 lead to better text-aligned images at the cost of less diversity. Compared to the baseline trained only  
 231 with the denoising task loss, distillation-based pretraining leads to much better trade-off curves.

232 **5.2 Personalized T2I with DreamBooth**

233 Table 5 compares the results of DreamBooth finetuning [50] with different pretrained models. BK-  
 234 SDM-Small can preserve over 97% performance of the original SDM with the reduced finetuning  
 235 time and number of parameters. Figure 8 depicts that our models can accurately capture the subject  
 236 details and generate various scenes. Over the models pretrained with a batch size of 64, we observe  
 237 the impact of KD pretraining on personalized synthesis. The baselines without KD fail to generate  
 238 the subjects entirely or cannot maintain the identity details.



Table 5: Personalized generation with finetuning over different pretrained models. Our compact models can preserve subject fidelity (DINO and CLIP-I) and prompt fidelity (CLIP-T) of the original SDM with reduced finetuning (FT) time and fewer parameters.

Pretrained Model	DINO $\uparrow$	CLIP-I $\uparrow$	CLIP-T $\uparrow$	FT Time <sup>†</sup>	# Params
SDM-v1.4 [46, 47]	0.728	0.725	0.263	881.3s	1.04B
BK-SDM-Base (Ours)	0.723	0.717	0.260	622.3s	0.76B
BK-SDM-Small (Ours)	0.720	0.705	0.259	603.6s	0.66B
BK-SDM-Base, Batch Size 64	0.718	0.708	0.262	622.3s	0.76B
- Without KD & Random Init.	0.594	0.465	0.191	622.3s	0.76B
- Without KD & Pretrained Init.	0.716	0.669	0.258	622.3s	0.76B

<sup>†</sup> Per-subject finetuning time for 800 iterations on NVIDIA GeForce RTX 3090 GPU.

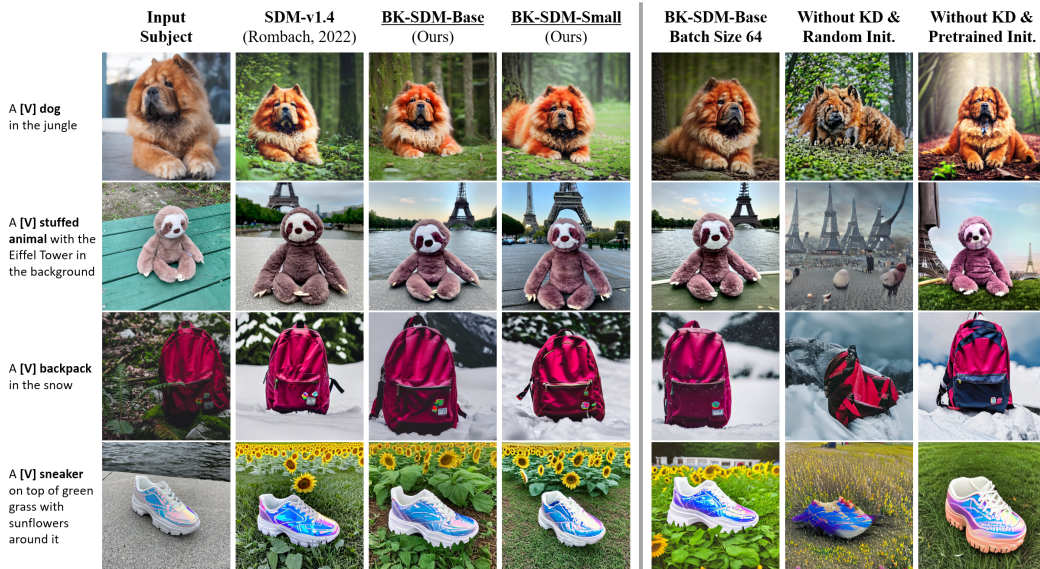


Figure 8: Visual results of personalized generation. Each subject is marked as “a [identifier] [class noun]” (e.g., “a [V] dog”). Similar to the original SDM, our compact models can synthesize the images of input subjects in different backgrounds while preserving their appearance.

## 239 6 Conclusion and discussion

240 This study uncovers the potential of architectural compression for general-purpose text-to-image  
 241 synthesis with a renowned model, Stable Diffusion. Our block-removed lightweight models are  
 242 effective for zero-shot generation, achieving competitive performance against large-scale baselines.  
 243 Distillation is a key aspect of our method, leading to effective pretraining even under very constrained  
 244 resources. Moreover, our smaller and faster pretrained models are successfully applied in personalized  
 245 generation. Our work is orthogonal to previous directions for efficient diffusion models, e.g., enabling  
 246 fewer sampling steps, and can be readily combined with them. We hope our study can facilitate future  
 247 research on structural compression of large diffusion models.

248 **Limitations and future works.** Our compact models inherit the capability of the source model for  
 249 high-fidelity image generation, but they have shortcomings such as inaccurate generation of full-body  
 250 human appearance. While we show that distillation pretraining is powerful even with very limited  
 251 resources, increasing the volume of data and analyzing its effects would be promising.

252 **Negative social impacts.** Because recent large generative models are capable of creating high-quality  
 253 plausible content, they also involve potential risks of malicious use. To avoid causing unintended  
 254 social bias, researchers should take steps to ensure the appropriateness of training data. Moreover,  
 255 the release of resulting models should be accompanied by strong and reliable safeguards.

## References

- 256
- 257 [1] A. Blattmann, R. Rombach, H. Ling, T. Dockhorn, S. W. Kim, S. Fidler, and K. Kreis. Align your latents:  
258 High-resolution video synthesis with latent diffusion models. In *CVPR*, 2023.
- 259 [2] T. Brooks, A. Holynski, and A. A. Efros. Instructpix2pix: Learning to follow image editing instructions.  
260 In *CVPR*, 2023.
- 261 [3] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging properties in  
262 self-supervised vision transformers. In *ICCV*, 2021.
- 263 [4] Y.-H. Chen, R. Sarokin, J. Lee, J. Tang, C.-L. Chang, A. Kulik, and M. Grundmann. Speed is all you need:  
264 On-device acceleration of large diffusion models via gpu-aware optimizations. In *CVPR Workshop*, 2023.
- 265 [5] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers  
266 for language understanding. In *NAACL*, 2019.
- 267 [6] P. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. In *NeurIPS*, 2021.
- 268 [7] M. Ding, Z. Yang, W. Hong, W. Zheng, C. Zhou, D. Yin, J. Lin, X. Zou, Z. Shao, H. Yang, et al. Cogview:  
269 Mastering text-to-image generation via transformers. In *NeurIPS*, 2021.
- 270 [8] M. Ding, W. Zheng, W. Hong, and J. Tang. Cogview2: Faster and better text-to-image generation via  
271 hierarchical transformers. In *NeurIPS*, 2022.
- 272 [9] P. Esser, R. Rombach, and B. Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*,  
273 2021.
- 274 [10] O. Gafni, A. Polyak, O. Ashual, S. Sheynin, D. Parikh, and Y. Taigman. Make-a-scene: Scene-based  
275 text-to-image generation with human priors. In *ECCV*, 2022.
- 276 [11] Z. Hao, J. Guo, D. Jia, K. Han, Y. Tang, C. Zhang, H. Hu, and Y. Wang. Learning efficient vision  
277 transformers via fine-grained manifold distillation. In *NeurIPS*, 2022.
- 278 [12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- 279 [13] B. Heo, J. Kim, S. Yun, H. Park, N. Kwak, and J. Y. Choi. A comprehensive overhaul of feature distillation.  
280 In *ICCV*, 2019.
- 281 [14] J. Hessel, A. Holtzman, M. Forbes, R. Le Bras, and Y. Choi. CLIPScore: A reference-free evaluation  
282 metric for image captioning. In *EMNLP*, 2021.
- 283 [15] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale  
284 update rule converge to a local nash equilibrium. In *NeurIPS*, 2017.
- 285 [16] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. In *NeurIPS Workshop*,  
286 2014.
- 287 [17] J. Ho and T. Salimans. Classifier-free diffusion guidance. In *NeurIPS Workshop*, 2021.
- 288 [18] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020.
- 289 [19] J. Hou and Z. Asghar. World’s first on-device demonstration of stable diffusion on an android phone.  
290 <https://www.qualcomm.com/news>, 2023.
- 291 [20] A. Jaegle, F. Gimeno, A. Brock, O. Vinyals, A. Zisserman, and J. Carreira. Perceiver: General perception  
292 with iterative attention. In *ICML*, 2021.
- 293 [21] X. Jiao, Y. Yin, L. Shang, X. Jiang, X. Chen, L. Li, F. Wang, and Q. Liu. Tinybert: Distilling bert for  
294 natural language understanding. In *Findings of EMNLP*, 2020.
- 295 [22] Q. Jin, J. Ren, O. J. Woodford, J. Wang, G. Yuan, Y. Wang, and S. Tulyakov. Teachers do more than teach:  
296 Compressing image-to-image models. In *CVPR*, 2021.
- 297 [23] B. Kavar, S. Zada, O. Lang, O. Tov, H. Chang, T. Dekel, I. Mosseri, and M. Irani. Imagic: Text-based real  
298 image editing with diffusion models. In *CVPR*, 2023.
- 299 [24] B.-K. Kim, S. Choi, and H. Park. Cut inner layers: A structured pruning strategy for efficient u-net gans.  
300 In *ICML Workshop*, 2022.

- 301 [25] N. Kumari, B. Zhang, R. Zhang, E. Shechtman, and J.-Y. Zhu. Multi-concept customization of text-to-image  
302 diffusion. In *CVPR*, 2023.
- 303 [26] H. Li, A. Kadav, I. Durdanovic, H. Samet, and H. P. Graf. Pruning filters for efficient convnets. In *ICLR*,  
304 2017.
- 305 [27] M. Li, J. Lin, Y. Ding, Z. Liu, J.-Y. Zhu, and S. Han. Gan compression: Efficient architectures for  
306 interactive conditional gans. In *CVPR*, 2020.
- 307 [28] X. Li, L. Lian, Y. Liu, H. Yang, Z. Dong, D. Kang, S. Zhang, and K. Keutzer. Q-diffusion: Quantizing  
308 diffusion models. *arXiv preprint arXiv:2302.04304*, 2023.
- 309 [29] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft  
310 coco: Common objects in context. In *ECCV*, 2014.
- 311 [30] L. Liu, Y. Ren, Z. Lin, and Z. Zhao. Pseudo numerical methods for diffusion models on manifolds. In  
312 *ICLR*, 2022.
- 313 [31] Y. Liu, Z. Shu, Y. Li, Z. Lin, F. Perazzi, and S.-Y. Kung. Content-aware gan compression. In *CVPR*, 2021.
- 314 [32] C. Lu, Y. Zhou, F. Bao, J. Chen, C. Li, and J. Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic  
315 model sampling in around 10 steps. In *NeurIPS*, 2022.
- 316 [33] C. Lu, Y. Zhou, F. Bao, J. Chen, C. Li, and J. Zhu. Dpm-solver++: Fast solver for guided sampling of  
317 diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022.
- 318 [34] C. Meng, R. Gao, D. P. Kingma, S. Ermon, J. Ho, and T. Salimans. On distillation of guided diffusion  
319 models. In *NeurIPS Workshop*, 2022.
- 320 [35] C. Meng, R. Gao, D. P. Kingma, S. Ermon, J. Ho, and T. Salimans. On distillation of guided diffusion  
321 models. In *CVPR*, 2023.
- 322 [36] S. Mo, M. Cho, and J. Shin. Freeze the discriminator: a simple baseline for fine-tuning gans. In *CVPR  
323 Workshop*, 2020.
- 324 [37] C. Mou, X. Wang, L. Xie, J. Zhang, Z. Qi, Y. Shan, and X. Qie. T2i-adapter: Learning adapters to dig out  
325 more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023.
- 326 [38] A. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen. Glide:  
327 Towards photorealistic image generation and editing with text-guided diffusion models. In *ICML*, 2022.
- 328 [39] W. Park, D. Kim, Y. Lu, and M. Cho. Relational knowledge distillation. In *CVPR*, 2019.
- 329 [40] J. Pinkney. Small stable diffusion. [https://huggingface.co/OFA-Sys/  
330 small-stable-diffusion-v0](https://huggingface.co/OFA-Sys/small-stable-diffusion-v0), 2023.
- 331 [41] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al. Language models are unsupervised  
332 multitask learners. *OpenAI blog*, 2019.
- 333 [42] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin,  
334 J. Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- 335 [43] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever. Zero-shot  
336 text-to-image generation. In *ICML*, 2020.
- 337 [44] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen. Hierarchical text-conditional image generation  
338 with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 339 [45] Y. Ren, J. Wu, X. Xiao, and J. Yang. Online multi-granularity distillation for gan compression. In *ICCV*,  
340 2021.
- 341 [46] R. Rombach and P. Esser. Stable diffusion v1-4. [https://huggingface.co/CompVis/  
342 stable-diffusion-v1-4](https://huggingface.co/CompVis/stable-diffusion-v1-4), 2022.
- 343 [47] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with  
344 latent diffusion models. In *CVPR*, 2022.
- 345 [48] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio. Fitnets: Hints for thin deep  
346 nets. In *ICLR*, 2015.

- 347 [49] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation.  
348 In *MICCAI*, 2015.
- 349 [50] N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman. Dreambooth: Fine tuning  
350 text-to-image diffusion models for subject-driven generation. In *CVPR*, 2023.
- 351 [51] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes,  
352 B. Karagol Ayan, T. Salimans, et al. Photorealistic text-to-image diffusion models with deep language  
353 understanding. In *NeurIPS*, 2022.
- 354 [52] T. Salimans and J. Ho. Progressive distillation for fast sampling of diffusion models. In *ICLR*, 2022.
- 355 [53] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for  
356 training gans. In *NeurIPS*, 2016.
- 357 [54] V. Sanh, L. Debut, J. Chaumond, and T. Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper  
358 and lighter. In *NeurIPS Workshop*, 2019.
- 359 [55] C. Schuhmann and R. Beaumont. Laion-aesthetics. <https://laion.ai/blog/laion-aesthetics>,  
360 2022.
- 361 [56] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta,  
362 C. Mullis, M. Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation  
363 image-text models. In *NeurIPS Workshop*, 2022.
- 364 [57] H. Shen, P. Cheng, X. Ye, W. Cheng, and H. Abidi. Accelerate stable diffusion with intel neural compressor.  
365 <https://medium.com/intel-analytics-software>, 2022.
- 366 [58] C. Shu, Y. Liu, J. Gao, Z. Yan, and C. Shen. Channel-wise knowledge distillation for dense prediction. In  
367 *ICCV*, 2021.
- 368 [59] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. In *ICLR*, 2021.
- 369 [60] L. D. Y. B. S. P. Sourab Mangrulkar, Sylvain Gugger. Peft: State-of-the-art parameter-efficient fine-tuning  
370 methods. <https://github.com/huggingface/peft>, 2022.
- 371 [61] Z. Sun, H. Yu, X. Song, R. Liu, Y. Yang, and D. Zhou. Mobilebert: a compact task-agnostic bert for  
372 resource-limited devices. In *ACL*, 2020.
- 373 [62] M. Tao, B.-K. Bao, H. Tang, and C. Xu. Galip: Generative adversarial clips for text-to-image synthesis. In  
374 *CVPR*, 2023.
- 375 [63] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jegou. Training data-efficient image  
376 transformers amp; distillation through attention. In *ICML*, 2021.
- 377 [64] A. Van Den Oord, O. Vinyals, et al. Neural discrete representation learning. In *NeurIPS*, 2017.
- 378 [65] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin.  
379 Attention is all you need. In *NeurIPS*, 2017.
- 380 [66] P. von Platen, S. Patil, A. Lozhkov, P. Cuenca, N. Lambert, K. Rasul, M. Davaadorj, and T. Wolf. Diffusers:  
381 State-of-the-art diffusion models. <https://github.com/huggingface/diffusers>, 2022.
- 382 [67] H. Wang, X. Du, J. Li, R. A. Yeh, and G. Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d  
383 diffusion models for 3d generation. In *CVPR*, 2023.
- 384 [68] J. Z. Wu, Y. Ge, X. Wang, W. Lei, Y. Gu, W. Hsu, Y. Shan, X. Qie, and M. Z. Shou. Tune-a-video: One-shot  
385 tuning of image diffusion models for text-to-video generation. *arXiv preprint arXiv:2212.11565*, 2022.
- 386 [69] Z. Xie, L. Zhu, L. Zhao, B. Tao, L. Liu, and W. Tao. Localization-aware channel pruning for object  
387 detection. *Neurocomputing*, 403:400–408, 2020.
- 388 [70] S. Zagoruyko and N. Komodakis. Paying more attention to attention: Improving the performance of  
389 convolutional neural networks via attention transfer. In *ICLR*, 2017.
- 390 [71] L. Zhang and M. Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint*  
391 *arXiv:2302.05543*, 2023.
- 392 [72] L. Zhang, X. Chen, X. Tu, P. Wan, N. Xu, and K. Ma. Wavelet knowledge distillation: Towards efficient  
393 image-to-image translation. In *CVPR*, 2022.



- 394 [73] Q. Zhang and Y. Chen. Fast sampling of diffusion models with exponential integrator. In *ICLR*, 2023.
- 395 [74] Y. Zhou, R. Zhang, C. Chen, C. Li, C. Tensmeyer, T. Yu, J. Gu, J. Xu, and T. Sun. Towards language-free  
396 training for text-to-image generation. In *CVPR*, 2022.
- 397 [75] L. Zhu. Thop: Pytorch-opcounter. <https://github.com/Lyken17/pytorch-OpCounter>, 2018.